

1 **Title: The HERV-K HML-2 integration within *RASGRF2* is associated with**  
2 **intravenous drug abuse and modulates transcription in a cell-line model**

3

4

5

6 +Correspondence should be addressed to:

7

8 aris.katzourakis@zoo.ox.ac.uk

9 gmagi@med.uoa.gr

10

11

12 **This file includes:**

13 Materials and Methods

14 Supplementary Figures

15 **Other Supplementary Materials for this manuscript includes the following:**

16 Dataset S1

17

## 18 **Materials and Methods**

19

### 20 **Wet Lab**

21

#### 22 *Design of target enrichment and high throughput sequencing*

23 We developed a custom library-preparation method for the targeted  
24 enrichment of HK2 LTR and their flanking sequences. To do this, we designed  
25 custom 120bp-long biotinylated probes from IDT to perform a capture of the  
26 first and the last ~360bp of the HK2-LTR sequence (5 probes overlapping by  
27 60bp at each LTR side). The libraries were prepared with a target insert size  
28 at library Quality Control (QC) of 1200bp. From trial runs, we knew that a  
29 library QC showing an insert of approximately 1200bp would produce  
30 sequencing data with a considerably shorter insert. The prepared libraries  
31 were validated and multiplexed before 2X300bp paired-end sequencing on an  
32 Illumina MiSeq instrument.

33

#### 34 *Cells and reagents*

35 General cell culture reagents were purchased from Sigma-Aldrich  
36 (Poole, U.K.) The human embryonic kidney HEK 293 cell line was purchased  
37 for this specific project from the European Collection of Authenticated Cell  
38 Cultures (ECCAC; Porton Down, U.K.) and the teratocarcinoma cell line  
39 NCCIT was obtained from the American Type Culture Collection (ATCC, LGC  
40 Standards, Teddington, U.K.). Both of the cell lines came with quality  
41 certificates for their identity and being free from mycoplasma contamination.

42 The cells were incubated in a humidified atmosphere at 37°C with 5% CO<sub>2</sub>  
43 and cultured as described below.

44 HEK 293 cells were cultured in Dulbecco's Modified Eagle's Medium  
45 (DMEM) with high glucose (4500mg/dL), L-glutamine and sodium  
46 bicarbonate. It was supplemented with 10% heat-inactivated foetal bovine  
47 serum, 1% penicillin-streptomycin and 2mM L-glutamine. The cells were  
48 allowed to grow to 100% confluence in a T75 flask and were typically  
49 subcultured twice per week using trypsin-EDTA to dislodge them from the  
50 substratum. The cell suspension was diluted 1:20 into a new flask. Two days  
51 after seeding a flask, the media was replaced and the spent media filtered  
52 through a 0.2µM syringe filter. Aliquots of the filtered media were stored at -  
53 80°C until needed for use as conditioned media in genome engineering.

54 The NCCIT cells were cultured in RPMI growth medium supplemented  
55 with 10% v/v foetal bovine serum, 1% v/v penicillin-streptomycin and 2mM L-  
56 glutamine. Two days after seeding the flask, the media was replaced with  
57 fresh growth media due to the elevated respiration rate. They were allowed to  
58 grow to 100% confluence in a T75 flask and were typically subcultured every  
59 four days using trypsin-EDTA to dislodge them from the substratum. On each  
60 passage, approximately  $1 \times 10^6$  cells were reserved for later gDNA extraction.

61 SCR7 was obtained from Stratech Scientific Limited (Newmarket,  
62 U.K.). Puromycin was purchased from Cayman Chemicals (Cambridge  
63 Bioscience, Cambridge, U.K.)

64

65 *Genome Engineering*

66

67 *-CRISPR/Cas9 plasmids*

68 The Cas9 and sgRNA plasmids were designed by Andrew Bassett and Joey  
69 Riepsaame of Genome Engineering Oxford (GEO) to introduce the HERV  
70 sequence into the RASGRF2 gene. We designed the HDR plasmid and this  
71 was synthesised using GeneArt (Thermo Fisher Scientific). Subsequently, for  
72 the HAP1 transfection experiments we sub-cloned the HDR construct into a  
73 new backbone to remove the CMV promoter and increase targeting  
74 efficiencies. This construct was converted into a minicircle (Supplementary  
75 Figure 3).

76

77 *-Transfection of cell-lines with CRISPR/Cas9 plasmids*

78 HEK 293 cells were seeded into a 6-well plate at a density of  $1 \times 10^6$   
79 cells/well. The following day, the cells were co-transfected with  $1 \mu\text{g}$  each of  
80 the sgRNA and HDR plasmids using calcium phosphate. The transfection  
81 mixture was left on the cells overnight and then replaced with fresh growth  
82 media. The cells were allowed to recover for 24h, then the selection agent  
83 (puromycin) was added. Selection was conducted for 7d, with the media  
84 (containing selection agent) replaced every 1-2d. The cells were then  
85 examined by PCR for the presence of the insertion.

86 Subsequently, cells that were positive for the integration were subjected to  
87 limiting dilution in order to create clonal cultures. Briefly, the cells were  
88 detached using trypsin, counted and a dilution made to produce a  
89 concentration of 100 cells/mL in 80% normal growth media and 20%  
90 conditioned media. These cells were then used to seed a 96 well plate  
91 ( $100 \mu\text{L}$ /well). The next day, the plate was examined and wells containing a

92 single cell were marked. The plate was then returned to the incubator and  
93 monitored for 7-10 days. When the positive wells reached a high density, they  
94 were subcultured and used to seed cluster plates with larger wells. This was  
95 repeated until the cells were growing in 6wp. Finally, once the density of the  
96 6-well plate was high enough, some cells were removed for gDNA extraction  
97 (Qiagen DNeasy Blood, Cells and Tissue kit) and PCR analysis.

98 A refinement of the above method involved using the NHEJ inhibitor  
99 SCR7. Briefly, the cells were seeded as above for transfection. The following  
100 day, the cells were treated with SCR7 2h before transfection. Importantly, the  
101 SCR7-containing media was not removed prior to the transfection. The  
102 transfection procedure was then conducted as described above.

103

#### 104 - *CRISPR/Cas9 editing of eHAP1 cells*

105 Since HEK293s select the WT alleles over those with the RASGRF2-int, we  
106 tried to edit a haploid cell line (eHAP1), using the same HDR synthesised  
107 plasmid as for the HEK293s. For the integration screening, we used two  
108 different PCRs. Each consisted of a primer binding within the  
109 newly introduced DNA, the other binding outside (recognizing the WT  
110 sequence). We succeeded in screening individually picked and expanded  
111 clones but unfortunately, none of these survived upon expansion to have  
112 integrated the HDR construct (Supplementary Figure 3). In our 2nd attempt,  
113 we first sub-cloned the HDR construct into a new backbone to remove the  
114 promoter and increase targeting efficiencies. This construct was converted  
115 into a minicircle, devoid of any bacterial backbone (improves HDR  
116 efficiencies) (more info at: [www.systembio.com/products/gene-expression-](http://www.systembio.com/products/gene-expression-)

117 systems/minicircle-technology/) eHap1 cells were targeted once more, using  
118 200.000 cells - 5 ug Cas9 protein - 1 ug sgRNA - 1 ug minicircle DNA. We  
119 Nucleofected the cells following the Horizon Discovery protocol (1575V - 10  
120 ms - 3 pulses). After nucleofection, cells were allowed to recover for 48h and  
121 then seeded at 5000 cells per 10 cm petri dish to grow into  
122 individual colonies. The remainder of cells (bulk) was used for  
123 triple genotyping PCR analysis. Unfortunately, the screening of the bulk DNA  
124 showed no sign of integrated DNA

125

#### 126 *-Controls of normal RASGRF2 transcription*

127 We used as normal controls for the RASGRF2 the unedited HEK293  
128 as well as edited HEK293s that “lost” the RASGRF2-int 48-days post-  
129 transfection (see also Figure 1). The first one is a “wild type” control, while the  
130 second one is a control that shows that transcription of RASGRF2 is restored  
131 when the integration is lost.

132

#### 133 *PCR screening for RASGRF2-int patient samples and qPCR on cell-lines*

134 To estimate the prevalence of the RASGRF2-polymorphic integration in  
135 PWIDs and non-PWIDs, we designed two pairs of primers (one common  
136 reverse), which specifically amplify either the allele with the integration alone  
137 (primers LTR\_splc\_F - dn\_intg\_R) or both the wild type and the integrated  
138 alleles, giving distinct PCR products (primers up\_intg\_F - dn\_intg\_R) (Dataset  
139 S1). This method allows the screening and the genotyping of multiple samples  
140 simultaneously (Figure 1, a-d).

141 We further investigated the impact of the HK2-LTR integration on the  
142 expression of the host RASGRF2 gene. We used CRISPR/Cas9 technology  
143 to edit the HEK293 cell line and to manually insert the LTR sequence in the  
144 pre-integration intronic site between exons 17 and 18 (see above). Before  
145 proceeding with the CRISPR/Cas9 engineering experiments we genotyped  
146 our 293s cell line (ECACC 85120602) and we found it negative for the  
147 RASGRF2 LTR integration. Using publicly available data from  
148 *www.hek293genome.org/v1/data.php* we retrieved the rs26907 SNP  
149 genotypes for most of the available 293s derivatives and we found them  
150 heterozygous. In addition, we did not find rs26907 in published Neanderthal  
151 genomes whereas RASGRF2-int has been identified, while the frequency of  
152 rs26907 in East Asian populations is similar to other populations, where  
153 RASGRF2-int has not been identified (see main text). These observations do  
154 not indicate linkage between rs26907 and RASGRF2-int..

155 The LTR and the flanking host sequences were previously confirmed  
156 with Sanger step-sequencing in patient samples positive for the integration.  
157 We confirmed the integration and we genotyped the cell line in multiple post-  
158 clonal-selection time-points using allele-specific PCR designs. The edited  
159 allele was present in the heterozygous state in the selected clones, as the  
160 wild-type allele-specific PCR was also positive after the transfection. The  
161 forward primer (intg\_conf\_F) of the primer-set used for the confirmation of the  
162 integration, was annealing upstream to the left homologous arm used in the  
163 HDR plasmid, allowing us to distinguish the edited alleles from the HDR  
164 plasmid (Figure 1).

165 The expression levels of RASGRF2 exons were evaluated using  
166 multiple SYBR-Green qPCRs (primers RAS2-3, 16-17, 17-18, 18-19). We  
167 found decreased expression levels of exons around the integration, which  
168 were reversed to normal after the selection of the wildtype alleles, 48 days  
169 post-transfection (T75 flask).

170

#### 171 *Assessment of HK2 LTR promoter in plasmids*

172 The HERVK-113 LTR *in vitro* synthesized dsDNA fragment (Invitrogen)  
173 was cloned into the HincII restriction site of pUC19 resulting HK-113-  
174 LTR/pUC19 recombinant plasmid. The HK-113-LTR XhoI excised fragment  
175 from HK-113-LTR/pUC19 plasmid was cloned into the XhoI restriction site of  
176 pGEM-luc vector (Promega) both in forward and reverse orientation resulting  
177 plasmids HK-113-LTR-Luc/F, HK-113-LTR-Luc/R respectively.

178 The human cell lines Huh7 (hepatocarcinoma) and HeLa (cervical  
179 carcinoma) were each maintained in Dulbecco's modified Eagle medium  
180 supplemented with 2 mM glutamine, 10% (v/v) heat-inactivated fetal calf  
181 serum and 100 U/ml penicillin/streptomycin.

182 For DNA transfections, Huh7 and HeLa cells were transiently  
183 transfected using JetPEI reagent (PolyPlus) according to the manufacturer's  
184 protocol. Following transfection, cells were harvested, seeded into 48-well  
185 plates, collected at the appropriate time points and subjected to luciferase  
186 assays with a commercially available kit (Promega). Luciferase activity was  
187 normalized to total cell protein in order to yield relative luciferase activity  
188 (RLA).

189 Huh7 and HeLa cells were transfected with pGEM-luc vector (control)  
190 and HK-113-LTR Luc/F, HK-113-LTR-Luc/R plasmids. 48h post-transfection,  
191 cells were lysed and luciferase activity was measured; luciferase readouts  
192 were normalised to total cell protein.

193

#### 194 *Expression patterns of read-through transcripts*

195



196 We found significantly increased expression levels over the LTR-ends  
197 compared to the LTR-starts (Wilcoxon,  $p < 0.05$ ) (Figure 2). We also found that  
198 elements with increased expression levels over the LTR-start would present  
199 decreased expression (or would be inactive) over the LTR-end and vice-  
200 versa. (Figure 2)

201 We tested the expression levels in combination with the type but also  
202 the estimated age of each element as described in Subramanian et al (1). We  
203 found that full-length elements are significantly more active compared to the  
204 solo-LTR integrations, over their LTR-end edges (Wilcoxon,  $p < 0.05$ ) (Figure  
205 2). We asked for at least 5 reads per fragment to evaluate a HK2-edge as  
206 active and we estimated their expression (active vs. not-active) and the  
207 magnitude ( $< 500$  vs.  $> 500$  reads) per edge, in combination with their age  
208 (younger than 5my, 5-20my and older than 20my). We found that the  
209 youngest elements (younger than 5my) are more likely to be active versus the  
210 very old ones (older than 20my) over the LTR-ends. The same trend was  
211 observed for the magnitude of the expression, with the youngest elements to  
212 be more expressed versus the older ones, both over the LTR-start and the  
213 LTR-end regions (z-test, one-tailed,  $p < 0.05$ ) (Figure 2).

214

215

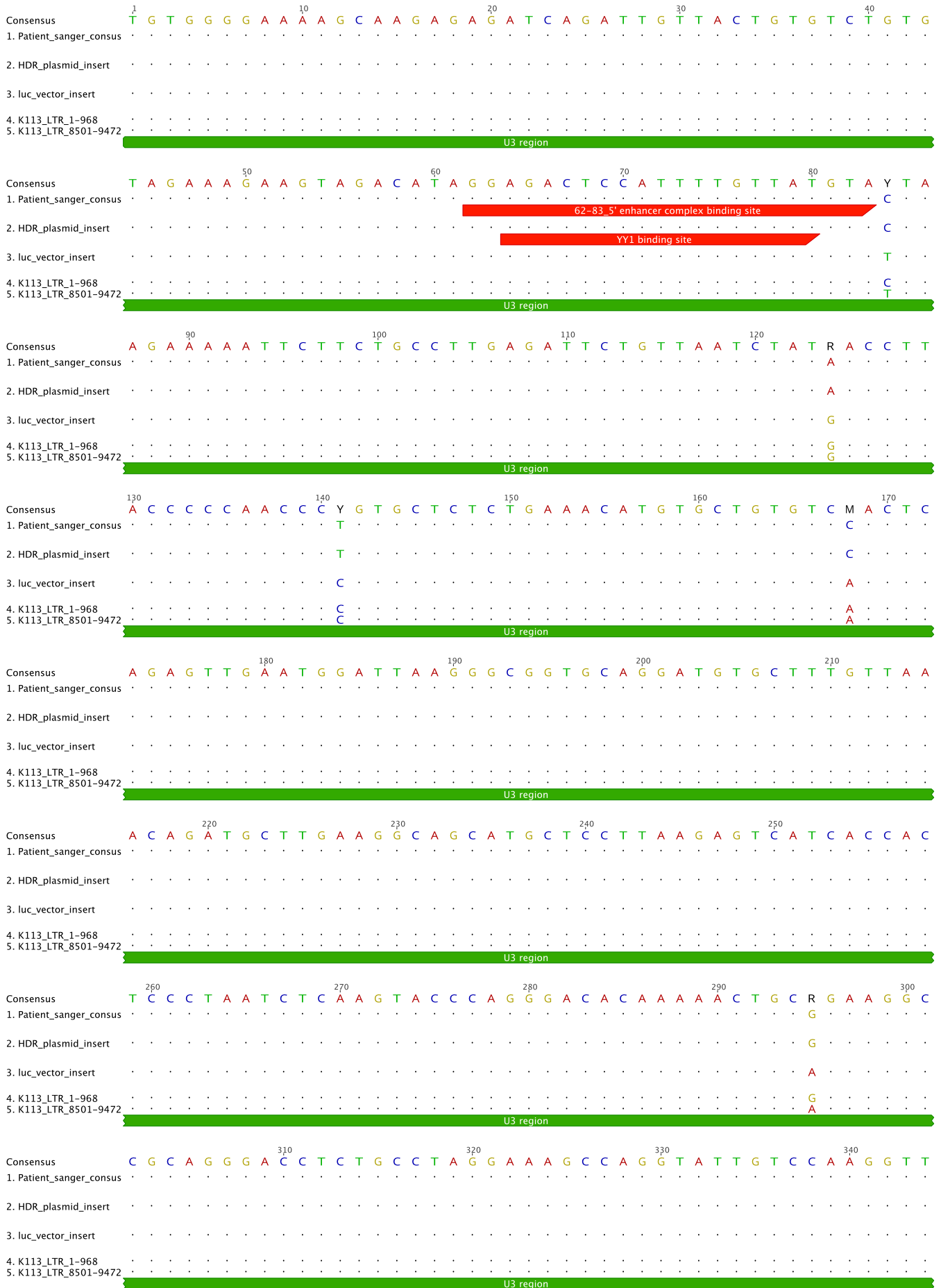
216

217 **Supplementary Figures**

218

219 **Supplementary Figure 1:** Nucleotide alignment and transcription regulatory  
220 motifs annotation of: (1) Sanger step sequencing consensus of the RASGRF2  
221 – integrated solo LTR. (2) Synthetic fragment inserted into the HDR plasmid,  
222 which was used for the CRISPR/Cas9-mediated editing of the HEK293 cells  
223 (identical to 1). (3) Synthetic fragment used for the Luciferase expression  
224 vector construction. (4, 5) HERV-K 113 LTRs, regions: 1-968bp and 8501-  
225 9472bp, respectively. The K113 LTR sequences only differs by <1%  
226 compared to the RASGRF2 solo LTR integration, while none of the observed  
227 sequence differences is located within promoter/regulatory sites, e.g. the four  
228 GC-boxes which are targets of the transcription factors Sp1 and Sp3 (2), the  
229 transcription initiation site (Inr), the TATA boxes -like and the Beta globin  
230 distal promoter -like sequence or the binding site of the transcription factor  
231 OTF-2 (3) or the 5'-enhancer complex binding site (nt 62-83) and the YY1  
232 transcription factor (nt 64-80) (4). The U3, R and U5 LTR regions are  
233 indicated in green, red and blue arrows respectively.

234



Consensus T C T C C <sup>350</sup> C C A T G T G A T A <sup>360</sup> G T C T G A A A T A <sup>370</sup> T G G C C T C G T G <sup>380</sup> G G A A G G G A

1. Patient\_sanger\_consus . . . . .

2. HDR\_plasmid\_insert . . . . .

3. luc\_vector\_insert . . . . .

4. K113\_LTR\_1-968 . . . . .

5. K113\_LTR\_8501-9472 . . . . .

U3 region

Consensus A A G A C C T G A C C G <sup>400</sup> T C C C C C A G C C C <sup>410</sup> G A C A C C C G T A <sup>420</sup> A A G G G T C T G T <sup>430</sup>

1. Patient\_sanger\_consus . . . . .

2. HDR\_plasmid\_insert . . . . .

3. luc\_vector\_insert . . . . .

4. K113\_LTR\_1-968 . . . . .

5. K113\_LTR\_8501-9472 . . . . .

U3 region

Consensus G C T G A G G A G <sup>440</sup> G A T T A G T A W <sup>450</sup> A A G A G G A A G G A <sup>460</sup> A T G C C T C T T G <sup>470</sup> C A G T

1. Patient\_sanger\_consus . . . . .

2. HDR\_plasmid\_insert . . . . .

3. luc\_vector\_insert . . . . .

4. K113\_LTR\_1-968 . . . . .

5. K113\_LTR\_8501-9472 . . . . .

U3 region

Consensus T G A G A C <sup>480</sup> A A G A G G A A G G C <sup>490</sup> A T C T G T C T C <sup>500</sup> C T S C C T G T C C <sup>510</sup> C T G G C A

1. Patient\_sanger\_consus . . . . .

2. HDR\_plasmid\_insert . . . . .

3. luc\_vector\_insert . . . . .

4. K113\_LTR\_1-968 . . . . .

5. K113\_LTR\_8501-9472 . . . . .

U3 region

Consensus A T G G A A T G T C T C G <sup>520</sup> G T A T A A A C C C <sup>530</sup> G A T T G T A T G C T C C A T C T A C <sup>540</sup> <sup>550</sup>

1. Patient\_sanger\_consus . . . . .

TATA box-like

2. HDR\_plasmid\_insert . . . . .

3. luc\_vector\_insert . . . . .

4. K113\_LTR\_1-968 . . . . .

5. K113\_LTR\_8501-9472 . . . . .

U3 region

Consensus T G A G A T A G G G <sup>560</sup> A A A A C C G C C T Y A G G G C T G G <sup>570</sup> A G G T G G A C C T G C <sup>580</sup> <sup>590</sup> <sup>600</sup>

1. Patient\_sanger\_consus . . . . .

GC-box1

2. HDR\_plasmid\_insert . . . . .

3. luc\_vector\_insert . . . . .

4. K113\_LTR\_1-968 . . . . .

5. K113\_LTR\_8501-9472 . . . . .

R region

Consensus G G G C A G C <sup>610</sup> A A T A C T G C T T G T A A A G C A Y <sup>620</sup> T G A G A T G T T T <sup>630</sup> <sup>640</sup>

1. Patient\_sanger\_consus . . . . .

GC-box?

2. HDR\_plasmid\_insert . . . . .

3. luc\_vector\_insert . . . . .

4. K113\_LTR\_1-968 . . . . .

5. K113\_LTR\_8501-9472 . . . . .

R region

Consensus A T G C <sup>650</sup> A T A T C T A A A A G <sup>660</sup> C A C A G C A C T T A A T C C T T T A C <sup>670</sup> A T T G T C T A <sup>680</sup>

1. Patient\_sanger\_consus . . . . .

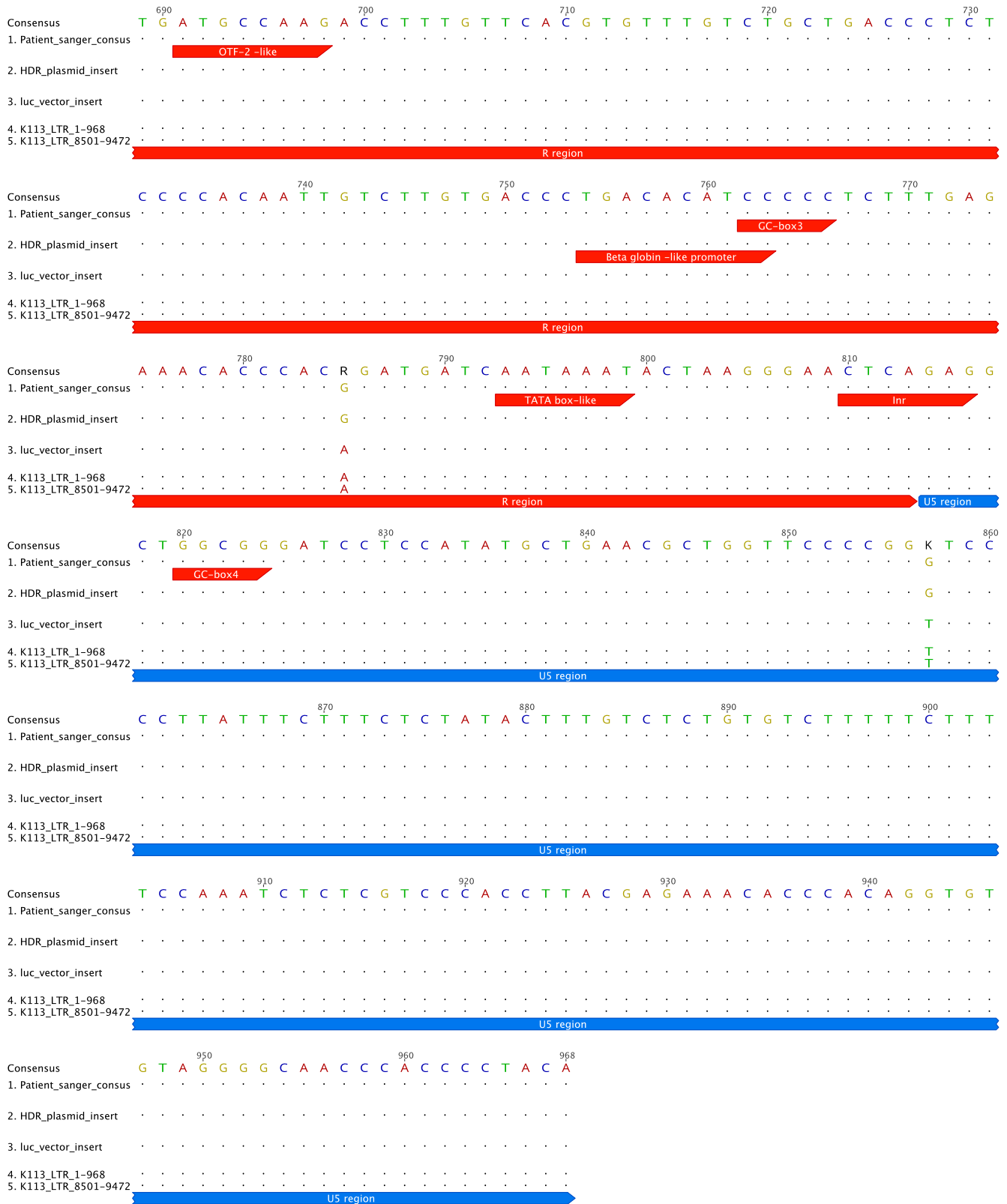
2. HDR\_plasmid\_insert . . . . .

3. luc\_vector\_insert . . . . .

4. K113\_LTR\_1-968 . . . . .

5. K113\_LTR\_8501-9472 . . . . .

R region



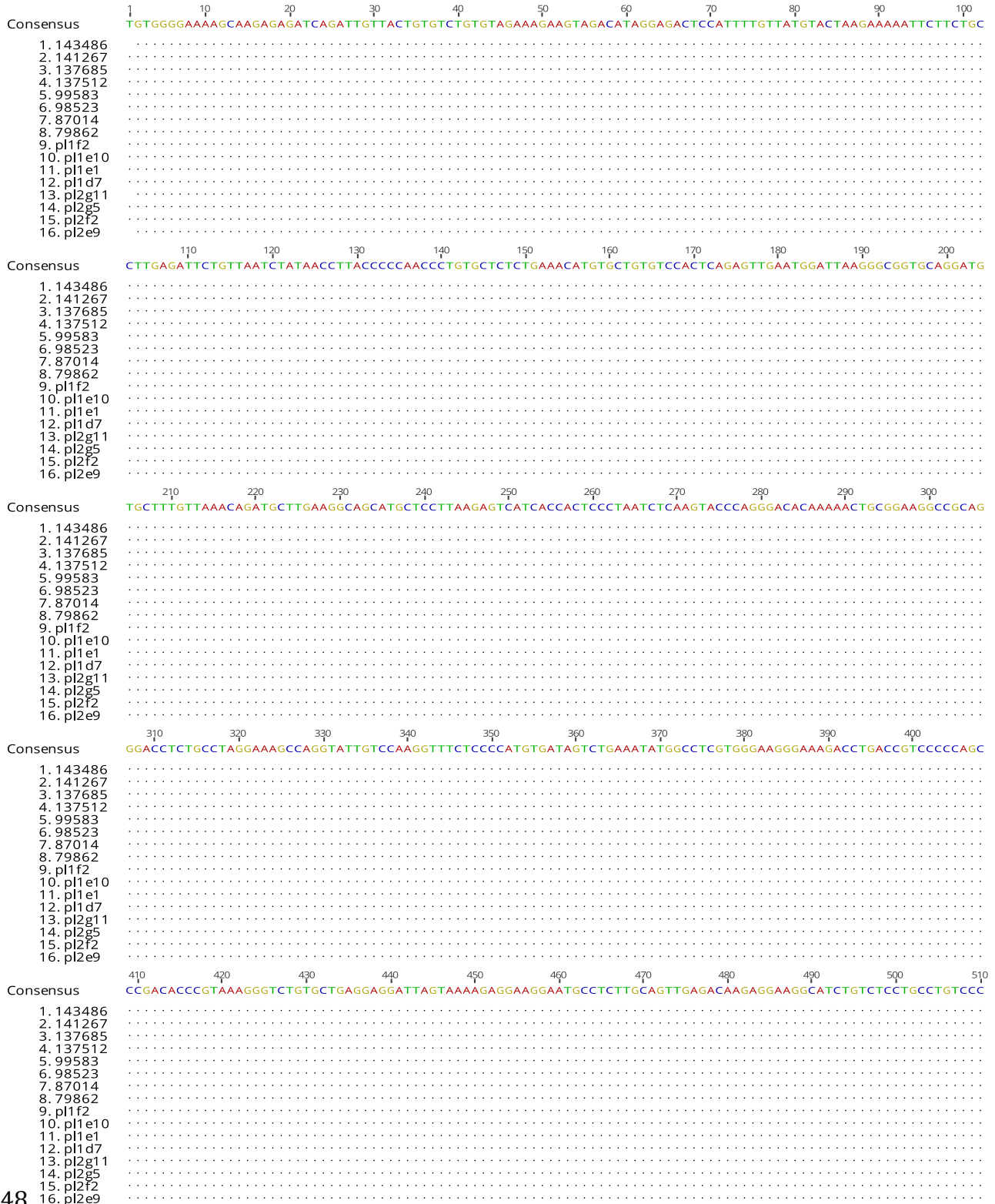
238

239

240

241

242 **Supplementary Figure 2:** Nucleotide alignment of 16 sanger-sequenced  
 243 RASGRF2-int LTR positive patients, 4 from each group: 1-4: GR PWID  
 244 (Greek Persons Who Inject Drugs),5-8: GR non-PWID (Greek non-PWID), 9-  
 245 12: UK PWID (British Persons Who Inject Drugs),13-16: UK non-PWID. All  
 246 samples were found identical.  
 247



248  
 249

Consensus TGGGCAATGGAATGCTCTCGGTATAAAACCCGATTGTATGCTCCATCTACTGAGATAGGGAAAAACCGCCTTAGGGCTGGAGGTGGGACCTGC GGGCAGCAAT

1. 143486  
2. 141267  
3. 137685  
4. 137512  
5. 99583  
6. 98523  
7. 87014  
8. 79862  
9. pl1f2  
10. pl1e10  
11. pl1e1  
12. pl1d7  
13. pl2g1.1  
14. pl2g5  
15. pl2f2  
16. pl2e9

Consensus ACTGCTTTGTAAGCACTGAGATGTTTATGTTGATGCATATCTAAAAGCACAGCAGCCTTAATCCCTTTACATTGTTCTATGATGCCAAGACCTTTGTTTACAGTGT

1. 143486  
2. 141267  
3. 137685  
4. 137512  
5. 99583  
6. 98523  
7. 87014  
8. 79862  
9. pl1f2  
10. pl1e10  
11. pl1e1  
12. pl1d7  
13. pl2g1.1  
14. pl2g5  
15. pl2f2  
16. pl2e9

Consensus TTGTCTGCTGACCCCTCTCCCAACAATTGTCTTGTGACCCGTGACACATCCCCCTCTTTGAGAAACACCCACGGATGATCAATAAATACTAAGGGAACCTCAGAG

1. 143486  
2. 141267  
3. 137685  
4. 137512  
5. 99583  
6. 98523  
7. 87014  
8. 79862  
9. pl1f2  
10. pl1e10  
11. pl1e1  
12. pl1d7  
13. pl2g1.1  
14. pl2g5  
15. pl2f2  
16. pl2e9

Consensus GCTGGCGGGATCCTCCATATGCTGAACGCTGGTCCCGGGTCCCTTATTTCTTTCTCTATACTTTGTCTCTGTGTCTTTTTCTTTTCCAAATCTCTCGTC

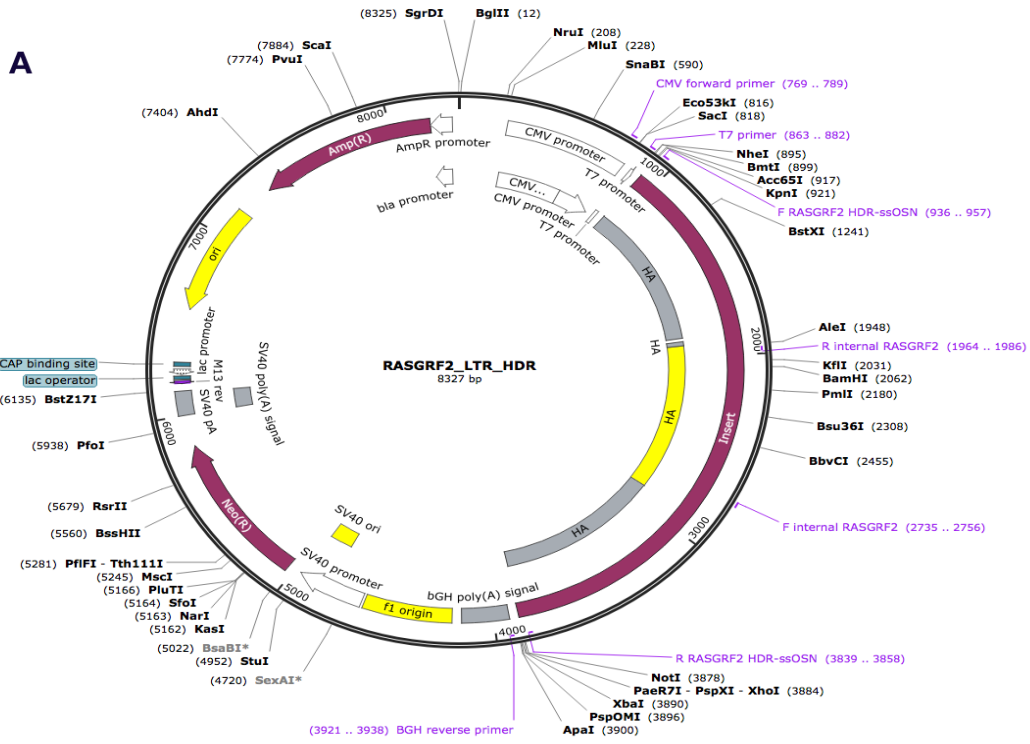
1. 143486  
2. 141267  
3. 137685  
4. 137512  
5. 99583  
6. 98523  
7. 87014  
8. 79862  
9. pl1f2  
10. pl1e10  
11. pl1e1  
12. pl1d7  
13. pl2g1.1  
14. pl2g5  
15. pl2f2  
16. pl2e9

Consensus CCACCTTACGAGAAACACCCACAGGTGTGTAGGGGCAACCCACCCCTACA

1. 143486  
2. 141267  
3. 137685  
4. 137512  
5. 99583  
6. 98523  
7. 87014  
8. 79862  
9. pl1f2  
10. pl1e10  
11. pl1e1  
12. pl1d7  
13. pl2g1.1  
14. pl2g5  
15. pl2f2  
16. pl2e9

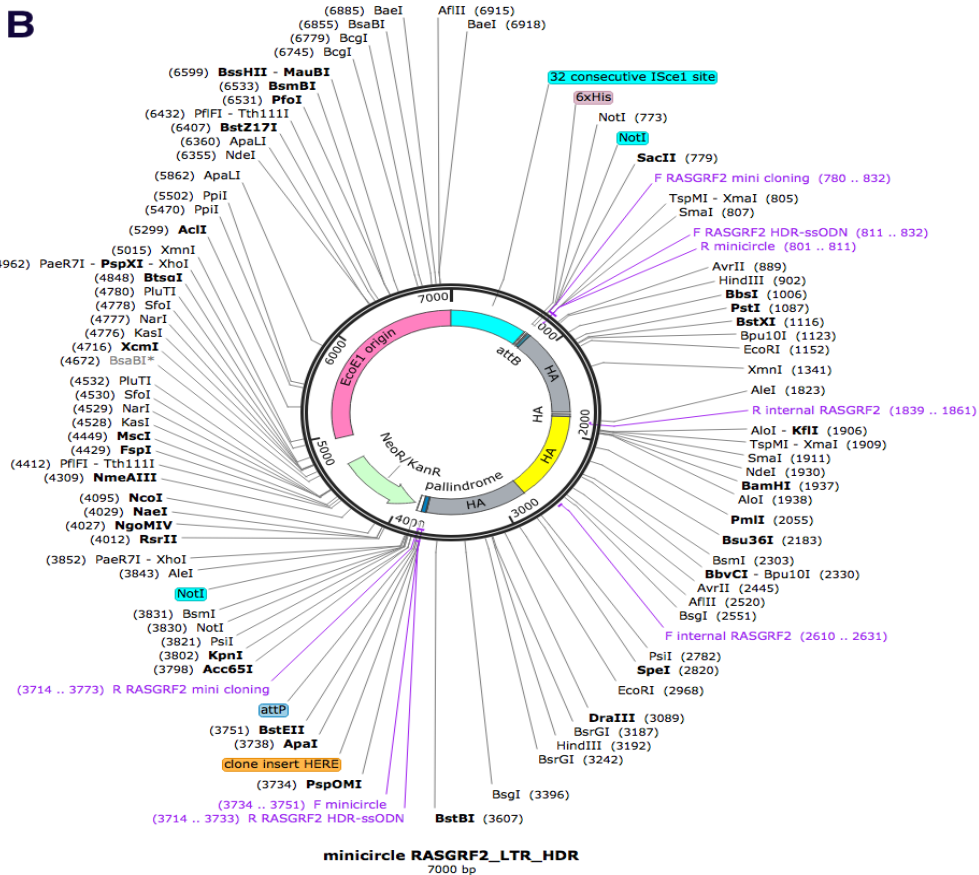
252 **Supplementary Figure 3:** Plasmid map of (A) the HDR vector used for the  
253 transfection of HEK293s and the first attempt of eHAP1s and (B) the  
254 minicircle construct used for the second transfection attempt of eHAP1s.

Created with SnapGene®



255

Created with SnapGene®



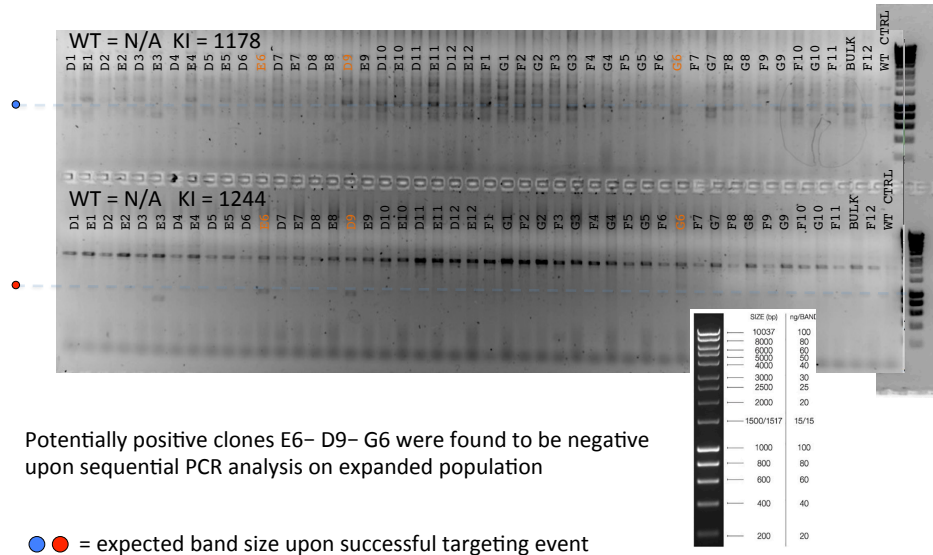
256



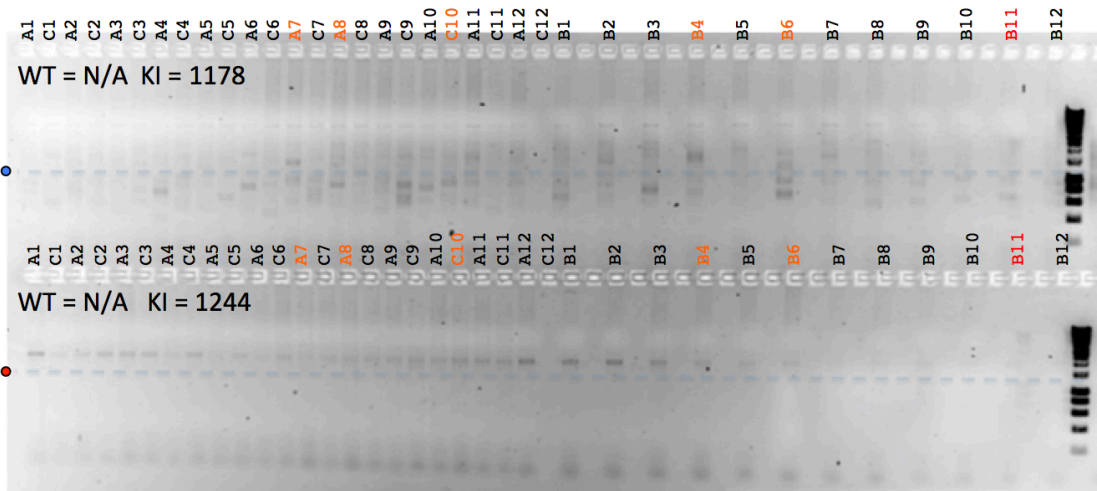
257 **Supplementary Figure 4: CRISPR/Cas9 editing of eHAP1 cells using (A) the**  
258 **same HDR vector as for the HEK293s and (B) the minicircle sub-construct**

**(A) RASGRF2 HERVK Integration (KI) PCR (KAPA2G 4kb) – 1<sup>st</sup> round**

Conditions: 1 ul gDNA per 25 ul PCR – 1.25 ul of each 10 uM primer  
Performed two different PCRs per sample



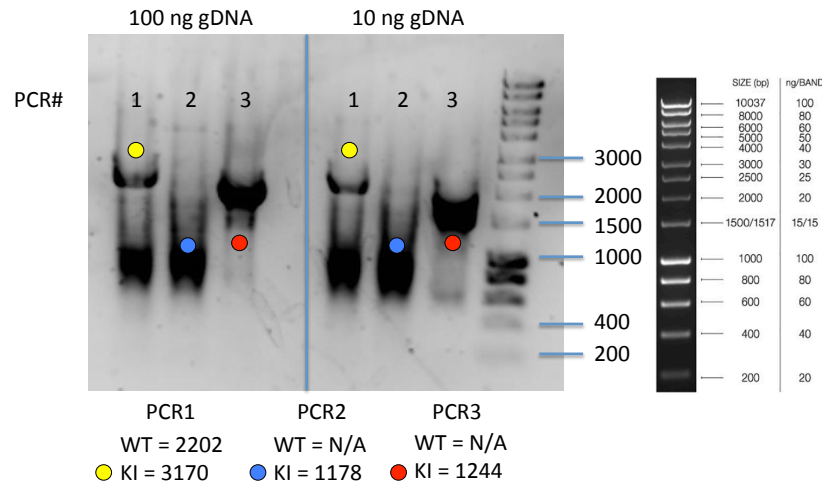
259  
260



261  
262  
263

**(B) RASGRF2 HERVK Integration (KI) PCR (KAPA2G 4kb) – 2<sup>nd</sup> round (minicircle)**

Conditions: 1 ul gDNA (bulk of targeted cells) per 25 ul PCR – 1.25 ul of each 10 uM primer  
 Tm = 60 deg



● ● ● = expected band size upon successful targeting event

264  
 265

266 **References**

267

- 268 1. Subramanian RP, Wildschutte JH, Russo C, & Coffin JM (2011)  
269 Identification, characterization, and comparative genomic distribution of  
270 the HERV-K (HML-2) group of human endogenous retroviruses.  
271 *Retrovirology* 8.
- 272 2. Fuchs NV, *et al.* (2011) Expression of the human endogenous retrovirus  
273 (HERV) group HML-2/HERV-K does not depend on canonical promoter  
274 elements but is regulated by transcription factors Sp1 and Sp3. *J Virol*  
275 85(7):3436-3448.
- 276 3. Kovalskaya E, Buzdin A, Gogvadze E, Vinogradova T, & Sverdlov E (2006)  
277 Functional human endogenous retroviral LTR transcription start sites are  
278 located between the R and U5 regions. *Virology* 346(2):373-378.
- 279 4. Knossl M, Lower R, & Lower J (1999) Expression of the human  
280 endogenous retrovirus HTDV/HERV-K is enhanced by cellular  
281 transcription factor YY1. *J Virol* 73(2):1254-1261.

282