



HAL
open science

The C-terminal domain of piggyBac transposase is not required for DNA transposition

Laura Helou, Linda Beauclair, Hugues Dardente, Peter Arensburger, Nicolas Buisine, Yan Jaszczyszyn, Florian Guillou, Thierry Lecomte, Alex Kentsis, Yves Bigot

► To cite this version:

Laura Helou, Linda Beauclair, Hugues Dardente, Peter Arensburger, Nicolas Buisine, et al.. The C-terminal domain of piggyBac transposase is not required for DNA transposition. *Journal of Molecular Biology*, 2021, 433 (7), pp.1-20. 10.1016/j.jmb.2020.166805 . hal-03115098

HAL Id: hal-03115098

<https://hal.science/hal-03115098v1>

Submitted on 26 Nov 2021

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

1 The C-terminal domain of *piggyBac* transposase is not required for DNA transposition

2

3

4 Laura Helou¹, Linda Beauclair¹, Hugues Dardente¹, Peter Arensburger², Nicolas Buisine³, Yan
5 Jaszczyszyn⁴, Florian Guillou¹, Thierry Lecomte⁵, Alex Kentsis^{6,7,8} and Yves Bigot^{1,*}

6

7

8 ¹ PRC, UMR INRAE 0085, CNRS 7247, Centre INRAE Val de Loire, 37380 Nouzilly, France

9 ² Biological Sciences Department, California State Polytechnic University, Pomona, CA 91768,

10 United States of America

11 ³ UMR CNRS 7221, Muséum National d'Histoire Naturelle, 75005 Paris, France

12 ⁴ Université Paris-Saclay, CEA, CNRS, Institute for Integrative Biology of the Cell (I2BC),

13 91198, Gif-sur-Yvette, France

14 ⁵ EA GICC 7501, CHRU de Tours, 37044 Tours Cedex 09, France

15 ⁶ Molecular Pharmacology Program, Sloan Kettering Institute, Memorial Sloan Kettering Cancer

16 Center, New York, New York, USA

17 ⁷ Weill Cornell Medical College, Cornell University, New York, New York, USA

18 ⁸ Department of Pediatrics, Memorial Sloan Kettering Cancer Center, New York, New York, USA

19

20

*Corresponding author address: PRC, UMR INRA 0085, CNRS 7247, 37380 Nouzilly, France. Tel:

+33 2 47 42 75 66, e-mail: yves.bigot@inrae.fr

22

23 **Highlights**

1 24 The C-terminal CRD in *pble* transposases is not essential for transposition

2 25 Two CRD-deficient *pble* transposases trigger transposition of *Ifp2*

3
4 26 Proper and improper insertions occur when CRD-deficient transposases mediate mobility

5
6 27 CRD-deficient and full-length *pble* transposases do not insert transposons at random

7
8 28 Features of the domesticated transposase PGBD5 originate from wild type transposase

9 29

10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

30 **Abstract**

31 *PiggyBac*(PB)-like elements (*pble*) are members of a eukaryotic DNA transposon family. This
32 family is of interest to evolutionary genomics because *pble* transposases have been domesticated at
33 least 9 times in vertebrates. The amino acid sequence of *pble* transposases can be split into three
34 regions: an acidic N-terminal domain (~100 aa), a central domain (~400 aa) containing a DD[D/E]
35 catalytic triad, and a cysteine-rich domain (CRD; ~90 aa). Two recent reports suggested that a
36 functional CRD is required for *pble* transposase activity. Here we found that two CRD-deficient
37 *pble* transposases, a PB variant and an isoform encoded by the domesticated PB-derived vertebrate
38 transposase gene 5 (*pgbd5*) trigger transposition of the *Ifp2 pble*. When overexpressed in HeLa cells,
39 these CRD-deficient transposases can insert *Ifp2* elements with proper and improper transposon
40 ends, associated with deleterious effects on cells. Finally, we found that mouse CRD-deficient
41 transposase *Pgbd5*, as well as PB, do not insert *pbles* at random into chromosomes. Transposition
42 events occurred more often in genic regions, in the neighbourhood of the transcription start sites and
43 were often found in genes predominantly expressed in the human central nervous system.

44
45
46
47
48
49
50
51 **Keywords:** transposon / DNA cleavage / neuron / insertion preference / vertebrate

53 **Abbreviations:**

54 CRD: cysteine-rich domain

55 gDNA: genomic DNA

56 GFP: green fluorescent protein

57 ISCR: insertion sites-containing regions

58 NeoR : neomycin resistance

59 NLS: nuclear localisation signal

60 ORF: open reading frame

61 PB: *piggyBac* transposase

62 *pble*: *piggyBac*-like element

63 PGBD or *pgbd*: “*piggyBac* derived transposase” protein or gene

64 STIR: sub-terminal inverted repeats

65 SV40: simian virus 40

66 TE: transposable element

67 TIR: terminal inverted repeats

68 TSD: target site duplication

69

70 Introduction

71 Transposable elements (TEs) gather diverse discrete DNA sequences of prokaryotic and eukaryotic
72 origins that use a wide range of mobility mechanisms to transpose within the genome of their hosts
73 [1-4]. *PiggyBac-like* elements (*pble*) consist of a family of DNA transposons that have so far only
74 been found in animal genomes and with copy numbers that vary widely between host species [5].
75 *Pbles* are able to jump from one chromosomal locus to another using cut-and-paste transposition
76 which is enzymatically catalysed by the transposase they encode. The first *pble* to be identified was
77 *Ifp2* (a.k.a. *piggyBac*) from the cabbage looper moth *Trichoplusia ni* (Lepidoptera) [6]. It is the
78 reference element in the *piggyBac* family for academic research purposes [7]. The *Ifp2* DNA
79 sequence is 2476 bp in length and contains an open reading frame (ORF) coding a 594 amino acid
80 transposase named PB. The *Ifp2* DNA sequence is flanked by 13 bp long terminal inverted repeats
81 (TIR) and by 19 bp long subterminal inverted repeats (STIR) located internally at 3 and 31
82 nucleotides of distance of 5' and 3' TIR inner ends, respectively. PB excises the transposon from its
83 donor chromosomal locus and reinserts it into a TTAA motif, which gets duplicated upon insertion.
84 Sequence analyses revealed that PB contains at least 3 domains. The first domain spans from
85 residues 1 to 116, shows no overt structural features and displays an acidic pI of 4.41. The second
86 domain is a macro domain that extends from residues 117 to 535 and has a basic pI (9.29); it contains
87 several highly conserved residues, including the predicted catalytic residues (D268, D346 and
88 D447) that are required for all transposition steps. On both sides of the catalytic domain (residues
89 263 to 457), recent structures obtained by cryo-electron-microscopy [8] revealed two DNA binding
90 sub-domains (residues 117 to 263 and 457 to 535, respectively) that bind to *Ifp2* TIRs. Finally, the
91 third domain spans from residues 559 to 594, and contains a cysteine-rich domain (CRD, pI=9.07)
92 for which the atomic structure was first solved by nuclear magnetic resonance [9]. This CRD was
93 shown to bind to a 5'-TGCGT-3'/3'-ACGCA-5' motif that is contained within the 19 bp subterminal
94 inverted repeat between positions 178 and 199 of the *Ifp2* sequence [9]. The CRD was also found
95 to be vital for the dimerization of PB [10], the removal of the last seven residues being sufficient to
96 yield a monomeric protein that still binds to *Ifp2* ends *in vitro*. Finally, this CRD contains a nuclear
97 localisation signal (NLS) that is required to mediate PB nuclear localisation [11].
98 Two studies [9,10] proposed non-exclusive roles for the CRD and the ability of PB to mediate
99 transposition. In the first study, the authors concluded that the CRD was essential for *Ifp2*
100 transposition because the use of a CRD-deficient PB (PB.1-558) in an integration assay performed
101 in mammalian cells did not lead to an increase in integration activity compared to controls done in
102 the absence of the transposase [9]. The second study proposed that PB dimerisation might serve to
103 prevent excessive transposition of *Ifp2* since removal of the CRD led the monomeric PB to be more
104 active in transposition excision [10]. It is also possible that the absence of the CRD may cause

105 cytosolic retention since it contains a nuclear localization sequence (NLS).

106 Here we aimed to determine whether PB.1-558 fused to a simian virus 40 (SV40) NLS was able to
107 mediate the transposition of *Ifp2* into human cell chromosomes. First, we showed that both PB.1-
108 558 and PB.NLS-1-558 had a negative effect on obtaining clones in integration assays. Second, we
109 found that PB.NLS-1-558 was able to carry out the transposition of *Ifp2* elements. Third, we
110 compared the properties of PB.NLS-1-558 with those of another CRD-deficient *piggyBac* protein,
111 the domesticated murine and human PGBD5 protein. Fourth, we evaluated the quality of *Ifp2* ends
112 neo-integrated into chromosomes by PB and the two CRD-deficient proteins, PB.NLS-1-558 and
113 PGBD5. Finally, we determined whether the three proteins used similar or different pools of
114 chromosomal insertion sites, whatever the state of *Ifp2* transposon ends.

115

116 **Results**

117 **PB.1-558 needs an NLS to locate into nuclei.**

118 We made a CRD-deficient PB mutant (PB.1-558) and a construct in which its N-terminal end was
119 fused with a SV40 NLS (PB.NLS-1-558). Such a position for the NLS does not modify the
120 transposition activity of PB and preserves the activity of the added localization motif or protein
121 domain [12,13]. To assess cellular localization of both proteins we made two more constructs in
122 which the green fluorescent protein (GFP) was C-terminally fused to PB.1-558 and PB.NLS-1-558.
123 An expression vector encoding GFP was used as a diffusion control within the cytoplasm and
124 nucleus (Fig. 1a-c) and a vector encoding a PB-GFP fusion was used as a control for active import
125 into the nuclei (Fig. 1d-f) [11]. Our data revealed that PB.1-558 is not enriched in the nucleus (Fig.
126 1g-i), which contrasts with its PB.NLS-1-558 counterpart that is almost completely nuclear (Fig. 1j-
127 l).

128

129 **PB variants display cytotoxicity**

130 Prior to assaying the ability of the full-length PB, PB.1-558, and PB.NLS-1-558 to trigger
131 transposition of an *Ifp2* source, we checked whether these proteins impacted random integration
132 rates into HeLa cell chromosomes. We used a DNA plasmid containing a gene cassette coding for
133 a neomycin resistance (NeoR) without the *Ifp2* sequence, the pBSK-NeoR plasmid. This was
134 performed using a classic integration assay (see material and methods section). We observed that
135 the integration rate of pBSK-NeoR into chromosomes was significantly lower in the presence of
136 each of the three variants compared to a control GFP sequence (Fig. 2a; ~1.25, 1.31, 2.67 folds
137 (1/fold change) for PB, PB.1-558, and PB.NLS-1-558, respectively). This indicated that PB and its
138 variants have a deleterious effect on cells, including those that display random and stable
139 integrations of NeoR into their chromosomes. The cytotoxicity of PB.NLS-1-558 is higher than that

140 of PB and PB.1-558.

141 The results of integration assays (Fig. 2b) performed with the *Ifp2*-NeoR transposon donor plasmid
142 confirmed that PB.1-558 has a negative effect on obtaining NeoR clones. The number of clones was
143 1.5 times lower than that obtained with the GFP control and 60 times lower than that obtained with
144 PB. This indicated that integration assays performed with PB.1-558 are affected by both: 1) the
145 cytotoxicity of PB.1-558, and 2) the integration of the NeoR cassette into chromosomes.

146 The number of NeoR clones obtained with PB.NLS-1-558 was 2.5 times higher than that with the
147 GFP control. Due to the cytotoxicity of PB.NLS-1-558, this number was likely underestimated.
148 After correcting for PB.1-558 toxicity rate (Fig. 2c) we estimate the number of integration events to
149 be ~7 times higher what is found with the GFP control. Overall, this means that PB.NLS-1-558 is
150 roughly 8 times less efficient than PB for obtaining NeoR clones in an integration assay done in
151 HeLa cells. Because cytotoxicity hampered our ability to directly evaluate integration rates and was
152 likely dose-dependent as observed with PGBD5 [14], we focused our investigation on the ability of
153 PB.NLS-1-558 to trigger *Ifp2* transposition by characterizing integration events into chromosomes
154 by NeoR clones.

155 156 **Features of sites targeted by PB and PB.NLS-1-558 when integrating *Ifp2* into chromosomes**

157 To verify the presence of transposition events and to determine their sequence features, we produced
158 fragment populations corresponding to *Ifp2*-chromosome junctions. These were made by LAM-
159 PCR using genomic DNA (gDNA) of NeoR clone populations that were sequenced using Illumina
160 Miseq technology. To prepare gDNA samples we used ~60000 clones from integration assays done
161 with *Ifp2*-NeoR and PB, and ~1000 clones from integration assays done with *Ifp2*-NeoR and
162 PB.NLS-1-558. Previous results of integration assays performed in HEK293 cells found that the
163 rate of *Ifp2* integration into chromosomes by proper transposition (i.e. with a perfect duplication of
164 the “TTAA” TSD and conservation of the TIR sequence) was about 96-98% when PB was used as
165 a transposase source [15,16].

166 Using DNA sequence alignments, we characterized 7623 *Ifp2*/chromosome junctions resulting from
167 integration events mediated by PB and 516 junctions mediated by PB.NLS-1-558, with *Ifp2*-NeoR
168 as a transposon source (supplementary Table 1a and 2a). Sequenced junctions at the 5' and the 3'
169 of *Ifp2* ends were not equally represented in the sequence data, likely because of efficiency
170 differences at certain steps of DNA fragment amplification during the LAM-PCR. Sequence
171 junctions were further examined taking into account the conservation of TSD and TIR sequences,
172 two features that were required to keep the capacity of neo-inserted elements to be efficiently
173 remobilized during excision and insertion, i.e. to remain “active in transposition” [15,16]. Four
174 kinds of junctions were observed: those displaying i) a full TIR sequence and a TTAA TSD (red

175 bars from positions 101 to 104 and 2222 to 2225 in Fig. 3), ii) a region containing an intact TIR and
176 a TTAA TSD juxtaposed to a little piece of plasmid backbone (black bars from positions 1 to 100
177 and 2226 to 2301 in Fig. 3a and b), iii) no TTAA TSD but a full TIR sequence (blue bars from
178 positions 102 to 105 and 2218 to 2221 in Fig. 3), and iv) no TIR sequence lacking one or several
179 nucleotides at its outer end (black bars from positions 107 to 178 and 2147 to 2217 in Fig. 3). The
180 summary of results in Table 1 indicates that the rate of proper events when PB was used as a
181 transposase was similar to that previously observed in other cell types, but 3.3% of junctions
182 nevertheless displayed improper TSD, TIRs or both. Interestingly, 19.0% of junctions mediated by
183 PB.NLS-1-558 were found to be proper, thus demonstrating that this variant is able to trigger
184 canonical transposition events even though less efficiently than PB. Our results also indicated that
185 PB.NLS-1-558 integrated *Ifp2* into non-canonical TSDs approximately 25 times more often than
186 PB. Furthermore, TIRs were damaged or accompanied by a piece of backbone sequence of variable
187 length juxtaposing the transposon in the transposon donor plasmid in about 65% of
188 *Ifp2*/chromosome junctions (while they represented only 2.25% of junctions among integration
189 events triggered by PB). These observations suggested that the observed junctions resulted from
190 both proper transposition events and improper integration events that could be mediated by both PB
191 variants, but the rates of each kind of integration events were dramatically different between the two
192 proteins.

193 Next, we identified 5' and 3' junctions for events that occurred exactly at the same chromosomal
194 insertion sites in both datasets. We observed that 7446 chromosomal sites were used and found 177
195 unambiguous insertion sites in our Lumpy raw file that were occupied several times. In these 177
196 sites, integration events occurred in both *Ifp2* orientations when mobility was mediated by PB
197 (supplementary Table 1b). A careful examination of the resulting bam file with IGV [17] revealed
198 11 cases of putative single integration events (supplementary Table 1b, case highlighted in cyan
199 blue). Among them, three corresponded to *Ifp2* transposons displaying at least one TIR damaged at
200 its outer end, and one TIR was inserted into a duplicated TSD corresponding to a duplicated CATG
201 motif. As previously described [18,19] we also found four sites in which both integration events
202 occurred into non-canonical TSD (CATG, TATC, ACAT, TTCC; supplementary Table 1b) and 16
203 sites where two events occurred with the insertion of transposons with at least one improper end.
204 These results suggest that virtually any type of non-canonical integration can be found at a very low
205 frequency when *Ifp2* was transposed by PB. In data resulting from the transposition of *Ifp2* by
206 PB.NLS-1-558, we found that 516 chromosomal sites were used and identified 32 unambiguous
207 insertion sites for which integration events occurred in both orientations of the transposon. Four of
208 these putatively corresponded to single integration events (supplementary Table 2b, case
209 highlighted in cyan blue, two would correspond to canonical integrations by transposition and two

210 with non-canonical TIR or TSD). The main difference with PB is that improper events were
211 dramatically more frequent when transposition was mediated by PB.NLS-1-558.

212

213

4213 **PGBD5 a natural domesticated CRD-deficient *pble* transposase**

5
6214 We compared the transposition features of the PB.NLS-1-558 variant to those of murine and human
7
8215 orthologues of the oldest domesticated *piggyBac* transposase since the origin of vertebrates, PGBD5
9
10216 (Mm523 and Hs524) [20]. Alignment of three protein sequences (Fig. 4) revealed that both CRD-
11
12217 deficient proteins displayed an acidic N-terminal domain and a second domain with a basic pI (~9.2)
13
14218 containing an apparent catalytic triad composed of 3 acidic amino acid residues that were essential
15
16219 for transposition activity [20]. Another shared feature was their ability to trigger *Ifp2* transposition
17
18220 [14,21]. This transposition ability was rather unexpected for PGBD5 compared to PB.NLS-1-558
19
20221 because the PGDB5 catalytic triad was not located at the same positions as in *pble* transposases
21
22222 (Fig. 3, bold residues highlighted in yellow). PGBD5 acquired a new putative NLS that is centrally
23
24223 located in the sequences of Mm523 and Hs524 (Fig. 3, RKRKKRK motif typed in green and
25
26224 underlined). In agreement with the literature [14] we observed that the ectopic expression of murine
27
28225 PGBD5 isoform of 523 amino acids (Mm523) reduced the apparent efficiency of obtaining NeoR
29
30226 clones (Fig. 5a) that is close to that of PB.NLS-1-558 in HeLa cells (Fig. 2a). In integration assays
31
32227 done with the *Ifp2*-NeoR transposon donor plasmid under experimental conditions similar to those
33
34228 used above for PB.NLS-1-558, the rate of NeoR clones obtained with Mm523 (Fig. 5b) was similar
35
36229 to that obtained with the GFP control. In order to verify whether this was due to PGBD5
37
38230 cytotoxicity, we used a second cellular system developed in human rhabdoid tumor G401 cells and
39
40231 in which the endogenous expression of PGBD5 (Hs524) was found to have little impact on cell
41
42232 viability [14]. Under these experimental conditions, we found that the rate of NeoR clones was
43
44233 sevenfold higher than that of the GFP control (Fig. 5b). Together, this indicates that the expression
45
46234 rate of CRD-deficient *pble* transposases strongly impact the outcome of integration assays.

47
48235 The sequence features of integration events were studied through *Ifp2*-chromosome junctions
49
50236 obtained with Mm523 and Hs524. We prepared gDNA samples from ~1800 and 1600 NeoR clones
51
52237 obtained from integration assays done with *Ifp2*-NeoR and, Mm523 or Hs524, respectively. Using
53
54238 the Mm523 gDNA sample, we obtained 1461 transposon/chromosome junctions that were analyzed
55
56239 as described above (supplementary Table 3a). The profiles of transposon/chromosome junctions
57
58240 were found to be similar between integration events mediated by PB.NLS-1-558 and Mm523 (Fig.
59
60241 3c, and Table 2 versus last row in Table1). This was also verified by examining chromosomal sites
61
62242 where we found integration events in both orientations within the 1461 chromosomal sites used
63
64243 (supplementary Table 3a and b). When the junctions were categorized and analyzed in terms of
65
66244 percentages at each *Ifp2* end for PB, PB.NLS-1-558 and Mm523 we observed that: i) proper

245 junctions occurred more often at the 3' end than at the 5' end (Fig. 6, red bars), and ii) among
246 improper junctions, those without a canonical TSD and those located within TIR and juxtaposed
247 with transposon sequences (Fig. 6, blue bars and internal black bars, i.e. wounds at transposition
248 ends as exemplified in [22]) occurred more often than those located within the plasmid backbone
249 sequences juxtaposed near the TSD and TIR of the donor plasmid (Fig. 3, flanking black bars).
250 Using the Hs524 gDNA sample, we obtained 1051 transposon/chromosome junctions
251 (supplementary Table 4a). The junction profile was overall similar to those of both CRD-deficient
252 proteins (Fig. 3c and 6d), but it displayed a marked difference in that there were fourfold and twofold
253 less proper insertion events by transposition than in those obtained with PB.NLS-1-558 and Mm523,
254 respectively. Unexpectedly, this suggests that PGBD5 is more prone to trigger improper integration
255 in rhabdoid tumor G401 cells, consistent with the proposal that PGBD5 exhibits aberrant activities
256 in human rhabdoid tumors [14]. Since we observed that PB.NLS-1-558, Mm523 and Hs524 display
257 similar junction profile, we wondered if their insertion site preferences might be similar.

258 259 **Features of insertion sites targeted by PB variants and PGBD5 when integrating *Ifp2* into** 260 **chromosomes**

261 PB and PGBD5 have been shown to integrate *Ifp2* into intragenic regions more frequently than
262 expected by chance, specifically within transcription start site (TSS) regions flanking (± 5 kbp)
263 protein-coding genes [23-26]. Using our junction data, we observed that PB, PB.NLS-1-558,
264 Mm523, Hs524 did not distributed *Ifp2* integrations at random into intergenic and intragenic regions
265 (Fig. 7a; Chi2, $p = 2.08 \times 10^{-95}$, 2.21×10^{-10} , 0.0026, and 3.72×10^{-81} respectively), but with a significant
266 enrichment for intragenic regions (hypergeometric test, $p \ll 0.01$ for the four proteins). Similar
267 investigations were also done within regions flanking TSSs of five types of genes coding for: i)
268 proteins, ii) non-coding RNA (ncRNA), iii) micro RNA (miRNA), or being annotated in hg38 as
269 iv) pseudogenes or v) uncharacterized genes. We also found that *Ifp2* was integrated more frequently
270 than expected by chance within regions flanking TSSs in the 5 types of genes, except for the two
271 CRD-deficient proteins into uncharacterized genes (Fig. 7b; Chi2, $p = 2.58 \times 10^{-9}$, 2.16×10^{-8} , 8.27×10^{-12} ,
272 and 3.43×10^{-9} respectively), and with a significant enrichment in each type of genes
273 (hypergeometric test, $p \ll 0.01$ for the four proteins), except for the miRNA and ncRNA genes in
274 the PB.NLS-1-558 and Mm523 datasets, respectively (hypergeometric test, $p = 0.043$ and 0.051).
275 In addition to these global distribution features, a striking feature was that our Lumpy raw files,
276 after manual investigation using IGV, contained 166 (i. e. 177-11), 24 (44-20), 23 (26-3) and 11
277 (13-2) chromosomal sites, each displaying a fragment containing the *Ifp2* element inserted in both
278 orientations, i.e. inserted at least twice into these sites when integration events were mediated by
279 PB, PB.NLS-1-558, Mm523 and Hs524, respectively. We also found common insertions sites

280 among the PB, PB.NLS-1-558, Mm523 and Hs524 datasets (Table 3, lines 1, 3 and 5). These
281 insertion events occurred at the same nucleotide position site but this was not due to sample
282 contamination since they resulted from *Ifp2* integration events in different orientations and in some
283 cases from properly and improperly integrated *Ifp2* transposons. The number of these observations
284 was increased when using a 1000 bp window on both sides at each chromosomal insertion site
285 (Table 3, lines 2, 4 and 6; regions called below insertion sites-containing regions (ISCR)).
286 The choice of an insertion site by any *pble* transposase does not fully occur at random since a TTAA
287 motif is used. In the human genome model hg38, there are 18,713,270 TTAA motifs. Public data
288 about DNase I hypersensitivity mapping revealed that 98% of them are located in open chromatin
289 in HeLa cells (but also in HEK cells), i.e. accessible to DNA binding proteins such as transposases.
290 This means that the probability of integrating an *Ifp2* transposon twice into a single target site lies
291 about 1.8×10^{-7} in hg38. Taking into account the size of datasets used herein, to find several
292 insertions by chance into the same site is therefore unexpected under our experimental conditions.
293 Given the putative impact of some specific genomic features of HeLa cells such as their aneuploidy
294 [27-30], we further investigated ISCR features in other cell lines taking advantage of public datasets.
295 We used three of them that were produced from integration assays performed with *Ifp2* and PB in
296 HEK293 [23], in HCT116 [15] and in CD4+ [26] cells (21,967, 172,866 and 8954 chromosomal
297 sites, respectively (Table S4a, b, c)).
298 First, we confirmed that the rate of insertions mediated by PB into intragenic regions in each of the
299 four cell lines (HeLa, HEK293, HCT116 and CD4+; Table 4a, column 4) is 7 to 18% higher than
300 expected by chance (51.6%; Table 4b, column 1). This preference could not be explained by the
301 numbers of TTAA motifs in intragenic regions (53.1%; Table 3c, column 2), which is close to that
302 expected by chance. In spite of variations of aneuploidy and chromatin profiles between the four
303 cell lines, the insertion preference into intragenic regions does not appear to correlate to these
304 features.
305 In order to verify the statistical consistency of insertion sites shared between datasets, all pairs of
306 datasets were compared taking into account the variation of TTAA motif distribution between intra
307 and intergenic regions (Table 4c). P-values indicated that the number of commonly used
308 chromosomal insertion sites was significantly more elevated than expected by chance (Table 3; rows
309 1 to 16) whatever the window used around the insertion sites (0 or 1000 bp). We also observed that
310 18 ISCR were shared by the four datasets obtained with *Ifp2* and PB in the four cell lines.
311 We also examined insertion datasets obtained with the transposon *sleeping beauty* in HEK293 and
312 CD4+ cells (28490 and 8290 insertion sites, respectively [15,26]). Taking into account the
313 distribution of its TA targets in hg38 (Table 4c), results in Table 3 revealed that this transposon does
314 not display a significant preference between available putative TA target sites. It displayed higher

315 rates of insertion into intragenic regions (~62.5%, Table 4a) than predicted by chance (51.6%, Table
316 4b). In contrast to PB, this can however be correlated with TA density that is dramatically increased
317 in these regions (60.4%, Table 4c) compared to intergenic ones.

318 In all, our data reveal that PB, PB.1-558, Mm523 and Hs524 insert *Ifp2* preferentially into intragenic
319 regions with some level of site preference between available TTAA target sites.

320 321 **Features of genes targeted by *Ifp2* insertions mediated by PB variants and PGBD5**

322 We wondered whether insertion site preferences of *pble* transposases might also be seen at the level
323 of some intragenic regions. We postulated that experimental conditions of transposition assays are
324 conducive to forced integration of transposons into chromosomes. Therefore, we predicted that
325 insertions should be enriched among ISCR shared by several datasets than among those unique to
326 each dataset.

327 In the ontology analysis done with the 410 genes overlapped by ISCR and shared by at least two
328 datasets among those obtained with PB, PB.NLS-1-558, Mm523 and Hs524 in HeLa cells (Figure
329 6a), we found that 35/50 significant terms were directly related to the nervous system (Table 5).
330 This issue was therefore further investigated by verifying whether there was an enrichment of
331 “neuron genes” among ISCR overlapping with genes. For this purpose, we used the 6854 “neuron
332 genes” identified in hg19 based on the expression properties of 29165 genes (protein-, ncRNA- and
333 miRNA-coding plus some uncharacterized genes and pseudogenes) in 216 distinct human brain
334 structures [31]. We found that neuronally expressed genes were significantly enriched in each of the
335 PB, PB.NLS-1-558, Mm523 and Hs524 datasets obtained in HeLa cells (Table 4a, columns 5,6,7).
336 They were again enriched among the 697 genes overlapped by ISCR that were shared by at least
337 two of four datasets. This enrichment in neuron genes was 31.7-38.8% in each of the four datasets
338 and 37.2% (259/697) among the 697 shared genes. These results therefore support the notion that
339 neuron genes are preferred regions for *pble* transposases to insert *Ifp2*.

340 These observations were confirmed using datasets obtained in HEK293, HCT116 and CD4+ cells.
341 First, we found that neuronally expressed genes were significantly enriched in each of the four
342 datasets (Table 4a, columns 5,6,7). This did not appear to be related to the target density since the
343 percentage of ISCR in those genes was found to be 3 to 8% more elevated than the rate of TTAA
344 motifs in those genes (Table 4a, column 4 versus Table 4c, column 3). We found that genes
345 overlapped by ISCR and shared by at least two of four datasets did not display a significant
346 enrichment in neuron genes (3884/12618; 30.8%; Figure 6b). However, a very strong enrichment
347 was found when only ISCR shared by the four datasets were kept for the analysis (705/1100, 64%
348 neuron genes) and a strong depletion in neuron genes was found among genes occurring in only one
349 dataset (2273/12186; 18.7%).

350 Finally, we evaluated whether the insertion preferences into neuronally expressed genes were
351 specific to PB by analysing the same features in datasets obtained with two unrelated transposons
352 [15,26]. Data obtained with *sleeping beauty* in HCT116 and CD4 cells and with a *TcBuster* in
353 HCT116 cells indicated that both transposons also inserted more frequently into neuronally
354 expressed genes than in other genes (Table 4a, column 4,5,6). Their insertion preferences were also
355 about 6-8% above the density in their respective target motif (Table 4c, column 3). For *sleeping*
356 *beauty*, we also found that there was an enrichment in neuronally expressed genes among genes that
357 overlapped by ISCR and shared by both datasets (1143 neuron genes for 2813 genes; 40.63%) and
358 a depletion in those which were only found in one dataset (2625/9791; 26.8%).
359 Altogether, these last results reveal that *pble* transposases insert *Ifp2* more often than expected by
360 chance into neuronally expressed genes. However, this apparent preference is also displayed by
361 *sleeping beauty* and very likely by *TcBuster*, indicating that it is not specific of *pble* transposases.
362 This is not related to the gene size and the number of TTAA and TA because the densities in target
363 motifs are very close in neuron and non-neuron genes (TTAA: 5.318 ± 0.029 and 5.114 ± 0.015
364 motifs/kbp, respectively; TA: 44.68 ± 0.156 and 44.07 ± 0.083 motifs/kbp, respectively). Therefore,
365 this might result from the enhanced accessibility of neuronally expressed genes or their association
366 with DNA repair and chromatin remodeling factors that support DNA transposition, a property that
367 would be shared by multiple cell lines as exemplified here.

368 369 **Discussion**

370 This study generated two sets of novel insights. First, the two CRD-deficient transposases PB.NLS-
371 1-558 and PGBD5 (Mm523 and Hs524) mediate canonical *Ifp2* transposition, but also non-
372 canonical events that may result from events of improper transposition, transposase-dependent
373 integration by recombination and random integration. We assume that these CRD-deficient
374 transposases operate at a reduced efficiency than the "wild-type" or "full-length" *piggyBac*
375 transposase but their cytotoxicity on host cells and the possibility that they do integration by
376 transposase-dependent recombination indicated that they have nuclease activity, as previously
377 suggested for PGBD5 isoforms [13]. Furthermore, cytotoxicity issues often arise from the balance
378 between the level of expressed protein and the efficiency of mechanisms responsible for the
379 maintenance of genome integrity, which varies widely from one cell line to another. Such effects
380 could also be related to the cell cycle. This might explain why PB.NLS-1-558, which was always
381 located in the nucleus, had a more negative effect than PB.1-558 (Fig. 2a), which was mainly
382 cytoplasmic during most of the cell cycle and in contact with chromosomes only during the cell
383 division phase.

384 The second insight from this work concerns the ability of PB and both CRD-deficient *pble*

385 transposases to trigger integration events that did not seem to occur at random into chromosomes.
386 However, we cannot assume strict insertion site specificity of *pble* transposases since only a small
387 part of the observed insertions events in datasets are the same, down to the same nucleotide position.
388 However, our data demonstrated that these transposases displayed real preferences for insertion into
389 regions containing genes and frequently close to their TSS. In addition, our results supported that
390 *pble* transposases frequently targeted their *pble* insertions into genes committed to the central
391 nervous system function. This last point will need further experimental confirmation. Indeed, the
392 lengths of genes involved in nervous system function were, on average, longer than those of other
393 genes (for review see [32]). An alternative interpretation might be that the insertion sites we found
394 were preferentially located in neuronally expressed genes due to their size. However, our results
395 take into account the genome coverage and the density in TTAA motif and support that the observed
396 insertion preferences are not related these factors.

397

398 **Contribution to the understanding of *piggyBac* transposition.**

399 Previous studies suggested two roles for the CRD of PB. The CRD might be an essential component
400 of DNA-binding to the ends of *Ifp2*, mandatory for transposition [9]. This was confirmed in another
401 study, which also indicated that the CRD was essential for the assembly of the transposase dimer,
402 the active oligomer form for transposition [10]. Here we demonstrate that two CRD-deficient
403 *piggyBac* transposases were able to trigger proper *pble* transposition. Therefore, the CRD is not
404 essential for transposition but seems necessary for triggering proper transposition events, probably
405 by driving precise DNA cleavages at *pble* ends and directing a strict choice of TSD.

406

407 **Evolutionary reasons for domesticating a CRD-deficient *pble* transposase**

408 These results, as well as previously published data [14,20,21], have highlighted two properties of
409 PB, and perhaps of other *pble* transposases, that may have previously been underestimated in the
410 context of transposase-coding gene domestication. First, while PB mostly mediates proper *Ifp2*
411 transposition, it is also sometimes responsible for improper transposition that leads to neo-inserted
412 elements that are difficult or impossible to re-mobilise during new rounds of transposition. Second,
413 *pble* transposases might display strong preferences for insertion and genome rearrangements.
414 Indeed, we noted that PGBD5 is highly expressed in the mammalian nervous system [20] and that
415 current publicly available data (<https://www.gtexportal.org/home/gene/PGBD5>;
416 <https://www.proteinatlas.org/search/PGBD5>) widely support this conclusion. However, these data
417 concern mRNA expression and data regarding mRNA translation and protein expression will be
418 required. Nevertheless, if acquiring a mechanism for triggering irreversible DNA rearrangements in
419 the early steps of vertebrate evolution in the nervous system can be considered advantageous, then

420 PGBD5 seems to possess all the properties required to play such a role. The evolutionary history of
421 the RAG1/RAG2 proteins [33] suggests that each time a domesticated transposase has emerged
422 during evolution its domestication was concurrent with the domestication of its transposon targets.
423 In the PGBD5 context, verifying that *pbles* are domesticated and are used as binding targets by
424 PGBD5 for genome rearrangements will be challenging. Indeed, while PGBD5 is a highly
425 conserved protein in vertebrates, the *pble* landscape in these genomes varies drastically from one
426 host species to the next. In the human genome three *pbles* unrelated to PGBD5 are annotated:
427 MER75, MER85 and *Looper*. In the mouse genome only one *pble* closely related to human *Looper*
428 is annotated. In the zebrafish seven *pbles* that were not related to PGBD5 and to those present in
429 human and mouse genomes, have been annotated. In the chicken genome no *pble* has been found
430 so far [34,35]. It is possible that PGBD5 has been domesticated in order to mobilise multiple *pbles*
431 for recombination. Its protein sequence conservation in chicken and the absence of *pbles* in this
432 species suggests that it binds to other DNA binding targets, which may be related to the PGBD5-
433 specific signal (PSS) sequences observed in human rhabdoid tumors [14].

434

435

435 **Materials and methods**

436

436 **cDNA cloning of PGBD5 murine isoforms.**

437

437 A single mouse brain (strain C57Bl6) was used for total RNA extraction using Tri-reagent (Sigma-
438 Aldrich, St-Louis, MO, USA). cDNA synthesis was carried out using Omniscript RT kit and oligo
439 dT primers (Qiagen, Valencia, CA, USA). PCR primers with appropriate flanking restriction sites
440 were synthesized by Eurofins Genomics, Ebersberg, Germany. PCR was performed with Phusion
441 High-Fidelity PCR Master Mix (ThermoScientific). Following agarose gel electrophoresis, PCR
442 fragments were extracted (QIAquick gel extraction kit, Qiagen), submitted to enzyme restriction
443 (*EcoRI/XbaI* for the long N-term isoform and *EcoRI/XhoI* for the short N-term isoform), purified
444 (QIAquick PCR purification kit, Qiagen) and kept for cloning. Their sequence identity was verified
445 by Sanger sequencing (Eurofins Genomics, Ebersberg, Germany). The primers used to amplify
446 Mm523 (Accession N°: XM_006530804.1) isoforms are supplied in supplementary data 1a.

447

448

448 **Integration assay.**

449

449 **Plasmid expression for transposases.** The plasmids pCS2-PB and pCS2-PB.NLS-1-558 encode the
450 V5 tagged PB transposases. Each cDNA was inserted into the multi-cloning site of the pCS2+ vector
451 (Life Technologies, Paisley, UK) as described [36]. The plasmid pCS2-Mm523 encodes a two myc
452 tagged PGBD5 isoform 524 amino acid residues in size. Mm523 cDNA was inserted into the multi-
453 cloning site of a modified pCS2 vector with an in-frame N-term 5XMyC tag [37]. The plasmid pCS2-
454 GFP plasmid was built by cloning the gene coding for the green fluorescent protein gene into the

455

456

457

458

459

455 multi-cloning site of the pCS2+ vector. pCS2-GFP was used as a negative control of transposition
456 (i.e. absence of transposase expression).

1
2457
3
4458 **Plasmids donor of transposon.** The plasmid pBSK-IFP2-TIR5'-NeoR-TIR3' (supplementary data
5
6459 1b) was built by introducing the IFP2 5' and 3' terminal regions (262 and 400 bp, respectively) into
7
8460 the pBluescript SK plasmid (pBSK). A cassette (NeoR) containing a SV40 promoter, the neomycin
9
10461 phosphotransferase ORF and a sv40 terminator was cloned between transposon ends as described
11
12462 [36]. NeoR was cloned in its middle using a BamHI site that was added to its sequence during DNA
13463 synthesis. The plasmid pBSK-NeoR was built by cloning the NeoR cassette into the multi-cloning
14
15464 site of a pBSK plasmid as described [38].

16
17465
18
19466 **Integration assay in HeLa cells.** Assays were monitored as described [36]. Briefly, each sample of
20
21467 100000 cells in a well of a 24-well plates of plaque assays was co-transfected with JetPEI (Polyplus-
22
23468 transfection, Illkirch-Graffenstaden) and 400 ng DNA plasmid and with equal amounts of donor of
24469 NeoR cassette included or not within a transposon and transposase sources (1:1 ratio). Two days
25
26470 post-transfection, each cell sample was transferred to a cell culture dish (100 mm diameter) and
27
28471 selected with a culture medium containing 800 µg/mL G418 sulfate (Eurobio Scientific, Les Ulis)
29
30472 for 15 days. After two washing with 1X saline phosphate buffer, cell clones were fixed and stained
31
32473 overnight with 70% EtOH-0.5% methylene blue and colonies > 0.5 mm in diameter were counted.
33474 Experiments were performed at least twice in triplicate.

34
35475
36
37476 **Integration assay in G401 cells.** Assays were monitored as described [14]. Two clonal cell G401
38
39477 lines were used. The first line was lentivirally transduced to constitutively express specific shRNA
40
41478 suppressing the expression of Hs524 PGBD5 [14]. The second line was modified as a control to
42
43479 constitutively express shRNA to target GFP which is not expressed, thereby preserving the
44480 endogenous expression of Hs524 PGBD5. Briefly, each sample of 100000 cells in a well of a 24-
45
46481 well plates of plaque assays was transfected with jetOptimus and 500 ng DNA plasmid pBSK-IFP2-
47
48482 TIR5'-NeoR-TIR3' as recommended by the supplier (Polyplus- transfection, Illkirch-
49
50483 Graffenstaden). Two days post-transfection, each cell sample was transferred to a cell culture dish
51
52484 (100 mm diameter) and selected with a culture medium containing 2 mg/mL G418 sulfate (Eurobio
53
54485 Scientific, Les Ulis) for 15 days. After two washing with 1X saline phosphate buffer, cell clones
55486 were fixed and stained overnight with 70% EtOH-0.5% methylene blue and colonies > 0.5 mm in
56
57487 diameter were counted. Experiments were performed at least twice in triplicate.

58
59488
60
61489 **Cellular localization of green fluorescent protein-fusion proteins**

490 **Plasmid expression for transposase-GFP fusions.** The plasmids pCS2-PB-GFP, pCS2-PB.1-558,
491 pCS2-PB.NLS-1-558 and pCS2-Mm523-GFP were made as described [39].

2492
3
4493 **Cell manipulation.** HeLa cells were plated at a density of 5×10^4 cells per well in 1 cm² Lab-Tek™
5 chamber slides (Fisher Scientific, Waltham, MA, USA) and grown in DMEM (Gibco/Life
6494 Technologies, Paisley, UK) supplemented with 10 % heat inactivated fetal bovine serum (FBS,
7
8495 Eurobio, France) at 37 °C in a humidified atmosphere containing 5% CO₂ for 48 h. Cells were
9
10496 transfected with 500 ng plasmid DNA and jetPEI™ (Polyplus Transfection, Illkirch, France) at an
11
12497 N/P ratio of 5 in DMEM 10% FBS following the Manufacturer's instructions. Cells were then
13498 incubated with the complexes for 4 h. The transfection medium was then discarded and replaced by
14
15499 fresh DMEM supplemented with 10% FBS before being incubated for 48 hours at 37°C.

16
17500
18
19501
20
21502 **Imaging.** Cells on slides were fixed in 1X PBS/2% paraformaldehyde at RT for 15 min, and then
22503 permeabilised with PBS/1% (w/v) Triton-X100 for 10 min. The slides were washed three times for
23
24504 5 min with 1x PBS. Nuclei were stained using Vectashield Vibrance "Antifade Mounting Medium
25
26505 (hardening) + DAPI" (Vector Laboratories, Burlingame CA, USA). All images of fluorescence were
27
28506 collected with an LSM 700 laser scanning microscope and the associated Zen software (Carl Zeiss,
29
30507 Oberkochen, Germany). All images shown correspond to one focal plane (0.5 μm). Images to be
31
32508 used for figures were pseudocolored by LSM Image browser software (Carl Zeiss, Thornwood, NY)
33
34509 and Photoshop (Adobe Systems, San Jose, CA) was on the resulting tiff files only to adjust for
35510 brightness and contrast.

36
37511
38
39512 **Recovery of integration sites.**

40
41513 **LAM-PCR and Illumina libraries.** Integration assays were done to produce cell populations
42
43514 containing integrated copies of the donor transposon. Fifteen days post-transfection, cell clones
44515 were harvested for genomic DNA preparation using the DNeasy kit (Qiagen, Hilden, Germany).
45
46516 Linear amplification-mediated PCR (LAM-PCR) was performed to amplify the vector-genomic
47
48517 DNA junctions of *Ifp2* vectors as described [40]. All PCR were done using the high fidelity Q5
49
50518 DNA Polymerase (New England Biolabs, Ipswich, MA). For both approaches, 1 μg DNA was used
51
52519 for twice 50 rounds of linear amplification using a biotinylated primer anchored near one end of the
53
54520 NeoR cassette to enrich DNA species containing transposon-chromosomal DNA junctions (for
55
56521 sequences of (B)-NeoR 5' and 3' primers, see supplementary data 1c). One reaction was done per
57522 ends. The single-stranded products were immobilized on streptavidin-coated magnetic beads
58
59523 (Dynabeads M-280 Streptavidin, Invitrogen, Carlsbad, CA). All subsequent steps were performed
60
61524 on the magnetic bead-bound DNA. Two washes with water followed each step. Second strand

62
63 Helou *et al.*
64
65

525 synthesis was performed with random hexamer primers (Roche, Basel, Switzerland) using Klenow
526 DNA polymerase (New England Biolabs, Ipswich, MA). The double-stranded DNA was split in
527 two batches and subjected to restriction digests with *DpnI* for the first one and *PciI*, *NcoI* and *BspHI*
528 for the second one using restriction enzymes. The DNA fragments with a CG-3' or a CATG-3'
529 overhang ends were ligated to linkers displaying appropriate overhang ends and made from annealed
530 oligonucleotides (supplementary data 1c).

531 To increase the specificity of the full process, an initial PCR was done using one biotinylated primer
532 anchored within the 5' or 3' region of the transposon donor and one primer anchored within the
533 linker (for sequences of (B)-TIR-UTR 5' and 3', and LC1 primers, see supplementary data 1c). PCR
534 products were immobilized on streptavidin-coated magnetic beads and purified as described above.
535 Next, the bead-bound DNA was subjected to a nested PCR using nested primers anchored within
536 transposon ends and within linkers (supplementary data 1c). Final PCR products were purified,
537 quantified and gathered in equimolar DNA amounts for each transposon vector (4 populations of
538 LAM-PCR products) before being used to make Illumina libraries using NEBNext® Ultra™ II
539 DNA Library Prep Kit for Illumina® and NEBNext Multiplex Oligos for Illumina (New England
540 Biolabs, Ipswich, MA). Fragment size selection, library quality control and Illumina sequencing
541 (MiSeq 250 nucleotides, TruSeq SBS Kit v3) were achieved at the Plateforme de Séquençage Haut
542 Débit I2BC (Gif-sur-Yvette, France). DNA quantities were monitored at various steps in the
543 procedure with the Qubit® dsDNA (Molecular Probes, Eugene, USA).

544
545 **Computer analysis.** Trimmomatic [41] was used to filter Miseq reads using default parameters,
546 except for SLIDINGWINDOW:5:20 and MINLEN:100. The purpose of the following steps was to
547 recover chromosome-inserted DNA fragment junctions taking into account the plasmid backbone
548 regions located 100-bp upstream and downstream the *Ifp2*-NeoR transposon. Filtered reads were
549 first mapped to the sequence of plasmid backbone minus the 100-bp regions flanking on both sides
550 the *Ifp2*-NeoR transposon with bwa-mem using default parameters [42]. Unmapped reads were then
551 extracted reads using SAMtools view with parameters -b -f 4 [43] and bamToFastq from the
552 BEDTools suite using default parameters [44]. Recovered unmapped reads were aligned using bwa-
553 mem against a bwa bank gathering the sequences of hg38 chromosomes plus those of the *Ifp2*-NeoR
554 transposon flanked by the 100-bp plasmid backbone regions on both sides (supplementary data 1d).
555 Default parameters were used excepted for -w 1 and -r 1. The bam files resulting from each dataset
556 alignment were analysed with Lumpy in order to identify split reads [45]. The parameters were -e -
557 mw 2 -tt 0.0 and back_distance:20,weight:1,id:lumpy_v1,min_mapping_threshold:20. Structural
558 variants (SV) characterized by “BND” for the broken end notations and displaying for each of them
559 an SV with two positions, one genomic and one on the transposon, were extracted using a house

560 python program (<https://github.com/Leelouh/lumpy2site>). Results were filtered taking into account
561 a difference below 3 between the transposon breakpoint calculated by Lumpy and the maximal
562 spread of read alignments in the transposon donor sequence for each integration event. Each TSD
563 nucleotide motif at insertion site was obtained after extracting 10-bp sequences before and after the
564 breakpoint in the chromosome sequences.

565 Gene ontology (GO) analyses were focused mostly on protein coding genes and those encoding long
566 non-coding RNA (lncRNA). We used hg38 gene annotations from UCSC. Gene ontology was first
567 investigated using DAVID (<https://david.ncifcrf.gov/>) and AmiGO2
568 (<http://amigo.geneontology.org/amigo>) to assess term enrichment. This was followed up by the
569 Cytoscape plugin ClueGO [46,47].

570

571 **Access of publicly available data**

572 Sequences corresponding to *Ifp2*, *Sleeping Beauty* and *TcBuster* insertion sites in HEK 293 cells
573 [23] were downloaded from public databases using accession numbers JS717545 to JS799249.

574 Sequences corresponding to *Ifp2* insertion sites in HCT116 cells [15] were recovered at
575 <https://www.genetics.org/content/suppl/2012/01/03/genetics.111.137315.DC1>. Sequences

576 corresponding to *Ifp2* and *Sleeping Beauty* insertion sites in CD4+ cells [26] were recovered in the
577 GSE58744 at <https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSM1419000> and

578 <https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSM1419001>. For the last two sources, the
579 positions of insertion sites were transformed in hg38 using liftOver at [https://genome.ucsc.edu/cgi-](https://genome.ucsc.edu/cgi-bin/hgLiftOver)

580 [bin/hgLiftOver](https://genome.ucsc.edu/cgi-bin/hgLiftOver). All sites mapped in hg38 were supplied in supplementary Table 4.

581 PWMTrain [48] at <https://ccg.epfl.ch/pwmtools/pwmtrain.php> was used to calculate the position-
582 specific weight matrix of *TcBuster* insertion sites using available data [23]. The numbers and the

583 positions of putative insertion sites in hg38 for *pbles* (TTAA), *sleeping beauty* (TA) and *TcBuster*
584 were calculated using PWMScan [48] at <https://ccg.epfl.ch/pwmtools/pwmscan.php>.

585 The list of 6985 neuron genes in hg19 was recovered in supplemental data of [31] and was updated
586 to hg38. 131 genes were removed. They corresponded to artefactual genes coding nc RNA that were

587 withdrawn in hg 38. DNase1 map for HeLa and HEK293 cells were recovered at
588 <http://hgdownload.soe.ucsc.edu/wgEncodeAwgDnaseUwdukeHelas3UniPk.narrowPeak.gz> golden

589 [Path/hg19/encodeDCC/wgEncodeOpenChromDnase/](http://hgdownload.soe.ucsc.edu/goldenPath/hg19/encodeDCC/wgEncodeOpenChromDnase/), files and
590 <http://hgdownload.soe.ucsc.edu/goldenPath/hg19/encodeDCC/wgEncodeAwgDnaseUniform/wgE>

591 [ncodeAwgDnaseUwdukeHelas3UniPk.narrowPeak.gz](http://hgdownload.soe.ucsc.edu/goldenPath/hg19/encodeDCC/wgEncodeAwgDnaseUniform/wgEncodeAwgDnaseUwdukeHelas3UniPk.narrowPeak.gz). They were updated in hg38 using liftOver.

592 Values in graphs were medians, quartiles 1 and 3 and spread of experiments done at least twice in
593 triplicate. Shapiro-Wilk tests were used to confirm the normality of each set of samples, t-test to

594 analyse distribution differences between experimental samples, Chi2 test to analyse differences

595 between an experimental distribution and a theoretical one, and logarithmic distribution test to
596 analyse enrichments using free tools and tutorials available at <http://www.anastat.fr/outils.php>.
597 Permutation tests (10,000 per test) were computed using in-house bash programs that accounted the
598 distribution in TA and TTAA motifs in hg38. The normality of each distribution of permuted results
599 was verified using a Shapiro-Wilk test using free tools and tutorials available at
600 <http://www.anastats.fr/outils.php>. When the distributions were normal, probabilities were calculated
601 from Z score at <https://www.fourmilab.ch/rpkp/experiments/analysis/zCalc.html>. When they were
602 not normal, the distributions were used to determine the 1 and 0.1% thresholds at both tails and the
603 observed values were positioned in regards to those values.

604 605 **Data deposition.**

606 All raw and processed data are available through the European Nucleotide Archive under accession
607 number PRJEB36226, PRJEB36229, PRJEB41045 and PRJEB41053. Files describing the
608 annotation of insertion sites copies in the hg38 release are supplied as supplementary Tables 1, 2
609 and 3.

610

611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655

611 **Acknowledgments**

1 612 This work was supported by the C.N.R.S., the I.N.R.A., and the GDR CNRS 2157. It also received
2 613 funds from a research program grants from the Ligue Nationale Contre le Cancer, the Merck
3
4 614 foundation, and the French National Society of Gastroenterology. Laura Helou holds a PhD
5
6 615 fellowship from the Région Centre Val de Loire. We acknowledge the high-throughput sequencing
7
8 616 facility of I2BC for its sequencing and bioinformatics expertise. Alex Kentsis is a consultant for
9
10 617 Novartis and is supported by the National Cancer Institute grants R01 CA214812 and P30
11 618 CA008748. Yves Bigot, who was in charge of the achievement of this project does not have to thank
12
13 619 the French National Research Agency for its financial support but he kindly thanks it for the
14
15 620 excellent reviews embellished with arguments based on scientific and cultural novelties in the
16
17 621 expertise of his yearly application file during the last decade.

18
19 622

20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

623 **Appendix A: Supplemental data**

624 **Supplementary data 1. Transposon donor sequences and primers used in the study.**

625

626 **Supplementary Table 1. (a) Inventory of *Ifp2*-chromosome junctions resulting from**
627 **integration events mediated by PB.** Transposon breakpoints located at positions 102 and 2220,
628 103 and 2219, 104 and 2218, and 105 and 2217 had a TSD displaying a DTAA, DDAA, DDBA or
629 DDBB nucleotide motif (IUPAC nucleotide code; blue bars in Fig. 3a). Those at positions 101 and
630 2221, 100 and 2222, 99 and 2223, and 98 and 2224 (red bars in Fig. 3a) displayed a TTAA TSD
631 with upstream or downstream 0 nucleotide, one nucleotide, a dinucleotide or a trinucleotide that
632 were identical in the plasmid backbone and the chromosome. As the origin of this specific trait could
633 not be differentiated with a probability threshold below 1%, all junctions were considered in the
634 analysis as originating from proper integration events. **(b) Inventory of chromosomal site in which**
635 ***Ifp2* insertions were found several times.** Insertions corresponding to a single putative integration
636 event were highlighted in cyan blue and their duplicated TSD was highlighted in green. Integration
637 events corresponding to proper transposon ends were typed in black while those with improper ends
638 were typed in red.

639

640 **Supplementary Table 2. (a) Inventory of *Ifp2*-chromosome junctions resulting from**
641 **integration events mediated by PB.NLS-1-558.** Transposon breakpoints located at positions 102
642 and 2220, 103 and 2219, 104 and 2218, and 105 and 2217 had a TSD displaying a DTAA, DDAA,
643 DDBA or DDBB nucleotide motif (IUPAC nucleotide code; blue bars in Fig. 3b). Those at positions
644 101 and 2221, 100 and 2222, 99 and 2223, and 98 and 2224 (red bars in Fig. 3b) displayed a TTAA
645 TSD with upstream or downstream 0 nucleotide, one nucleotide, a dinucleotide or a trinucleotide
646 that were identical in the plasmid backbone and the chromosome. As the origin of this specific trait
647 could not be differentiated with a probability threshold below 1%, all junctions were considered in
648 the analysis as originating from proper integration events. **(b) Inventory of chromosomal site in**
649 **which *Ifp2* insertions were found several times.** Insertions corresponding to a single putative
650 integration event were highlighted in cyan blue and their duplicated TSD was highlighted in green.
651 Integration events corresponding to proper transposon ends were typed in black while those with
652 improper ends were typed in red.

653

654 **Supplementary Table 3. (a) Inventory of *Ifp2*-chromosome junctions resulting from**
655 **integration events mediated by Mm523.** Transposon breakpoints located at positions 102 and
656 2220, 103 and 2219, 104 and 2218, and 105 and 2217 had a TSD displaying a DTAA, DDAA,
657 DDBA or DDBB nucleotide motif (IUPAC nucleotide code; blue bars in Fig. 3c) displayed a TTAA

658 TSD with upstream or downstream 0 nucleotide, one nucleotide, a dinucleotide or a trinucleotide
659 that were identical in the plasmid backbone and the chromosome. As the origin of this specific trait
660 could not be differentiated with a probability threshold below 1%, all junctions were considered in
661 the analysis as originating from proper integration events. **(b) Inventory of chromosomal site in
662 which *Ifp2* insertions were found several times.** Insertions corresponding to a single putative
663 integration event were highlighted in cyan blue and their duplicated TSD was highlighted in green.
664 Integration events corresponding to proper transposon ends were typed in black while those with
665 improper ends were typed in red.

666
667 **Supplementary Table 4. (a) Inventory of *Ifp2*-chromosome junctions resulting from
668 integration events mediated by Hs524 in G401 cells.** Transposon breakpoints located at positions
669 102 and 2220, 103 and 2219, 104 and 2218, and 105 and 2217 had a TSD displaying a DTAA,
670 DDAA, DDBA or DDBB nucleotide motif (IUPAC nucleotide code; blue bars in Fig. 3c) displayed
671 a TTAA TSD with upstream or downstream 0 nucleotide, one nucleotide, a dinucleotide or a
672 trinucleotide that were identical in the plasmid backbone and the chromosome. As the origin of this
673 specific trait could not be differentiated with a probability threshold below 1%, all junctions were
674 considered in the analysis as originating from proper integration events. **(b) Inventory of
675 chromosomal site in which *Ifp2* insertions were found several times.** Insertions corresponding
676 to a single putative integration event were highlighted in cyan blue and their duplicated TSD was
677 highlighted in green. Integration events corresponding to proper transposon ends were typed in black
678 while those with improper ends were typed in red.

679
680 **Supplementary Table 5. Chromosomal positions mapped here in hg38 for *Ifp2* insertion sites
681 in HEK293 (a), HCT116 (b) and CD4+ (c) cells, sleeping beauty in HEK293 (d) and CD4+ (e),
682 and TcBuster in HEK293 (f).**

683

684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765

684 **References**

- 685 1. Piégu, B., Bire; S., Arensburger, P., Bigot, Y. (2016) A survey of transposable element
1
2686 classification systems--a call for a fundamental update to meet the challenge of their diversity
3
4687 and complexity. *Mol. Phylogenet. Evol.* 86, 90-109.
5
6688 2. Arensburger, P., Piégu, B., Bigot, Y. (2016) The future of transposable element annotation and
7
8689 their classification in the light of functional genomics - what we can learn from the fables of Jean
9
9690 de la Fontaine? *Mob. Genet. Elements.* 6, e1256852.
10
11691 3. Arkhipova, I.R. (2017) Using bioinformatic and phylogenetic approaches to classify transposable
12
13692 elements and understand their complex evolutionary histories. *Mob. DNA.* 8, 19.
14
15693 4. Goerner-Potvin, P., Bourque, G. (2018) Computational tools to unmask transposable elements.
16
17694 *Nat. Rev. Genet.* 19, 688-704.
18
19695 5. Bouallègue, M., Rouault, J.D., Hua-Van, A., Makni, M., Capy, P. (2017) Molecular Evolution of
20
21696 piggyBac Superfamily: From Selfishness to Domestication. *Gen. Biol. Evol.* 9, 323-339.
22697 6. Cary, L.C., Goebel, M., Corsaro, B.G., Wang, H.G., Rosen, E., Fraser, M.J. (1989) Transposon
23
24698 mutagenesis of baculoviruses: analysis of *Trichoplusia ni* transposon IFP2 insertions within the
25
26699 FP-locus of nuclear polyhedrosis viruses. *Virology.* 172, 156-169.
27
28700 7. Yusa, K. (2015) piggyBac Transposon. *Microbiol. Spectr.* 3, MDNA3-0028-2014.
29
30701 8. Chen, Q., Luo, W., Veach, R.A., Hickman, A.B., Wilson, M.H., Dyda, F. (2020) Structural basis
31
32702 of seamless excision and specific targeting by piggyBac transposase. *Nat Commun* 11, 3446.
33703 9. Morellet, N., Li, X., Wieninger, S.A., Taylor, J.L., Bischerour, J., Moriau, S., Lescop, E.,
34
35704 Bardiaux, B., Mathy, N., Assrir, N., Bétermier, M., Nilges, M., Hickman, A.B., Dyda, F., Craig,
36
37705 N.L., Guittet, E. (2018) Sequence-specific DNA binding activity of the cross-brace zinc finger
38
39706 motif of the piggyBac transposase. *Nucl. Acids. Res.* 46, 2660-2677.
40
41707 10. Sharma, R., Nirwal, S., Narayanan, N., Nair, D.T. (2018) Dimerization through the RING-
42
43708 Finger domain attenuates excision activity of the piggyBac transposase. *Biochemistry.* 57, 2913-
44709 2922.
45
46710 11. Keith, J.H., Fraser, T.S., Fraser, M.J.Jr. (2008) Analysis of the piggyBac transposase reveals a
47
48711 functional nuclear targeting signal in the 94 c-terminal residues. *BMC. Mol. Biol.* 9, 72.
49
50712 12. Hong, J.B., Chou, F.J., Ku, A.T., Fan, H.H., Lee, T.L., Huang, Y.H., Yang, T.L., Su, I.C., Yu,
51
52713 I.S., Lin, S.W., Chien, C.L., Ho, H.N., Chen, Y.T. (2014) A nucleolus-predominant piggyBac
53
54714 transposase, NP-mPB, mediates elevated transposition efficiency in mammalian Cells. *PLoS.*
55715 *One.* 9, e89396.
56
57716 13. Luo, W., Galvan, D.L., Woodard, L.E., Dorset, D., Levy, S., Wilson, M.H. (2017) Comparative
58
59717 analysis of chimeric ZFP-, TALE- and Cas9-piggyBac transposases for integration into a single
60
61718 locus in human cells. *Nucl. Acids. Res.* 45, 8411-8422.
62
63
64
65

- 719 14. Henssen, A.G., Koche, R., Zhuang, J., Jiang, E., Reed, C., Eisenberg, A., Still, E., MacArthur,
720 I.C., Rodríguez-Fos, E., Gonzalez, S., Puiggròs, M., Blackford, A.N., Mason, C.E., de Stanchina,
1 721 E., Gönen, M., Emde, A.K., Shah, M., Arora, K., Reeves, C., Socci, N.D., Perlman, E.,
2 722 Antonescu, C.R.; Roberts, C.W.M., Steen, H., Mullen, E., Jackson, S.P., Torrents, D., Weng, Z.,
3 4722 Armstrong, S.A., Kentsis, A. (2017) PGBD5 promotes site-specific oncogenic mutations in
5 6723 human tumors. *Nat. Genet.* 49, 1005-1014.
7 724
- 9 725 15. Wang, H., Mayhew, D., Chen, X., Johnston, M., Mitra, R.D. (2012) "Calling cards" for DNA-
10 726 binding proteins in mammalian cells. *Genetics.* 190, 941-949.
12
- 13 727 16. Li, M.A., Pettitt, S.J., Eckert, S., Ning, Z., Rice, S., Cadiñanos, J., Yusa, K., Conte, N., Bradley,
14 15728 A. (2013) The piggyBac transposon displays local and distant reintegration preferences and can
16 17729 cause mutations at noncanonical integration sites. *Mol. Cell. Biol.* 33, 1317-1330.
- 18 730 17. Thorvaldsdóttir, H., Robinson, J.T., Mesirov, J.P. (2013) Integrative Genomics Viewer (IGV):
19 731 high-performance genomics data visualization and exploration. *Brief. Bioinformatics.* 14, 178-
20 21 732 192.
23
- 24 733 18. Elick, T.A., Lobo, N., Fraser, M.J.Jr. (1997) Analysis of the cis-acting DNA elements required
25 26734 for piggyBac transposable element excision. *Mol. Gen. Genet.* 255, 605-610.
27
- 28 735 19. Mitra, R., Fain-Thornton, J., Craig, N.L. (2008) piggyBac can bypass DNA synthesis during cut
29 30 736 and paste transposition. *EMBO. J.* 27, 1097-1109.
- 31 737 20. Pavelitz, T., Gray, L.T., Padilla, S.L., Bailey, A.D., Weiner, A.M. (2013) PGBD5: a neural-
32 33738 specific intron containing piggyBac transposase domesticated over 500 million years ago and
34 35739 conserved from cephalochordates to humans. *Mob. DNA.* 4, 23-39.
36
- 37 740 21. Henssen, A.G., Henaff, E., Jiang, E., Eisenberg, A.R., Carson, J.R., Villasante, C.M., Ray, M.,
38 39741 Still, E., Burns, M., Gandara, J., Feschotte, C., Mason, C.E., Kentsis, A. (2015) Genomic DNA
40 41 742 transposition induced by human PGBD5. *Elife.* 4, e10565.
- 42 743 22. Lohe, A.R., Timmons, C., Beerman, I., Lozovskaya, E.R., Hartl, D.L. (2000) Self-inflicted
43 44744 wounds, template-directed gap repair and a recombination hotspot. Effects of the mariner
45 46745 transposase. *Genetics.* 154, 647-656.
47
- 48 746 23. Woodard, L.E., Li, X., Malani, N., Kaja, A., Hice, R.H., Atkinson, P.W., Bushman, F.D., Craig,
49 50 747 N.L., Wilson, M.H. (2012) Comparative analysis of the recently discovered hAT transposon
51 52 748 TcBuster in human cells. *PLoS. One.* 7, e42666.
- 53 749 24. Wilson, M.H., Coates, C.J., George, A.L.Jr. (2007) PiggyBac transposon-mediated gene transfer
54 55750 in human cells. *Mol. Ther.* 15, 139-145.
56
- 57 751 25. Huang, X., Guo, H., Tammana, S., Jung, Y.C., Mellgren, E., Bassi, P., Cao, Q., Tu, Z.J., Kim,
58 59 752 Y.C., Ekker, S.C., Wu, X., Wang, S.M., Zhou, X. (2010) Gene transfer efficiency and genome-
60 61 753 wide integration profiling of Sleeping Beauty, Tol2, and piggyBac transposons in human primary
62

- 754 T cells. *Mol. Ther.* 18, 1803-1813.
- 1 755 26. Gogol-Döring, A., Ammar, I., Gupta, S., Bunse, M., Miskey, C., Chen, W., Uckert, W., Schulz,
2 756 T.F., Izsvák, Z., Ivics, Z. (2016) Genome-wide profiling reveals remarkable parallels between
3 insertion site selection properties of the MLV retrovirus and the piggyBac transposon in primary
4 757 insertion site selection properties of the MLV retrovirus and the piggyBac transposon in primary
5 human CD4(+) T cells. *Mol Ther.* 24, 592-606.
- 6 758 27. Landry, J.J., Pyl, P.T., Rausch, T., Zichner, T., Tekkedil, M.M., Stütz, A.M., Jauch, A., Aiyar,
7 R.S., Pau, G., Delhomme, N., Gagneur, J., Korbel, J.O., Huber, W., Steinmetz, L.M. (2013) The
8 759 genomic and transcriptomic landscape of a HeLa cell line. *G3 (Bethesda)*. 3, 1213-1224.
- 9 760 28. Adey, A., Burton, J.N., Kitzman, J.O., Hiatt, J.B., Lewis, A.P., Martin, B.K., Qiu, R., Lee, C.,
10 761 Shendure, J. (2013) The haplotype-resolved genome and epigenome of the aneuploid HeLa
11 762 cancer cell line. *Nature*. 500, 207-211.
- 12 763 29. Lin, Y.C., Boone, M., Meuris, L., Lemmens, I., Van Roy, N., Soete, A., Reumers, J., Moisse,
13 764 M., Plaisance, S., Drmanac, R., Chen, J., Speleman, F., Lambrechts, D., Van de Peer, Y.,
14 765 Tavernier, J., Callewaert, N. (2014) Genome dynamics of the human embryonic kidney 293
15 766 lineage in response to cell biology manipulations. *Nat Commun.* 5, 4767.
- 16 767 30. Liu, Y., Mi, Y., Mueller, T., Kreibich, S., Williams, E.G., Van Drogen, A., Borel, C., Frank, M.,
17 768 Germain, P.L., Bludau, I., Mehnert, M., Seifert, M., Emmenlauer, M., Sorg, I., Bezrukov, F.,
18 769 Bena, F.S., Zhou, H., Dehio, C., Testa, G., Saez-Rodriguez, J., Antonarakis, S.E., Hardt, W.D.,
19 770 Aebersold, R. (2019) Multi-omic measurements of heterogeneity in HeLa cells across
20 771 laboratories. *Nat Biotechnol.* 37, 314-322.
- 21 772 31. Negi, S.K., Guda, C. (2017) Global gene expression profiling of healthy human brain and its
22 773 application in studying neurological disorders. *Sci. Rep.* 7, 897.
- 23 774 32. Zylka, M.J., Simon, J.M., Philpot, B.D. (2015) Gene length matters in neurons. *Neuron*. 86, 353-
24 775 355.
- 25 776 33. Zhang, Y., Cheng, T.C., Huang, G., Lu, Q., Surleac, M.D., Mandell, J.D., Pontarotti, P.,
26 777 Petrescu, A.J., Xu, A., Xiong, Y., Schatz, D.G. (2019) Transposon molecular domestication and
27 778 the evolution of the RAG recombinase. *Nature*. 569, 79-84.
- 28 779 34. Guizard, S., Piégu, B., Arensburger, P., Guillou, F., Bigot, Y. (2016) Deep landscape update of
29 780 dispersed and tandem repeats in the genome model of the red jungle fowl, *Gallus gallus*, using a
30 781 series of de novo investigating tools. *BMC Genomics*. 17, 659.
- 31 782 35. Kapusta, A., Suh, A. (2017) Evolution of bird genomes-a transposon's-eye view. *Ann. N. Y.*
32 783 *Acad. Sci.* 1389, 164-185.
- 33 784 36. Bire, S., Ley, D., Casteret, S., Mermoud, N., Bigot, Y., Rouleux-Bonnin, F. (2013) Optimization
34 785 of the piggyBac transposon using mRNA and insulators: toward a more reliable gene delivery
35 786 system. *PLoS. One*. 8, e82559.
- 36 787
37 788

789 37. Travnickova-Bendova, Z., Cermakian, N., Reppert, S.M., Sassone-Corsi, P. (2002) Bimodal
790 regulation of mPeriod promoters by CREB-dependent signaling and CLOCK/BMAL1 activity.
1 2791 Proc. Natl. Acad. Sci. USA. 99, 7728-7733.
3
4792 38. Bire, S., Casteret, S., Piégu, B., Beauclair, L., Moiré, N., Arensbuger, P., Bigot, Y. (2016)
5 6793 Mariner Transposons Contain a Silencer: Possible Role of the Polycomb Repressive Complex 2.
7 8794 PLoS. Genet. 12, e1005902.
9
9795 39. Demattei, M.V., Hedhili, S., Sinzelle, L., Bressac, C., Casteret, S., Moiré, N., Cambefort, J.,
10 11796 Thomas, X., Pollet, N., Gantet, P., Bigot, Y. (2011) Nuclear importation of Mariner transposases
12 13797 among eukaryotes: motif requirements and homo-protein interactions. PLoS One. 6, e23693.
14
15798 40. Bartholomae, C.C., Glimm, H., von Kalle, C., Schmidt, M. (2012) Insertion site pattern: global
16 17799 approach by linear amplification-mediated PCR and mass sequencing. Meth. Mol. Biol. 859,
18 19800 255-265.
20
21801 41. Bolger, A.M., Lohse, M., Usadel, B. (2014) Trimmomatic: a flexible trimmer for Illumina
22802 sequence data. Bioinformatics. 30, 2114-2120.
23
24803 42. Li, H., Durbin, R. (2010) Fast and accurate long-read alignment with Burrows–Wheeler
25 26804 transform, Bioinformatics. 26, 589-595.
27
28805 43. Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., Marth, G., Abecasis, G.,
29 30806 Durbin, R., and 1000 Genome Project Data Processing Subgroup. (2009) The Sequence
31807 alignment/map (SAM) format and SAMtools. Bioinformatics. 25, 2078-2079.
32
33808 44. Quinlan, A.R., Hall, I.M. (2010) BEDTools: a flexible suite of utilities for comparing genomic
34 35809 features. Bioinformatics. 26, 841-842.
36
37810 45. Layer, R.M., Chiang, C., Quinlan, A.R., Hall, I.M. (2014) LUMPY: a probabilistic framework
38 39811 for structural variant discovery. Gen. Biol 15, R84.
40
41812 46. Bindea, G., Mlecnik, B., Hackl, H., Charoentong, P., Tosolini, M., Kirilovsky, A., Fridman,
42813 W.F., Pagès, F., Trajanoski, Z., Galon, J. (2009) ClueGO: a Cytoscape plug-in to decipher
43 44814 functionally grouped gene ontology and pathway annotation networks. Bioinformatics. 25, 1091-
45 46815 1093.
47
48816 47. Mlecnik, B., Galon, J., Bindea, G. (2018) Comprehensive functional analysis of large lists of
49 50817 genes and proteins. J. Proteomics. 171, 2-10.
51
52818 48. Ambrosini, G., Groux, R., Bucher, P. (2018) PWMScan: A Fast Tool for Scanning Entire
53819 Genomes with a Position-Specific Weight Matrix. Bioinformatics. 34, 2483-2484.
54
55
56820
57
58
59
60
61
62
63
64
65

821 **Legends**

822 **Fig. 1. Cellular localization of GFP fusions in HeLa cells transiently transfected with a vector**
823 **expressing GFP (a, b, c), PB-GFP (d, e, f), PB.1-558-GFP (g, h, i) and PB.NLS-1-558-GFP (j,**
824 **k, l).** The left panels (a, d, g, j) show GFP fluorescence, the middle panels (b, e, h, k)) show the
825 nuclear genomic DNA staining by Hoechst 33342, the right panels (c, f, I, l) correspond to merge
826 pictures.

827
828 **Fig. 2. Box plot representations of integration assay results.** (a) impact of the two PB variants
829 (PB.1-558 and PB.NLS-1-558), PB and GFP on the rate of random integration of a NeoR cassette.
830 (b) rates of NeoR clones resulting from the integration of *Ifp2*-NeoR when recombination was
831 mediated by the two PB variants, PB and GFP. (c) rates of NeoR clones resulting from the
832 integration of *Ifp2*-NeoR when recombination was mediated by PB.NLS-1-558, PB and GFP and
833 corrected by the rate of toxicity of each protein calculated in (a). In (b) and (c), integration rates
834 were expressed in rate NeoR clones that were normalized using controls done with GFP. In each
835 plot, the red lines represented the median and the standard deviation, respectively.

836
837 **Fig. 3. Number and location of transposon breakpoints in *pble* sequences after transposition**
838 **into chromosomes.** Histogram distributions of *Ifp2*-NeoR extremities transposed by PB (a),
839 PB.NLS-1-558 (b), Mm523 (c) and Hs524 (d) (detailed in supplementary Tables 1 to 4). Red bars
840 indicated insertion events with perfectly conserved TSD and TIR while blue bars located those in
841 which TIR were perfectly conserved but the TSD did not correspond to a canonical TTAA at the
842 outermost extremities of pbles. Black bars represented breakpoints within the transposon sequence
843 and within plasmid backbone sequences juxtaposed to the transposon. Each bar corresponded to the
844 number of junctions found at a single nucleotide position. Green boxes located the position of
845 primers anchored within the transposon sequence and used at the last step of LAM-PCR. These
846 graphics described the relative importance of wounds at transposon ends under our experimental
847 conditions. However, they could not allow calculating wound rates at each of both ends due to the
848 fact that the final LAM-PCR products in each dataset came from the gathering of several LAM-
849 PCR reactions.

850
851 **Fig. 4. Sequence features of the *Ifp2* transposase (PB) variants and two Mm523-like PGBD5**
852 **isoforms.** (a) Protein sequence alignment of *Ifp2* transposase (PB) with two murine and human
853 domesticated PGBD5 proteins corresponding to the orthologous Hs524 and Mm523 isoforms. (b)
854 Sequence features of PB.NLS-1-558. Secondary structure predictions calculated with psipred
855 (<http://bioinf.cs.ucl.ac.uk/psipred/>) and Jpred4 (<http://www.compbio.dundee.ac.uk/jpred/>) were

856 highlighted in pink for α -helices and in orange for β -strands. The three proteins share two domains:
1 857 a N-terminal domain that was few structured, with an acid pI (boxed regions) and repeated acid
2 858 motifs (in red letters), a domain of ~400 amino acid residues that display a basic pI. PB contained a
3
4 859 third C-terminal domain, the CRD, that contains cysteins (highlighted in green) able to assemble
5
6 860 zinc finger folds. Aspartic residues inactivating the recombinase catalytic activity were bolded and
7
8 861 highlighted in yellow [17,21]. The PB NLS and the putative NLS in PGBD5 isoforms were
9
10 862 underlined and typed in green.

11 863
12

13 864 **Fig. 5. Graphic representations of integration assay results.** (a) impact of Mm523 on the rate of
14
15 865 random integration of a NeoR cassette. (b) rates of integration of an *Ifp2*-NeoR when recombination
16
17 866 was mediated by PB and Mm523 in HeLa cells, and Hs524 in G401 cells. Integration rates were
18
19 867 expressed in rate NeoR clones that were normalized using controls done with GFP (green). In each
20
21 868 plot, the red lines represented the median and the standard deviation, respectively.

22 869
23

24 870 **Fig. 6. Features of transposon breakpoints in *Ifp2* transposed by PB (a), PB.NLS-1-558 (b),**
25
26 871 **Mm523 (c) and Hs524 (d).** Black bars corresponded to percentages of breakpoints located within
27
28 872 the plasmid backbone flanking the transposon or those located within inner transposon regions (from
29
30 873 the position 2 in TIR to the primer used for the LAM-PCR). Red bars corresponded to those within
31
32 874 transposons that displayed intact canonical TTAA TSD and TIRs. Blue bars corresponded to *pbles*
33 875 displaying noncanonical TSD but intact TIRs.

34 876
35 877
36

37 877 **Fig. 7. Proportions of *Ifp2* insertions mediated by PB, PB.NLS-1-558, Mm523 and Hs524 in**
38
39 878 **intragenic regions (a) and regions containing TSS (b) taking into account the five gene**
40
41 879 **categories: protein-coding genes, ncRNA-coding genes, miRNA-coding genes, pseudogenes and**
42
43 880 **uncharacterized genes. Black bars indicated the expect percentage in a random distribution.**

44 881
45

46 882 **Fig. 8. Venn diagram representations of intragenic regions overlapped by insertion site-**
47
48 883 **containing regions (\pm 1000 bp) between transposition assays done with PB, PB.NLS-1-558,**
49
50 884 **Mm523 and Hs524 (a) and with PB in HeLa, HEK293, HCT116 and CD4+ cells (b). The**
51
52 885 **numbers of intragenic regions specific of datasets were italicized and typed in grey. Those shared**
53 886 **by all datasets were typed in orange.**

54
55
56
57
58
59
60
61
62
63
64
65

Table 1. Presence of canonic TTAA TSD and-or TIR among transposon/chromosome junctions resulting from LAM-PCR products and originating from events mediated by PB variants.

Transposase source	Transposon source	Proper TSD and TIR	Proper TSD improper TIR	Improper TSD/proper TIR	Improper TSD and TIR
PB	<i>Ifp2</i> -NeoR	96.7 % (7370)	1.05 % (78)	1.05 % (80)	1.2 % (95)
PB-NLS-1-558	<i>Ifp2</i> -NeoR	19.0 % (98)	2.5 % (13)	3.1 % (16)	75.4 % (388)

Table 2. Presence of canonic TTAA TSD and-or TIR among transposon/chromosome junctions resulting from LAM-PCR products and originating from events mediated by Mm523 and Hs524.

Transposase source	Transposon source	Proper TSD and TIR	Proper TSD improper TIR	Improper TSD/proper TIR	Improper TSD and TIR
Mm523 PGBD5	<i>Ifp2</i> -NeoR	10.0 % (147)	5.4 % (80)	4.3 % (64)	80.3 % (1188)
Hs524 PGBD5	<i>Ifp2</i> -NeoR	4.6 % (48)	2.1 % (22)	1.3% (14)	92.0 % (967)

Table 3. Number of chromosomal sites in common in each dataset pair and their statistical consistency in permutation tests.

Feature of dataset 1			Feature of dataset 2			Window around insertion sites (\pm) in both datasets	Number of sites observed	Random permutations features			Probability = H0 was no differences between obs. and exp.
Cell line 1	Transposon 1	Transposase 1	Cell line 2	Transposon 2	Transposase 2			Average number of sites expected per chance	Standard deviation	Z score	
HeLa	<i>Ifp2</i>	PB	HeLa	<i>Ifp2</i>	PB.NLS-1-558	0	1	N.A.	N.A.	N.A.	p<0.001
HeLa	<i>Ifp2</i>	PB	HeLa	<i>Ifp2</i>	PB.NLS-1-558	1000	37	N.A.	N.A.	N.A.	p<0.001
HeLa	<i>Ifp2</i>	PB	HeLa	<i>Ifp2</i>	Mm523-PGBD5	0	83	N.A.	N.A.	N.A.	p<0.001
HeLa	<i>Ifp2</i>	PB	HeLa	<i>Ifp2</i>	Mm523-PGBD5	1000	175	N.A.	N.A.	N.A.	p<0.001
HeLa	<i>Ifp2</i>	PB	G401	<i>Ifp2</i>	Hs524-PGBD5	0	24	N.A.	N.A.	N.A.	p<0.001
HeLa	<i>Ifp2</i>	PB	G401	<i>Ifp2</i>	Hs524-PGBD5	1000	42	N.A.	N.A.	N.A.	p<0.001
HeLa	<i>Ifp2</i>	PB	HEK293	<i>Ifp2</i>	PB	0	21	N.A.	N.A.	N.A.	p<0.001
HeLa	<i>Ifp2</i>	PB	HEK293	<i>Ifp2</i>	PB	1000	471	260.162	17.0215	12.3865	0
HeLa	<i>Ifp2</i>	PB	HCT116	<i>Ifp2</i>	PB	0	238	67.01	8.4329	20.2763	0
HeLa	<i>Ifp2</i>	PB	HCT116	<i>Ifp2</i>	PB	1000	2387	N.A.	N.A.	N.A.	p<0.006
HeLa	<i>Ifp2</i>	PB	CD4+	<i>Ifp2</i>	PB	0	24	N.A.	N.A.	N.A.	p<0.001
HeLa	<i>Ifp2</i>	PB	CD4-	<i>Ifp2</i>	PB	1000	213	97.684	10.2987	11,197	0
HEK293	<i>Ifp2</i>	PB	HCT116	<i>Ifp2</i>	PB	0	391	N.A.	N.A.	N.A.	p<0.001
HEK293	<i>Ifp2</i>	PB	HCT116	<i>Ifp2</i>	PB	1000	6600	5419.484	93.3269	12.0060	0
HEK293	<i>Ifp2</i>	PB	CD4+	<i>Ifp2</i>	PB	0	362	N.A.	N.A.	N.A.	p<0.001
HEK293	<i>Ifp2</i>	PB	CD4+	<i>Ifp2</i>	PB	1000	752	N.A.	N.A.	N.A.	p<0.001
HCT116	<i>Ifp2</i>	PB	CD4+	<i>Ifp2</i>	PB	0	16	81.825	9.0919	17.5073	0
HCT116	<i>Ifp2</i>	PB	CD4+	<i>Ifp2</i>	PB	1000	5235	2213.501	50.9891	59.2576	0
HEK293	<i>Sleeping Beauty</i>	SB	CD4+	<i>Sleeping Beauty</i>	SB	0	3	N.A.	N.A.	N.A.	p>0.083
HEK293	<i>Sleeping Beauty</i>	SB	CD4+	<i>Sleeping Beauty</i>	SB	1000	385	349.967	19.09	1.8350	p=0.034

N.A., not appropriated to use a Z-test as the distribution of the 1000 permutations results did not fulfil normality in a Shapiro-Wilk test. In the most right column is typed in red probabilities supporting that there was no difference at a threshold of 0.01.

Table 4a. Rate of regions containing insertion sites in genes and neuron genes in hg38

1, Dataset features: cells, transposon / transposase sources	2, No of different ISCR*	3, No & % of different ISCR overlapping a gene***	4, No of genes overlapped by ISCR**	5, No & % of different ISCR overlapping a neuron gene***	6, No & % of neuron genes among genes overlapped by ISCR**	7, p****
HeLa, <i>Ifp2</i> / PB	7,623	5,060 (66.4%)	4,663	2,231 (29.2%)	1,587 (34.0%)	2.1×10^{-82}
HeLa, <i>Ifp2</i> /PB.NLS-1-558	516	327 (63.2%)	376	123 (23.8%)	123 (37.6%)	2.5×10^{-6}
HeLa, <i>Ifp2</i> / Mm523	1,479	863 (58.4%)	956	321 (21.7%)	303 (31.7 %)	3.1×10^{-11}
G401, <i>Ifp2</i> / Hs524	1052	665 (63.2%)	740	279 (26.5%)	258 (38.8%)	1.3×10^{-14}
HEK293, <i>Ifp2</i> /PB [22]	21,967	15,226 (69.3%)	9,370	7,173 (32.6%)	3,032 (32.4%)	1.9×10^{-150}
HCT116, <i>Ifp2</i> /PB [14]	172,866	113,648 (67.7%)	23,578	49,688 (28.8%)	5,945 (25.3%)	1.3×10^{-88}
CD4+, <i>Ifp2</i> / PB [25]	8,954	6,940 (77.5%)	5,763	3,134 (35.0%)	1,903 (33.0%)	6.2×10^{-89}
HCT116, <i>sleeping beauty</i> /SB [14]	28,490	17,534 (61.54%)	10,460	7,776 (27.3%)	3,229 (30.9%)	2.5×10^{-128}
CD4+, <i>sleeping beauty</i> /SB [25]	8,290	5,441 (64.63%)	5,133	2,406 (29.0%)	1,750 (34.1%)	2.9×10^{-93}
HCT116, <i>TcBuster</i> /TCBUSTER [14]	17,227	11,841 (68.7%)	8,522	5,517 (32.0%)	2,820 (33.1%)	4.9×10^{-151}

Table 4b. Numbers and sequence coverages (%) of genes and neuron genes in hg38

1, Genes	2, Neuron genes
30077 (51.6%)	6854 (22.8%)

Table 4c. Numbers and rate of potential insertion sites for PB, SB and TcBuster in hg38

	1, Number of potential insertion sites in hg38	2, Number & % of potential insertion sites in genes in hg38	3, Number & % of potential insertion sites in neuron genes in hg38
<i>Ifp2</i> (TTAA)	18,713,270	9,943,117 (53.1%)	4,521,501 (24.2%)
<i>sleeping beauty</i> (TA)	152,412,514	92,065,399 (60.4%)	36,803,438 (24.1%)
<i>TcBuster</i> (NNNTANNN)	130,938	74,733 (57.1%)	33,940 (25.9%)

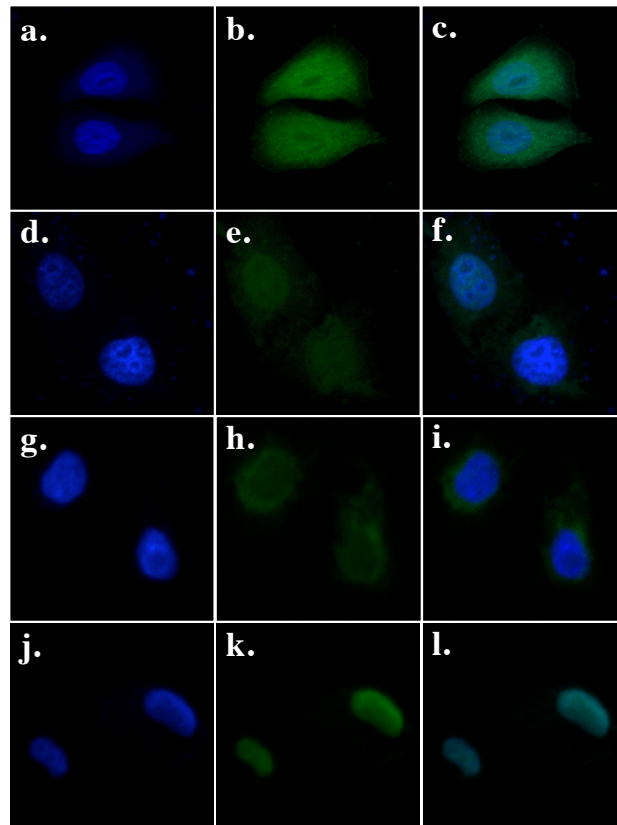
*, ISCR = insertion sites-containing regions, i.e. plus 1000 bp upstream and downstream each insertion site ; **, 2 genes can overlap an ISCR; ***, each genes and neurons genes = exons and introns of their transcriptional unit plus 5000 bp upstream and downstream; ****, hypergeometric test using H_0 = no enrichment in neuron genes among genes overlapped by an ISCR.

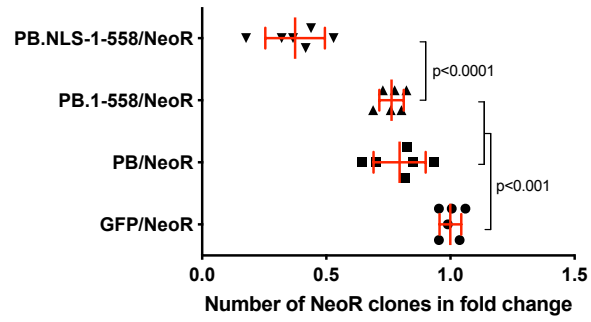
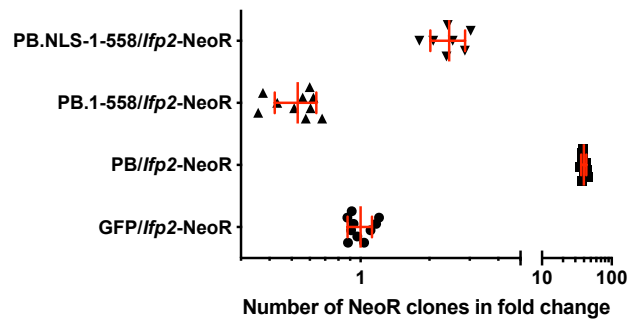
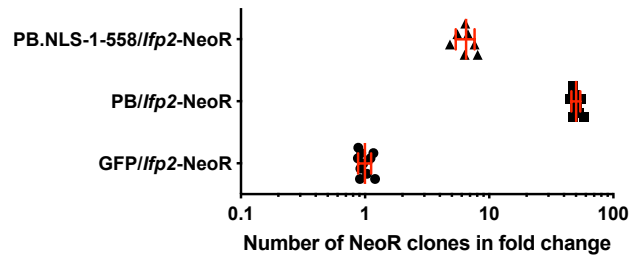
Table 5. Significant ontology terms resulting from the analysis of the 817 genes shared by PB, PB.NLS-1-558, Mm523 and Hs524.

GO ID	GO Term	P Value*	Nr. Genes
Fused GO terms: kinase activity (10.62% genes)			
GO:0016301	kinase activity	0,004222	80
Fused GO terms: establishment of cell polarity (2.26 % genes)			
GO:0030010	establishment of cell polarity	0,004364	17
Fused GO terms: cranial nerve morphogenesis (2.92% genes)			
GO:0021602	cranial nerve morphogenesis	0,000393	9
GO:0021953	central nervous system neuron differentiation	0,030239	19
Fused GO terms: regulation of signal transduction (25.76% genes)			
GO:0023051	regulation of signaling	0,008051	161
GO:0048583	regulation of response to stimulus	0,013282	183
GO:0010646	regulation of cell communication	0,005835	160
GO:0009966	regulation of signal transduction	0,001688	145
Fused GO terms: synapse organization (6,51% genes)			
GO:0034330	cell junction organization	0,001312	49
GO:0034329	cell junction assembly	0,023703	32
GO:0050808	synapse organization	0,000431	36
GO:0099173	postsynapse organization	0,005316	19
GO:0099084	postsynaptic specialization organization	0,013763	8
Fused GO terms: regulation of small GTPase mediated signal transduction (14.61% genes)			
GO:0044093	positive regulation of molecular function	0,018196	92
GO:0007264	small GTPase mediated signal transduction	0,004348	40
GO:0008047	enzyme activator activity	0,013464	37
GO:0043087	regulation of GTPase activity	0,000939	37
GO:0060589	nucleoside-triphosphatase regulator activity	0,000800	30
GO:0030695	GTPase regulator activity	0,000726	28
GO:0043547	positive regulation of GTPase activity	0,022595	30
GO:0051056	regulation of small GTPase mediated signal transduction	0,000102	31
GO:0005096	GTPase activator activity	0,000898	26
Fused GO terms:nervous system development (37.32% genes)			
GO:0007275	multicellular organism development	0,000027	238
GO:0009653	anatomical structure morphogenesis	0,000016	139
GO:0060322	head development	0,002424	53
GO:0000902	cell morphogenesis	0,000148	67
GO:0032989	cellular component morphogenesis	0,001750	52
GO:0048468	cell development	0,006713	107
GO:0048731	system development	0,000513	212
GO:0051128	regulation of cellular component organization	0,020339	114
GO:2000026	regulation of multicellular organismal development	0,048972	100
GO:0048513	animal organ development	0,035004	158
GO:0000904	cell morphogenesis involved in differentiation	0,004732	49
GO:0007399	nervous system development	0,000000	142
GO:0007417	central nervous system development	0,000026	69
GO:0022008	neurogenesis	0,000002	99
GO:0051960	regulation of nervous system development	0,001844	59
GO:0120036	plasma membrane bounded cell projection organization	0,000084	88
GO:0007420	brain development	0,005967	50
GO:0120035	regulation of plasma membrane bounded cell projection organization	0,023390	44
GO:0048699	generation of neurons	0,000001	96
GO:0120039	plasma membrane bounded cell projection morphogenesis	0,000663	48
GO:0030182	neuron differentiation	0,000000	90
GO:0048666	neuron development	0,000141	71

GO:0045664	regulation of neuron differentiation	0,047755	42
GO:0031175	neuron projection development	0,000310	64
GO:0048667	cell morphogenesis involved in neuron differentiation	0,001687	43
GO:0016358	dendrite development	0,043162	22
GO:0061564	axon development	0,013767	37
GO:0007409	axonogenesis	0,010504	35

*Term PValue corrected with Bonferroni step down. GO terms related to neurogenesis and neuron



a. Transposase toxicity**b. Observed integration rates****c. Toxicity-corrected integration rates**

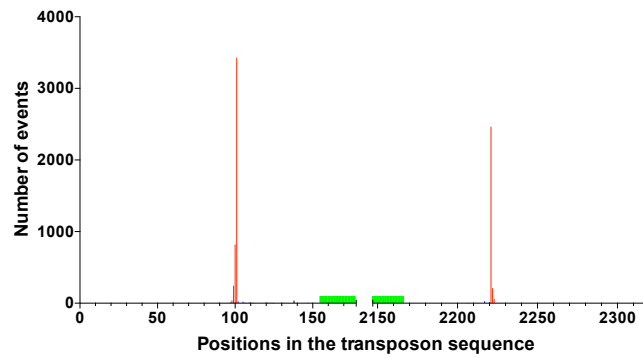
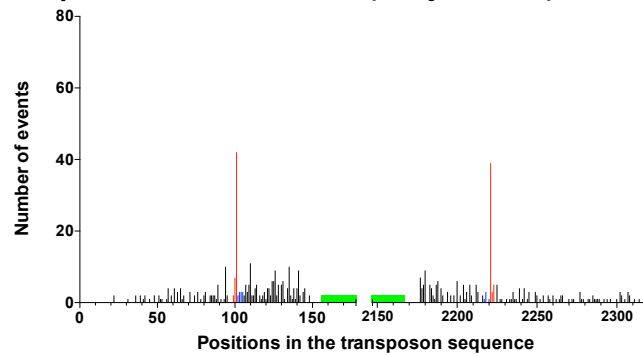
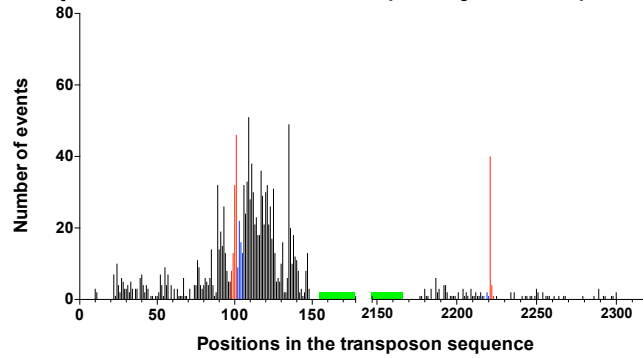
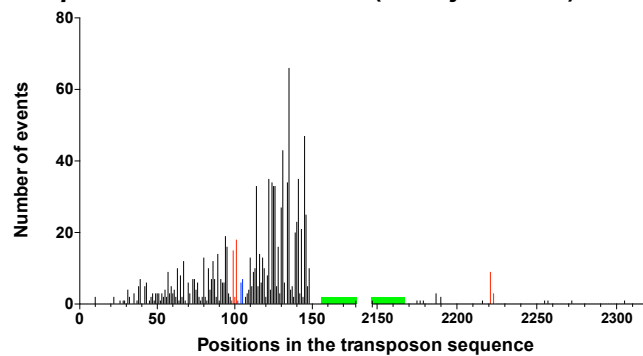
a. *lfp2*-NeoR/PB (7623 junctions)**b. *lfp2*-NeoR/PB.NLS-1-558 (516 junctions)****c. *lfp2*-NeoR/Mm523 PGBD5 (1461 junctions)****d. *lfp2*-NeoR/Hs524 PGBD5 (1051 junctions)**

Figure 4

a.

```

1 | pI= 4,46 | MGSSLDDEHILSALLOSDDELVGEDSDSEISDHVSEDDVQSDTEEAFIDEVHEVQPTSSGSEILDEQNVIQPGSSLASNRILT
2 | pI= 4,66 | MAEGGGGARRRAPALLEAARARYESLHI SDDVFGESGPDSSGNPFYSTSAASRSSSAAS SDDEREP-----PGPPGAAPP
3 | pI= 4,66 | MAEGGGGRRRAPALLEAARARYESLHI SDDVFGESGPDSSGNPFYSTSAASRSSSAAS SDDERER-----PAPPGTAPP

1 | LPQRTIRGKNKHCWSTSKSTRSRVSALNIVRSORGPTR-MCRNIYDPLLCFKLFTDEIISEIVKWTNAEISLKRRESMT-GATFRDTNEDEIYAFFGILVMTAVRKDN
2 | PPRAPDAQEP EDEEAGAGWSAALRDRPPRFEDTGGPTRKMP-SASAVDFQLFVPDNLVKNMVVQTNMYAKKQERFGSDGAWVEVTLTEMKAFGLYMISTSISHCES
3 | S-YAADPLELE EDEETGGGWSAVLRDRPSPRFEDTGGPTRKMP-SASAVDFQLFVPDNLVKNMVVQTNMYARKQERFGSDGAWVEVTLAEMKAFGLYVISTSVSHCES

1 | HMSTDDLFRDRLSMV-YVSVMSRDRDFLIRCLRMDDKSIRPTLRENDVFTPVKRIWDLFIHQCIQNYTPG----AHLTI DEQLLGFGRGCPFRMYIPNKPS KYGIKIL
2 | VLSIWSGGFYSN-RSLAL-VMSQARFEKILKYFHVVAFRSSQTTHG---LYKVQPFLDSLQNSFDSA FRPSQTQVLHEPLIDE DPVFIATCTERELRKRKRK FSLWVRQ
3 | VLSIWSGGFYSN-RSLAL-VMSQARFEKILKYFHVVAFRSSQTTHG---LYKVQPFLDSLQNSGFDAAFRPSQTQVLHEPLIDE DPVFIATCTERELRKRKRK FSLWVRQ

1 | MMCDSGTKYMINGMPYLGRGT----QTNGVPLGEYVVKELSKPVHGSCRNITC DNWFTSIPLAKNLLQEPYKLTIVGTVRSNKREIPEVLKNSRSP-VG TSMFCFDGP
2 | CSSTGFIIQIYVHLKEGGGPDGLDALKNKQQLHSMVARSLCRNAAGKNY IIF TGPSITSLTLFEEFEKQGIYCCGLLRARKSDCTGLPLSMLTNPATPPARG QYQIKMKG
3 | CSSTGFIIQIYVHLKEGGGPDGLDALKNKQQLHSMVARSLCRNAAGKNY IIF TGPSITSLTLFEEFEKQGIYCCGLLSSRKS DCTGLPPSMLTNPATPLARG QHQIRTKG

1 | LTLVSYKPKPAKMVYLLSSCDEASINESTGK----PQMVMYINQTKGGVD TLDQMSVMTCSRKTNRWPMALLYGMINIACINSFIIYSHNVSSKGEKVQSRKKFMR
2 | NMSLICWYKGFHFRFLTNAYSPVQ QGVI IKRKSGEIPCPLAVEAFAAHL SYICRYD DKYSKYFISHKPNKTWQQVFWFAISIAINNAYILYKMSDAYHVKRY SRAQFGER
3 | NMSLICWYKGFHFRFLTNAYSPVQ QGVI IKRKSGEIPCPLAVEAFAAHL SYICRYD DKYSKYFISHKPNKTWQQVFWFAISIAVNNAYILYKMSDAYHVKRY SRAQFGER

1 | NLYMSLTSSFMKRLEAPTLKRYLRDNISNILPNEVPGTSDDSTEPEVMKKRTY TY PSKIRKANAS KK KKVIRE NIDM QS F PB (Ifp2 transposase)
2 | LVRELLGLEDA SPTH----- PGBD5 Hs524
3 | LVRELLGLEDS SPAH----- PGBD5 Mm523

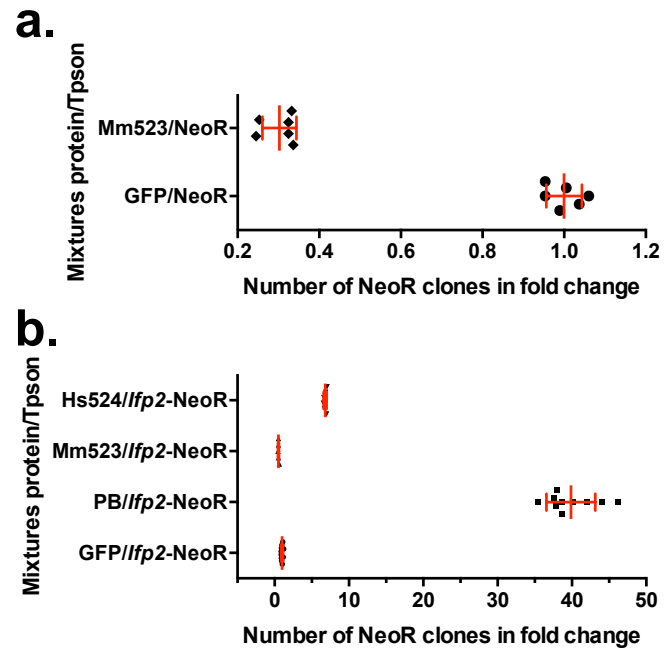
```

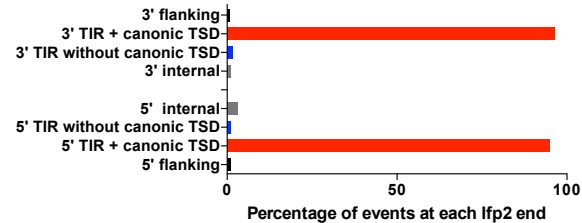
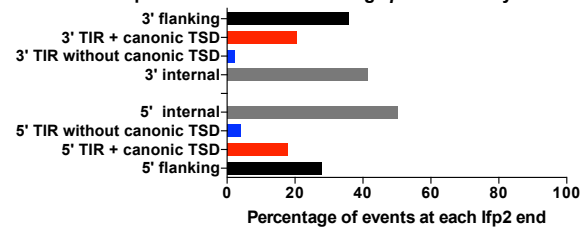
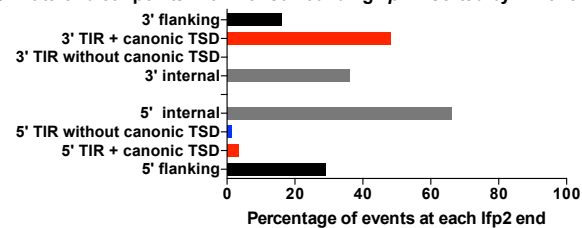
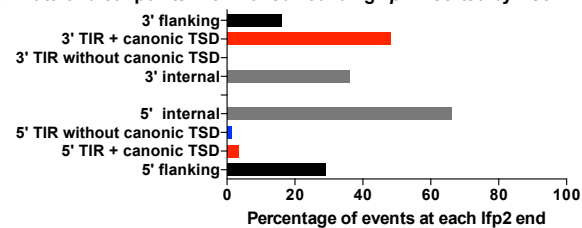
b. PB.NLS-1-558 sequence.

```

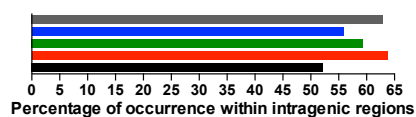
MPKKRKRKQRSAA TMGSSLDDEHILSALLOSDDELVGEDSDSEISDHVSEDDVQSDTEEAFIDEVHEVQPTSSGSEILDEQNVIQPGSSLASNRILTLPQRTIRGKNK
HCWSTSKSTRSRVSALNIVRSORGPTR-MCRNIYDPLLCFKLFTDEIISEIVKWTNAEISLKRRESMTGATFRDTNEDEIYAFFGILVMTAVRKDNHMSTDDLFRDRL
SMVYVSVMSRDRDFLIRCLRMDDKSIRPTLRENDVFTPVKRIWDLFIHQCIQNYTPGAHLTI DEQLLGFGRGCPFRMYIPNKPS KYGIKILMMCDSGTKYMINGMPYLG
RGTQTNGVPLGEYVVKELSKPVHGSCRNITC DNWFTSIPLAKNLLQEPYKLTIVGTVRSNKREIPEVLKNSRSPVGTSMFCFDGPI LTLVSYKPKPAKMVYLLSSCDEDA
SINESTGKQPQMVMYINQTKGGVD TLDQMSVMTCSRKTNRWPMALLYGMINIACINSFIIYSHNVSSKGEKVQSRKKFMRNLYMSLTSSFMKRLEAPTLKRYLRDNISN
ILPNEVPG

```



a. Rate of breakpoints within or surrounding *lfp2* inserted by PB**b. Rate of breakpoints within or surrounding *lfp2* inserted by PB.NLS-1-558****c. Rate of breakpoints within or surrounding *lfp2* inserted by Mm523****d. Rate of breakpoints within or surrounding *lfp2* inserted by Hs524**

a. Occurrence of *pble* insertion sites in intragenic regions



b. Occurrence of *pble* insertion sites around TSS (± 5 kpb)

