# Internal initiation of reverse transcription in a Penelope-like retrotransposon

Chris J. Frangieh<sup>1-6</sup>, Max E. Wilkinson<sup>1-5</sup>, Daniel Strebinger<sup>1-5</sup>, Jonathan Strecker<sup>1-5</sup>, Michelle Walsh<sup>1-5</sup>, Guilhem Faure<sup>1-5</sup>, Irina A. Yushenova<sup>7</sup>, Rhiannon K. Macrae<sup>1-5</sup>, Irina R. Arkhipova<sup>7\*</sup>, Feng Zhang<sup>1-5\*</sup>

Affiliations: (1) Howard Hughes Medical Institute, Cambridge, MA 02139, USA; (2) Broad Institute of MIT and Harvard, Cambridge, MA 02142, USA; (3) McGovern Institute for Brain Research, Massachusetts Institute of Technology, Cambridge, MA 02139, USA; (4) Department of Brain and Cognitive Science, Massachusetts Institute of Technology, Cambridge, MA 02139, USA; (5) Department of Biological Engineering, Massachusetts Institute of Technology, Cambridge, MA 02139, USA; (6) Department of Electrical Engineering and Computer Science, Massachusetts Institute of Technology, Cambridge, MA 02139, USA; (7) Josephine Bay Paul Center for Comparative Molecular Biology and Evolution, Marine Biological Laboratory, Woods Hole, MA, 02543, USA

<sup>\*</sup> Correspondence should be addressed to zhang@broadinstitute.org (F.Z.), iarkhipova@mbl.edu (I.A.).

## **ABSTRACT**

1

- 2 Eukaryotic retroelements are generally divided into two classes: long terminal repeat (LTR)
- 3 retrotransposons and non-LTR retrotransposons. A third class of eukaryotic retroelement, the
- 4 Penelope-like elements (PLEs), has been well-characterized bioinformatically, but relatively
- 5 little is known about the transposition mechanism of these elements. PLEs share some features
- 6 with the R2 retrotransposon from *Bombyx mori*, which uses a target-primed reverse transcription
- 7 (TPRT) mechanism, but their distinct phylogeny suggests PLEs may utilize a novel mechanism
- 8 of mobilization. Using protein purified from E. coli, we report unique in vitro properties of a
- 9 PLE from the green anole (Anolis carolinensis), revealing mechanistic aspects not shared by
- 10 other retrotransposons. We found that reverse transcription is initiated at two adjacent sites
- 11 within the transposon RNA that is not homologous to the cleaved DNA, a feature that is reflected
- in the genomic "tail" signature shared between and unique to PLEs. Our results for the first
- 13 active PLE in vitro provide a starting point for understanding PLE mobilization and biology.

### 14 Keywords: reverse transcriptase, GIY-YIG endonuclease, hammerhead ribozyme

### MAIN TEXT

15

16

17 Eukaryotic retroelements are generally divided into two classes: long terminal repeat (LTR)

18 retrotransposons and non-LTR retrotransposons. Penelope-like elements (PLEs), which are

19 found in fish, amphibians, reptiles, insects, protists, plants, and fungi, have phylogenetically and

- structurally been defined as a third class of eukaryotic retroelement <sup>1,2</sup>. Canonical PLEs encode a
- 21 GIY-YIG endonuclease (EN) domain alongside a telomerase-like reverse transcriptase (RT)
- domain in a single open reading frame (ORF) <sup>3</sup> (Fig. 1A-B). The presence of an EN and RT
- 23 domain suggests PLEs may utilize a target-primed reverse transcription (TRPT) mobilization
- 24 mechanism similar to the R2 retrotransposon from *Bombyx mori* (R2Bm), where the nicked
- 25 target DNA is used as a primer for reverse transcription of the transposon RNA<sup>4</sup>. Alternatively,
- since PLEs are phylogenetically distinct from both LTR and non-LTR retrotransposons, PLEs
- 27 may utilize a novel mechanism of mobilization. To date, no studies have purified an active, full-
- 28 length PLE protein, limiting our mechanistic understanding of these elements. Here we report the
- 29 in vitro characterization of a PLE from the green anole (PAc, Poseidon from Anolis

30 carolinensis), revealing mechanistic aspects not shared by other retrotransposons that explain the 31 unusual "tail" structure observed in most of the inserted genomic copies. 32 33 PLEs are organized in tandem or partial-tandem repeats in the genomes of their host organisms; 34 the tandem structure leads to the formation of repetitive flanking UTR sequences referred to as pseudo-LTRs (pLTRs) <sup>1,3,5</sup>. RNA sequencing from A. carolinensis confirmed the organization of 35 36 PAc in tandem repeats at the endogenous loci (Fig. 1B). The predicted domain structure of PAc is consistent with that of canonical PLEs from the Penelope/Poseidon clade <sup>3</sup>, consisting of a 37 putative N-terminal nucleic acid binding domain, a central RT domain, and a C-terminal GIY-38 YIG EN domain. RNA sequencing reads from A. carolinensis (Table S1) revealed that the PAc 39 40 element is expressed, with increased coverage at the pLTR-ORF junction as well as at the C-41 terminus of the ORF (Fig. 1B). 42 To detect transposition events in vitro, we developed an assay that does not introduce bias from 43 PCR amplification (Fig. 1C)<sup>4</sup>. Briefly, PAc protein purified from E. coli is incubated with a 44 45 template RNA and a dsDNA target where each strand is 5'-labelled with a different fluorophore, 46 corresponding to the PAc ORF-pLTR junction, to reconstitute an insertion reaction. Both strands 47 of the dsDNA target are blocked with a 3' inverted dT to prevent non-specific extension. The end-labeled dsDNA target is then analyzed on a denaturing gel to visualize cleavage and 48 49 extension activity separately for the top and bottom strands. Using this assay, we found that PAc 50 cleaves and extends both the top and bottom strands (Fig. 1C), unlike the R2Bm and human LINE-1 retrotransposons, which have only been shown to nick and extend the bottom strand<sup>4,6–8</sup>. 51 52 53 To determine important features of the RNA template for transposition, we performed 5' and 3' 54 truncations of the RNA template. Consistent with other retrotransposons, 5' truncations reduced 55 the length of the insertion product and showed that the full pLTR sequence is not required for insertion, with +20 being the shortest active 5' truncation point (numbered relative to the ORF 56 57 start codon) (Fig. 1D). 3' truncations defined the minimum RNA as requiring the first 172 nt of 58 the ORF (Fig. 1E). We therefore hypothesize that the protein binding motif is contained within the first 172 nt of the PAc ORF. When we added additional sequences to the 3' end of the 59 60 template RNA beyond +235 (i.e., +266, +289, +315, +347, +378), a low intensity band was

61 observed at the top of the gel, likely arising from priming off the 3' end of the template; 62 however, the dominant insertion band did not change in size with the truncations/additions, 63 suggesting that reverse transcription initiates somewhere internal to the 3' end of the template RNA (Fig. 1E). 64 65 66 We developed a strand-specific next generation sequencing (NGS) assay to more precisely probe 67 the PAc insertion products (Fig. S1A). Analysis of cleavage sites reveals a preference for AT-68 rich DNA, specifically a 5'-ATT-3' motif downstream of the ORF-pLTR junction, that 69 represented the most frequently cleaved site on each strand (Fig. 1F). Preference for AT-rich 70 DNA cleavage has previously been shown for a purified endonuclease domain of the Penelope 71 element from D. virilis<sup>9</sup>; however, in contrast to the D. virilis endonuclease domain, PAc cleaves 72 both the top and bottom strands. This may be due to our use of the polyprotein rather than the 73 endonuclease domain in isolation, as the PAc polyprotein is predicted to form a dimer by AlphaFold2 (**Fig. 2**) $^{10,11}$ . 74 75 76 Analysis of RT start sites mapped from the NGS data confirms that reverse transcription 77 preferentially initiates internal to the 3' end of the RNA template (Fig. 1G). There does not 78 appear to be a difference in RT start site between the top and bottom strands. Examining the 79 relationship between the EN cut site and RT start site reveals the most common enzymatic 80 activity as cleavage of a 5'-ATT-3' motif followed by reverse transcription beginning in an AC-81 rich stem loop located 111-118 bp into the start of the ORF (Fig. S1B). Surprisingly, RT initiates 82 internal to the 3' end of the RNA template at a site that is not homologous to the cleaved DNA, a feature that has not been seen with other retrotransposons. Previously, the D. virilis Penelope EN 83 84 was shown to recognize and cleave plasmids containing an equivalent Penelope fragment, and 85 limited tiling experiments showed that the so-called "tail" sequence, an optional 30-40 bp 86 extension found in a fraction of genomic copies, was required for EN recognition followed by cleavage, although the exact cleavage site was not defined<sup>9</sup>. This tail sequence clearly 87 88 corresponds to the sequence between the two observed RT start sites (Fig. 1G). In Fig. S2, we 89 document this correspondence by querying A. carolinensis genome assembly as well as publicly 90 available RNA-seq reads with the template RNA sequence, and visualizing the outputs at the

coverage and nucleotide level. The ends of the PAc genomic insertions (initially defining the

92 "tail" sequence, arrows) and the ends of the corresponding RNA-seq reads fall precisely onto the 93 RT initiation sites that are shown in Fig. 1F. 94 The PAc locus contains a hammerhead ribozyme (HHR) sequence, consistent with findings in 95 96 other PLEs, containing regions corresponding to stem I, stem II, and loop II of the canonical HHR structure<sup>12,13</sup>. The HHR encoded by PAc is conserved across several Anolis species (Fig. 97 98 3A). The PAc ORF begins at the start of stem II, and the HHR cut site is 26 bp downstream of 99 the start codon. The conserved residues C3 and G8 in the HHR catalytic core occur prior to the 100 start codon (Fig. 3A). The mutations C3G and G8C have been shown to abolish HHR cleavage activity, while the double mutant C3G/G8C partially restores ribozyme cleavage activity<sup>14</sup>. 101 102 Consistent with this, individual G8C and C3G mutations eliminate PAc HHR activity, while the 103 C3G/G8C double mutant partially restores cleavage activity (Fig. 3B). HHR catalytic mutants do 104 not affect RNA insertion activity and do not affect the size of the insertion band, however, 105 suggesting that HHR cleavage is not a necessary step for PAc-mediated transposition in vitro and 106 that unprocessed RNA is the dominant substrate used for insertion (Fig. 3B). In the GenBank 107 RNA-seq datasets, we failed to detect reads that begin or end at the presumptive HHR cleavage 108 site (Fig. S2B). 109 110 Interestingly, the HHR shows no cleavage activity without the addition of protein (Fig. 3B). To 111 confirm that RNA cleavage activity is not protein-catalyzed, we purified PAc catalytic mutants 112 G19A and Y17A (see Fig. 1A), which have been shown to eliminate GIY-YIG endonuclease activity<sup>15</sup>. GIY-YIG endonuclease mutants show no DNA cleavage activity as expected, but do 113 114 not impact RNA cleavage activity, suggesting that the PAc HHR is a protein-assisted ribozyme 115 (Fig. 3B). It is possible that in vitro binding of the PAc ORF to the RT initiation site on the RNA 116 template creates more favorable conditions for RNA folding into HHR-competent conformation 117 (Fig. S2C), and that in vivo such conditions may arise at the later stages of the transposition 118 cycle. 119 120 We next searched in vivo data from A. carolinensis to confirm our observations of AT-rich DNA 121 cleavage followed by reverse transcription initiation in an AC-rich stem loop. We searched 122 genome sequencing reads from the NCBI SRA database (Table S1) for motifs corresponding to

123 a PAc insertion that initiated in the identified AC-rich stem loop and homed to an AT-rich gDNA 124 site. For several reverse transcription initiation sites in the AC-rich stem loop, we find evidence 125 of AT-rich homing sites (Fig. S1B). 126 127 Our results, taken together, broadly outline the following scenario of PAc mobilization. The repetitive PAc locus is transcribed and translated, likely from a promoter in its pLTR<sup>5</sup>, leading to 128 129 an mRNA consisting of multiple PAc copies alongside the PAc protein. The PAc protein binds 130 its transposon RNA, and the PAc RNP complex then binds and nicks AT-rich gDNA, either in its 131 own pLTR or elsewhere in the genome, and initiates reverse transcription of its template RNA within an AC-rich stem loop. Following second strand synthesis and host repair, this process 132 133 eventually leads to a new PAc copy at an AT-rich gDNA sequence through additional steps that 134 are yet to be validated experimentally. The observation that reverse transcription is initiated at a 135 site within the transposon RNA that is not homologous to the cleaved DNA is unique to PLEs. Our system provides a controlled in vitro model to further our understanding of PLE 136 137 mobilization and biology. 138 List of abbreviations 139 EN, endonuclease; gDNA, genomic DNA; HHR, hammerhead ribozyme; NGS, next generation 140 sequencing; ORF, open reading frame; PAc, Poseidon A. carolinensis; PLE, Penelope-like 141 elements; pLTR, pseudo (or Penelope) long terminal repeat; RNP, ribonucleoprotein; RT, reverse transcription; SRA, short read archive; TPRT, target-primed reverse transcription. 142 143 FIGURE LEGENDS 144 Fig. 1. Cleavage and extension activity of the purified recombinant full-length PAc protein in vitro. (A) Amino acid sequence alignment of the core catalytic RT motifs 4 and 5<sup>16</sup> 145 146 (underlined) and the selected EN catalytic motifs (underlined, with mutagenized residues marked 147 by asterisks) for four representatives from the Penelope/Poseidon PLE clade (Repbase entries 148 from A. carolinensis, Petromyzon marinus, Branchiostoma floridae, Drosophila virilis); 149 numbering on the top corresponds to PAc ORF. (B) Locus map showing repetitive structure of 150 the PAc element. The schematic shows three example repeats with domains highlighted in blue 151 (nucleic acid binding), white (reverse transcriptase), and red (endonuclease). pLTR, pseudo long 152 terminal repeat; RT, reverse transcriptase; ORF, open reading frame; HHR, hammerhead

ribozyme motif. Read coverage for one copy of the PAc retroelement (pLTR-ORF-pLTR from the Repbase consensus, File S1) from publicly available A. carolinensis RNA sequencing data is shown below the locus map. The side panel shows the degree of PAc ORF purification (Methods). (C) Overview of the in vitro assay consisting of an end-labeled dsDNA target, RNA template, and PAc protein. The schematic shows the expected effect when adding protein and RNA (cleavage), and adding protein, RNA, and dNTPs (extension). Denaturing gel shown for one RNA template for both the bottom strand (Cy5 fluorophore) and top strand (fluorescein fluorophore). The expected significance of each band as it relates to the diagram is shown by the line drawing on the left-hand side of both gels. (D) The effect of truncating the RNA template from the 5' end on cleavage and extension activity using the in vitro assay. The 3' end of the template is held constant at +235 while the 5' end is tiled from -126 to +200 where bases are labeled relative to the start of the ORF. (E) The effect of changing the RNA template from the 3' end on cleavage and extension activity using the in vitro assay. The 5' end of the template is held constant at -145 while the 3' end is tiled from -63 to +378 where bases are labeled relative to the start of the ORF. (F) Cleavage frequency for both the top and bottom strands on a dsDNA target corresponding to the PAc ORF-pLTR junction. Bars indicate the percent of total reads mapping to a cleaved product at that location. (G) Reverse transcription initiation sites for a dsDNA target corresponding to the PAc ORF-pLTR junction and an RNA template spanning the PAc pLTR-ORF junction. Bars indicate the percent of total reads mapping to an insertion product that begins at that location.

173

174

175

176

177

178

179

180

181

182

183

153

154

155

156

157

158

159

160

161

162

163

164

165

166

167

168

169

170

171

172

Fig. 2. 3D structure modeling of the PAc ORF. The dimeric structure was predicted using AlphaFold2 within the Colabfold framework. Three independent replicates, each with 40 cycles, were used. Each replicate consistently yielded the identical dimeric conformation, represented in green and blue (shown in surface and cartoon representations). Two key regions of interaction, outlined by orange dashed circles, stabilize the dimer. The first region features an extended helix (positions 27-173) adopting a head-to-tail homodimeric interaction, predominantly sustained by an intricate hydrophobic network. The second region consists of two helices (positions 232-250 and 795-818) from each dimeric partner, interacting head-to-tail and forming a V shape. This interaction is characterized by two symmetrical salt bridges (indicated by red dashed lines) on the periphery of the interaction patch, encasing a hydrophobic core (highlighted in purple).

184	
185	Fig. 3. Structure and cleavage properties of the PAc hammerhead ribozyme (HHR). (A)
186	The HHR encoded by PAc is conserved across other Anolis species. Stem I, stem II, and loop II
187	from the canonical HHR structure along with the HHR cut site and PAc start codon are
188	annotated. pLTR, pseudo long terminal repeat; ORF, open reading frame. Visualization of the
189	PAc HHR secondary structure is shown next to the sequence alignment. Residues C3 and G8 in
190	the HHR catalytic core are boxed. (B) Visualizing the role of the HHR and GIY-YIG
191	endonuclease in RNA processing and extension in vitro. G8C and C3G are HHR mutants while
192	G8C + C3G is a partial rescue. G19A and Y17A are key mutations in the active site of the GIY-
193	YIG endonuclease.
194	
195	SUPPLEMENTARY FIGURE LEGENDS
196	
197	Fig. S1. Details of the NGS assay and cleavage/insertion site localization on PAc RNA. (A)
198	Schematic of the workflow for an NGS assay used for strand-specific sequencing of reaction
199	products. Biotinylated dsDNA is incubated with PAc protein, template RNA and dNTPs as
200	described in Methods, extracted from a denaturing gel, immobilized on streptavidin beads,
201	adapter ligated, and amplified prior to sequencing. (B) Predicted RNA secondary structure for
202	the conserved RNA sequence at the start of the PAc ORF. Shapes indicate the highest frequency
203	reverse transcription initiation sites as predicted from the NGS data shown in Fig. 1F, while pink
204	circles indicate all reverse transcription initiation sites (i.e. "tail" boundaries). WebLogo <sup>17</sup> plots
205	are shown for AT-rich homing sites for the reverse transcription initiation sites as calculated
206	from gDNA sequencing data from A. carolinensis (Table S1). Shapes correspond to the start sites
207	shown on the RNA secondary structure prediction by RNAfold <sup>23</sup> .
208	
209	Fig. S2. Visualization of gDNA and RNA sequences homologous to the N-terminal part of
210	PAc ORF. GenBank databases were queried with nt -100 to +200 relative to the PAc AUG
211	codon, so that the numbers shown on the top represent the consensus PAc numbering used
212	throughout the text plus 100, i.e. the approx. 34-nt "tail" sequence roughly spans nt 85-119 of the
213	consensus (vertical arrows). Screenshots display sequence alignments for a 150-bp window in
214	the NCBI MSA Viewer 1.25.0; plots on the top show query coverage and the number of reads at

215	peak coverage. (A) NCBI megablast search of the A. carolinensis WGS assembly AAWZ; (B)
216	Example of a megablast search of the $A.\ carolinensis\ 150$ -nt RNA-seq SRA reads $^{25}$ from Table
217	S1 (accession SRR14288908) with 'max target sequences' set at 5000. The rightmost part of the
218	alignment shows the reads extending from the "tail" region into the body of the element,
219	indicating ongoing transcription of full-length copies not visible in the plot in Fig. 1A due to
220	much higher coverage in the pLTR region. Position 30 (#130 on the figure) corresponds to the
221	expected HHR cleavage site (asterisk), however there are no reads beginning or ending at this
222	site, and the coverage plots on the top indicate no discontinuities in this region. Other SRA
223	accessions display similar patterns. (C) Alternative PAc RNA structure predictions using the
224	deep-learning-based MXfold2 server <sup>24</sup> . The top RNA and the RT start sites are the same as in
225	Fig. S1B, beginning with the presumed HHR cleavage site. The bottom RNA includes the
226	uncleaved HHR motif as shown in Fig. 3 (PAc nt -9 to +44) with folded stem-loop II, and
227	outlines the hypothetical interaction area with PAc RT moiety (cloud-like) in a large loop near
228	the first RT start site, which needs to undergo unfolding of the conserved HHR catalytic core (nt
229	-6 to 1) for reverse transcription to occur through it. Gray lines indicate base-pairing that would
230	be required to form HHR stem I. The expected cleavage site is indicated between C29-A30.
231	
232	
233	Table S1. SRA accessions used in preparing Fig. S1B (DNA), Fig. 1B and Fig. S2 (RNA).
234	
235	File S1. Consensus sequence of the Poseidon_Ac entry from Repbase <sup>19</sup> (in Genbank format).
236	The database requires subscription since 2019, however this entry was assembled from Sanger
237	reads and deposited by I.A. as part of ref. 20 under initial assumption that it would be freely
238	available to researchers.
239	DECLARATIONS
240	Ethics approval and consent to participate: Not applicable

Consent for publication: Not applicable

242	Availability of data and materials: All data generated or analysed during this study are
243	included in this published article and its supplementary information files. Materials can be
244	obtained from Addgene.
245	Competing interests: F.Z. is a scientific advisor and cofounder of Editas Medicine, Beam
246	Therapeutics, Pairwise Plants, Arbor Biotechnologies, Aera Therapeutics, and Moonwalk
247	Biosciences. F.Z. is a scientific advisor for Octant.
248	Funding: M.E.W is supported by a Helen Hay Whitney Foundation Postdoctoral Fellowship and
249	the Howard Hughes Medical Institute. F.Z. is supported by the Howard Hughes Medical
250	Institute; the Hock E. Tan and K. Lisa Yang Center for Autism Research at MIT; the Yang-Tan
251	Molecular Therapeutics Center at McGovern, the BT Charitable Foundation, and by the Phillips
252	family and J. and P. Poitras. I. A. and I.Y. are supported by NSF MCB-2139001.
253	Author contributions: C.J.F. and F.Z. conceived the project. C.J.F designed and performed
254	experiments with input from all authors. J.S., M.E.W., D.S., M.W., I.Y. and I.A. contributed to
255	the development of the project through experiments and discussion. C.J.F. and G.F. analyzed
256	genomic and structural data. F.Z. supervised the research and experimental design with support
257	from R.K.M. C.J.F, R.K.M, I.A. and F.Z. wrote the manuscript with input from all authors.
258	Acknowledgements: We thank all members of the Zhang lab for helpful discussions and
259	support.
260	METHODS
261	
262	Alignment of RNA and DNA sequencing reads from A. carolinensis
263	
264	RNA and DNA sequencing reads from A. carolinensis were downloaded from the NCBI
265	Sequence Read Archive (SRA) database in FASTQ file format <sup>18</sup> . SRA accession numbers are
266	included in <b>Supplemental Table S1</b> . Read files were aligned to the PAc locus map (Repbase ID:
267	Poseidon_Ac; Supplemental File 1) using Bowtie 2 in "sensitive local" mode 19,20. Alignment
268	files in .sam format were concatenated and parsed to isolate CIGAR strings. A match character
269	("M") was used to sum coverage across the PAc reference map for all CIGAR strings.

270 Alignment of representative Repbase sequences from the Penelope/Poseidon clade in Fig. 1A 271 was generated with MAFFT v.7 and visualized with Jalview v. 2.11.3.2 using the Clustal color scheme $^{21,22}$ . 272 273 274 Molecular cloning of plasmids 275 276 PLE sequences were downloaded from the Repbase database (PAc Repbase ID: Poseidon Ac) and cloned by Genscript<sup>23,24</sup>. Point mutations (e.g., G19A, Y17A) were generated by KLD 277 278 cloning (New England Biolabs). Several iterations of N and C-terminal solubility and affinity 279 tags were cloned by Gibson Assembly (New England Biolabs). All plasmid sequences were 280 verified via Tn5 tagmentation and next generation sequencing (Illumina)<sup>25</sup>. 281 282 Recombinant protein expression in *E. coli* 283 284 An inducible bacterial expression plasmid coding for the PAc ORF with an N-terminal 14x His-285 MBP tag and C-terminal TwinStrep tag was transformed into Rosetta(DE3) competent cells 286 (Millipore Sigma). A single colony was picked and placed in 33 mL of starter for overnight 287 incubation. 5 mL of starter was used to inoculate 1 L of media. Cultures were grown at 37°C in 288 Terrific Broth (Thermo Fisher Scientific) until  $OD_{280 \text{ nm}} = 0.8$ . Cultures were then cooled to 289 18°C, and protein expression was induced with the addition of Isopropyl β-D-1-290 thiogalactopyranoside (IPTG) (Gold Biotechnology) to a final concentration of 250 µM followed 291 by 18 h of growth at 18°C. Cells were harvested by centrifugation at 5000g for 10 min, and 292 pellets were resuspended in 50 mM Tris-HCl pH 7.5, 1 M NaCl, 10% glycerol, 5 mM beta-293 mercaptoethanol (BME), and cOmplete ULTRA EDTA-free protease inhibitor (Roche). Cells 294 were lysed using an LM20 microfluidizer (Microfluidics), and lysate was cleared by 295 centrifugation at 18,000 RPM for 20 min. Cleared lysate was bound to Strep-Tactin Superflow 296 Plus (Qiagen) resin at 4°C for 60 min with rotation. The resin was washed several times at 4°C, 297 and then washed one final time with elution buffer (20 mM HEPES-KOH pH 7.9, 500 mM KCl, 298 10% glycerol, and 1 mM tris carboxyl ethyl phosphene (TCEP)). Resin was eluted 10 times with 299 1 mL elution buffer supplemented with 5 mM d-Desthiobiotin (Millipore Sigma). Elutions were 300 run on an SDS-PAGE gel followed by Coomassie blue staining, and elutions containing high

concentrations of proteins were pooled, aliquoted in  $100 \,\mu\text{L}$ , snap frozen in liquid nitrogen, and stored at -80°C. From 12 L, we obtained ~2.5 mg purified Pac protein.

303304

301

302

## In vitro cleavage and extension assays

305

306 In vitro reactions were performed in 20 mM HEPES-KOH pH 7.9, 400 mM KAc, 100 mM KCl, 307 5% glycerol, 5 mM MgAc, 0.2 mM TCEP, and 1 mM d-Desthiobiotin. Reactions contained 500 308 nM PAc protein, 500 nM RNA template, and 10 nM labeled dsDNA target. RNA templates for 309 in vitro transcription (IVT) were produced by PCR with a T7 promoter added to the forward 310 primer. For IVT, PCR reactions were diluted 1/10 in the reaction mixture containing 4 mM each 311 NTP, 20 mM MgCl<sub>2</sub>, 40 mM Tris-HCl pH 8.0, 10 mM DTT, 1 mM spermidine, 85 µg/mL of homemade T7 RNA polymerase, and incubated at 37°C for 90 min. The pyrophosphate 312 313 precipitate was pelleted and removed, and reactions were treated with 1/100 vol of RNase-free 314 DNase I (NEB) at 37°C for 15 min. The 361-nt RNA template used in most assays contained the 315 126-nt UTR plus the first 235 nt of the PAc ORF from the Poseidon Ac Repbase consensus 316 317 uugcuuaguuuucuccauaccucacaaccucugaggAUGccugccauagaugugggcgaaacgucaggagagaaugcuucuggaacau 318 ggccacacagcccgaaagacauacaacacacugugaucccggccaugaaagccuucgacaacacauugaacaucucuacggggagaaauu 319 320 g) (full consensus in file S1). dsDNA target was prepared by annealing a labeled top strand (/56-321 FAM/acaaaggatgcccagaggaagaagaagacagcagataagcttttcaatgctaattaaagtgattaactacacaacatt/3Invd 322 T/) and a labeled bottom strand 323 (/5Cy5/aatgttgtgtagttaatcactttaattagcattgaaaagcttatctgctgtcttcttcctgcctctggggcatcctttgt/3InvdT/). 324 Reactions were incubated at 32°C for 90 min followed by mixing with equal volumes of 2X 325 TBE-Urea sample buffer (90 mM Tris base, 90 mM boric acid, 2 mM EDTA, 12% Ficoll Type 326 400, 7 M urea, and 0.02% bromophenol blue). The chosen temperature (32°C) corresponds to the average temperature experienced by anole lizards in their natural environments <sup>22</sup>. Reactions 327 328 were then boiled at 95°C for 3 minutes, followed by running on a precast 10% TBE-Urea gel 329 (Invitrogen) at 400V for 12 min. Fluorescent signals in all gels were visualized using a 330 ChemiDoc (BioRad). 331

#### 332 In vitro NGS assay 333 334 In vitro reactions were performed as described above, except the dsDNA target was biotinylated 335 on either the top strand or bottom strand rather than end-labeled with a fluorophore. Biotinylated 336 DNA primers were ordered (Integrated DNA Technologies) and used in a PCR with Q5 Hot 337 Start High-Fidelity polymerase (New England Biolabs) to create a dsDNA target. Reactions were 338 stopped with the addition of 5 µL RNase A (New England Biolabs) followed by incubation at 339 37°C for 30 min. A phenol chloroform extraction was then performed with UltraPure 340 Phenol:Chloroform:Isoamyl Alcohol (Thermo Fisher Scientific) followed by ethanol 341 precipitation and resuspension in 2X TBE-Urea sample buffer (90 mM Tris base, 90 mM boric 342 acid, 2 mM EDTA, 12% Ficoll Type 400, 7 M urea, and 0.02% bromophenol blue). Reactions were gel extracted from a 10% TBE-Urea gel and passed through a 1 mL syringe to shear gel 343 344 fragments. Sheared gel fragments were incubated overnight at 4°C in 0.3M sodium acetate pH 345 5.5 and 10 mM EDTA to allow diffusion of the DNA out of the gel. Another ethanol precipitation was performed prior to ligation. A 5' adenylated and 3' capped adapter 346 347 (/5rApp/CTGTCTCTTATACACATCTCCGAGCCCACGAGAC/3SpC3/) was ligated to 348 purified DNA with a thermostable 5' App DNA/RNA ligase (New England Biolabs) at 65°C for 349 16 h. Following proteinase K (New England Biolabs) treatment, ligated DNA was immobilized 350 on Dynabeads M-270 Streptavidin (Thermo Fisher Scientific) by incubation at room temperature 351 with agitation for 30 min in binding/wash buffer (5 mM Tris-HCl pH 7.5, 0.5 mM EDTA, and 1 352 M NaCl). The beads were washed several times with binding/wash buffer to remove excess 353 adapter. Beads were input directly into a PCR with KAPA HiFi HotStart ReadyMix (Roche) for 354 30 cycles. NGS libraries were gel extracted and quantified using Qubit Fluorometric 355 Quantification (Thermo Fisher Scientific). Libraries were sequenced on a MiSeq (Illumina) with 356 100 cycles read 1, 8 cycles index 1, 8 cycles index 2, and 100 cycles read 2 supplemented with 357 10% PhiX Control v3 (Illumina) for diversity. 358 359 **Processing NGS data** 360 361 NGS data was first trimmed for low quality reads (Q-score < 30), followed by removing reads 362 that did not contain the expected adapter sequence at the start of the read. Trimmed FASTQ files

- 363 were aligned to the dsDNA target reference map to identify the cleavage site using Bowtie 2 in "local" mode<sup>19,20</sup>. Sequences in the read beyond the dsDNA cleavage site were assumed to be 364 365 RNA-templated insertions, and these sequences were similarly mapped to an RNA template 366 reference map to determine the site of reverse transcription initiation. 367 368 In vitro HHR assay 369 370 In vitro HHR reactions were performed in 20 mM HEPES-KOH pH 7.9, 400 mM KAc, 100 mM 371 KCl, 5% glycerol, 5 mM MgAc, 0.2 mM TCEP, and 1 mM d-Desthiobiotin. Reactions contained 372 500 nM PAc protein and 500 nM RNA (-126/+235). Reactions were incubated at 32°C for 2 h 373 followed by mixing with equal volumes of 2X TBE-Urea sample buffer (90 mM Tris base, 90
- 374 mM boric acid, 2 mM EDTA, 12% Ficoll Type 400, 7 M urea, and 0.02% bromophenol blue).
- Reactions were then boiled at 95°C for 3 min, followed by running on a precast 10% TBE-Urea
- gel (Invitrogen) at 400V for 12 min. Previously, we observed self-cleavage in PLE HHR at Mg<sup>2+</sup>
- 377 concentrations varying from 3 mM to 25 mM <sup>26</sup>.

379 380 **REFERENCES** 

- 1. Evgen'ev, M. B. *et al.* Penelope, a new family of transposable elements and its possible role in hybrid dysgenesis in Drosophila virilis. *Proc. Natl. Acad. Sci. U. S. A.* **94**, 196–201 (1997).
- 383 2. Arkhipova, I. R., Pyatkov, K. I., Meselson, M. & Evgen'ev, M. B. Retroelements containing introns in diverse invertebrate taxa. *Nat. Genet.* 33, 123–124 (2003).
- 385 3. Craig, R. J., Yushenova, I. A. & Rodriguez, F. An Ancient Clade of Penelope-Like Retroelements with
  386 Permuted Domains Is Present in the Green Lineage and Protists, and Dominates Many Invertebrate Genomes.
- 387 Biology and Evolution (2021).
- Luan, D. D., Korman, M. H., Jakubczak, J. L. & Eickbush, T. H. Reverse transcription of R2Bm RNA is
  primed by a nick at the chromosomal target site: a mechanism for non-LTR retrotransposition. *Cell* 72, 595–605 (1993).
- Schostak, N. *et al.* Molecular dissection of Penelope transposable element regulatory machinery. *Nucleic Acids Res.* 36, 2522–2529 (2008).
- 393 6. Eickbush, T. H. & Eickbush, D. G. Integration, Regulation, and Long-Term Stability of R2 Retrotransposons.

- 394 *Microbiol Spectr*, *3*. (2015).
- Wilkinson, M. E., Frangieh, C. J., Macrae, R. K. & Zhang, F. Structure of the R2 non-LTR retrotransposon
- initiating target-primed reverse transcription. *Science* **380**, 301–308 (2023).
- 397 8. Thawani, A., Ariza, A. J. F., Nogales, E. & Collins, K. Template and target-site recognition by human LINE-1
- 398 in retrotransposition. *Nature* **626**, 186–193 (2024).
- 9. Pyatkov, K. I., Arkhipova, I. R., Malkova, N. V., Finnegan, D. J. & Evgen'ev, M. B. Reverse transcriptase and
- 400 endonuclease activities encoded by Penelope-like retroelements. *Proc. Natl. Acad. Sci. U. S. A.* 101, 14719–
- 401 14724 (2004).
- 402 10. Jumper, J. et al. Highly accurate protein structure prediction with AlphaFold. *Nature* **596**, 583–589 (2021).
- 403 11. Mirdita, M. et al. ColabFold: making protein folding accessible to all. Nat. Methods 19, 679–682 (2022).
- 404 12. Scott, W. G., Horan, L. H. & Martick, M. The hammerhead ribozyme: structure, catalysis, and gene regulation.
- 405 *Prog. Mol. Biol. Transl. Sci.* **120**, 1–23 (2013).
- 406 13. Cervera, A. & De la Peña, M. Eukaryotic penelope-like retroelements encode hammerhead ribozyme motifs.
- 407 *Mol. Biol. Evol.* **31**, 2941–2947 (2014).
- 408 14. Martick, M. & Scott, W. G. Tertiary contacts distant from the active site prime a ribozyme for catalysis. *Cell*
- **126**, 309–320 (2006).
- 410 15. Kowalski, J. C. et al. Configuration of the catalytic GIY-YIG domain of intron endonuclease I-TevI:
- 411 coincidence of computational and molecular findings. *Nucleic Acids Res.* 27, 2115–2125 (1999).
- 412 16. Xiong, Y. & Eickbush, T. H. Origin and evolution of retroelements based upon their reverse transcriptase
- 413 sequences. *EMBO J.* **9**, 3353–3362 (1990).
- 414 17. Crooks, G. E., Hon, G., Chandonia, J.-M. & Brenner, S. E. WebLogo: a sequence logo generator. *Genome Res.*
- **415 14**, 1188–1190 (2004).
- 416 18. Sayers, E. W. et al. Database resources of the National Center for Biotechnology Information. Nucleic Acids
- 417 Res. 49, D10–D17 (2021).
- 418 19. Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with Bowtie 2. Nat. Methods 9, 357–359 (2012).
- 419 20. Langmead, B., Wilks, C., Antonescu, V. & Charles, R. Scaling read aligners to hundreds of threads on general-
- 420 purpose processors. *Bioinformatics* **35**, 421–432 (2019).
- 421 21. Katoh, K. & Standley, D. M. MAFFT: iterative refinement and additional methods. *Methods Mol. Biol.* 1079,

- 422 131–146 (2014).
- 423 22. Waterhouse, A. M., Procter, J. B., Martin, D. M. A., Clamp, M. & Barton, G. J. Jalview Version 2--a multiple
- sequence alignment editor and analysis workbench. *Bioinformatics* **25**, 1189–1191 (2009).
- 425 23. Bao, W., Kojima, K. K. & Kohany, O. Repbase Update, a database of repetitive elements in eukaryotic
- 426 genomes. *Mob. DNA* **6**, 11 (2015).
- 427 24. Arkhipova, I. R. Distribution and phylogeny of Penelope-like elements in eukaryotes. Syst. Biol. 55, 875–885
- 428 (2006).
- 429 25. Schmid-Burgk, J. L. et al. Highly Parallel Profiling of Cas9 Variant Specificity. Mol. Cell 78, 794-800.e8
- 430 (2020).

- 431 22. Walguarnery, J.W., Goodman, R.M., Echternacht, A.C. Thermal biology and temperature selection in juvenile
- lizards of co-occurring native and introduced Anolis species. J. Herpetol. 2012; 46 (4): 620–624.
- 433 23. Lorenz, R., Bernhart, S.H., Höner zu Siederdissen, C., Tafer, H., Flamm, C., Stadler, P.F., Hofacker, I.L.
- ViennaRNA Package 2.0. Algorithms for Molecular Biology, 2011; 6 (1): 26.
- 435 24. Sato, K., Akiyama, M., Sakakibara, Y. RNA secondary structure prediction using deep learning with
- thermodynamic integration. *Nat Commun* 2021; 12: 941.
- 437 25. Kabelik, D., Julien, A. R., Ramirez, D., O'Connell, L. A. Social boldness correlates with brain gene expression
- 438 in male green anoles. *Horm. Behav.* 2021; 133:105007.
- 439 26. Arkhipova I.R., Yushenova I.A., Rodriguez F. Giant reverse transcriptase-encoding transposable elements at
- 440 telomeres. Mol. Biol. Evol. 2017; 34(9):2245-2257.