## Accelerated Article Preview

# Template and target site recognition by human LINE-1 in retrotransposition

Akanksha Thawani, Alfredo Jose Florez Ariza, Eva Nogales & Kathleen Collins

1  **Template and target site recognition by human LINE-1 in retrotransposition**
2
3  Akanksha Thawani[1,2,*], Alfredo Jose Florez Ariza[1,3], Eva Nogales[1,2,4,5,*] and Kathleen Collins[1,2,*]
4
5  [1]California Institute for Quantitative Biosciences (QB3), Berkeley, CA 94720, USA
6  [2]Department of Molecular and Cell Biology, University of California Berkeley, Berkeley, CA
7  94720, USA
8  [3]Biophysics Graduate Group, University of California Berkeley, Berkeley, CA, USA, 94720,
9  USA
10 [4]Howard Hughes Medical Institute, Chevy Chase, MD 20815, USA
11 [5]Molecular Biophysics and Integrated Bioimaging Division, Lawrence Berkeley National
12 Laboratory, Berkeley, CA, 94720, USA
13
14 *Corresponding authors. Email: athawani@berkeley.edu, enogales@lbl.gov,
15 kcollins@berkeley.edu

16    **Summary**

17

18    The Long Interspersed Element-1 (L1) retrotransposon has generated nearly one-third of the
19    human genome and serves as an active source of genetic diversity and human disease[1]. L1 spreads
20    via a mechanism termed target-primed reverse transcription (TPRT), in which the encoded enzyme
21    (ORF2p) nicks the target DNA to prime reverse transcription of its own or non-self RNAs[2]. Here,
22    we purified the full-length L1 ORF2p and biochemically reconstituted robust TPRT with template
23    RNA and target site DNA. We report cryo-electron microscopy structures of the human L1 ORF2p
24    bound to structured template RNAs and initiating cDNA synthesis. The template polyadenosine
25    tract is recognized in a sequence-specific manner by five distinct domains. Among them, a novel
26    RNA-binding domain bends the template backbone to allow engagement of an RNA hairpin stem
27    with the L1 ORF2p C-terminal segment. In addition, structure and biochemical reconstitutions
28    demonstrate a surprising target-site requirement: L1 ORF2p relies on upstream single-stranded
29    DNA to position adjacent duplex in the endonuclease active site for nicking of the longer DNA
30    strand, with a single nick generating a staggered DNA break. Our work provides key insights into
31    the mechanism of ongoing transposition in the human genome and informs the engineering of
32    retrotransposon proteins for gene therapy.

33    Non-long terminal repeat (non-LTR) retrotransposons are mobile genetic elements in the human
34    genome that are recognized as drivers of genome expansion and evolution[1]. The human genome
35    has one autonomously active retrotransposon from the Long Interspersed Element (LINE) family.
36    Human LINE-1 (L1) is present in an estimated 80-100 transposition competent copies[3] that are
37    sources of genetic diversity and ongoing somatic mosaicism[4], and contribute to more than 100
38    known human disease cases[5,6]. Bicistronic L1 encodes an ORF1 protein that binds RNA[7], and an
39    enzymatic ORF2 protein that has endonuclease (EN) and reverse-transcriptase (RT) activities[8,9]
40    (Fig. 1A). New L1 insertions initiate by target-primed reverse transcription (TPRT), in which
41    target site nicking creates a primer for cDNA synthesis directly into the genome[2,8,10,11]. L1 ORF2p
42    has generated more than 30% of the human genome via transposition and pseudogene synthesis[12].
43    Current efforts that seek to limit human disease by controlling L1 mobility[13], and to exploit non-
44    LTR retrotransposons and other RTs for genome engineering[14–17], provide an increasingly
45    compelling demand for mechanistic understanding of TPRT and stable cDNA incorporation into
46    the genome. However, much remains unknown, in large part because of experimental difficulties
47    in L1 ORF2p biochemical reconstitution and structural analyses.
48
49    The purification of active L1 ORF2p has been challenging due to the scarcity of L1
50    ribonucleoproteins (RNPs) in cells, as well as the heterogeneous association of L1 ORF2p with
51    L1 and other RNAs and many directly or indirectly interacting proteins[18–21]. Consequently,
52    biochemical assays for L1 activity have been limited, most relying on the cellular assembly of an
53    L1 ORF2p RNP[22,23]. Among the questions that remain to be addressed, understanding how L1
54    ORF2p recognizes template RNAs to initiate TPRT is particularly critical (Fig. 1B). The prevailing
55    model, termed *cis*-preference, proposes that L1 ORF2p co-translationally engages the
56    polyadenosine (polyA) tail of its encoding transcript to promote selective binding and cDNA
57    insertion of the L1 mRNA[24–26]. Yet, the most abundant insertions mediated by L1 ORF2p are the
58    non-autonomous Short Interspersed Nuclear Elements (SINEs), such as Alu SINEs[24,27,28]. In
59    another outstanding question, how the endonuclease (EN) domain of L1 ORF2p selects target sites
60    to nick for TPRT initiation, beyond the short consensus motif TTTTT/AA[8,29–32], remains poorly
61    understood (Fig. 1B). Robust biochemical reconstitutions and structural studies with the purified
62    L1 ORF2p are needed to understand the mechanisms of nucleic acid recognition for TPRT.
63
64    **Reconstitution of L1 ORF2p-mediated TPRT**
65    We expressed the full-length L1 ORF2p in insect cells and purified it to relative homogeneity
66    (Extended Data Fig. 1a-b). With an optimal target DNA structure (see below) containing a single
67    TTTTT/AA consensus for genomic L1 insertions[8,29–32], efficient nicking occurred at the intended
68    site, evident by the formation of a 16-nucleotide (nt) nicked product, and TPRT product was
69    synthesized by nick-primed reverse transcription of template RNA (Fig. 1c). We compared
70    template RNAs that are established native substrates of L1 ORF2p, including the L1 3'UTR and
71    Alu RNAs, each with a 3' 25A tail (Supplementary Table 1). All Alu RNAs, including the
72    evolutionarily youngest AluY RNA[28], a resurrected AluJ RNA[33,34], and a left-half monomer of the
73    Alu RNA tandem repeat sufficient for genome insertion[33] (AluJ half, AJh), were efficiently reverse
74    transcribed from the nicked primer (Fig. 1c, lanes AY, AJ and AJh). In contrast, TPRT of the L1
75    3'UTR RNA resulted in a lower amount of product synthesis with products predominantly
76    migrating faster than expected for full-length cDNA (Fig. 1c, lane L1). L1 3'UTR template lacking
77    nt 1-78 that form G-quadruplex[35] gave the expected cDNA length, which matched the length of
78    the shorter products from the full-length L1 3'UTR template (Fig. 1c, lanes L1 and L1Δ). Neither

3

79   L1 3'UTR template supported as much TPRT as Alu RNA (AJh), suggesting that the L1 3'UTR
80   is a suboptimal template for L1 ORF2. Using the optimal AJh template, we verified that neither a
81   control retroviral RT from Moloney Murine Leukemia Virus (M-MLV RT) nor an EN-dead L1
82   ORF2p mutant had nicking or TPRT activities (Fig. 1d and Extended Data Fig. 1c), yet both
83   showed robust RT activity as assayed by primer-extension on annealed RNA-DNA duplex
84   (Extended Data Fig. 1d-e). In contrast, an RT-dead L1 ORF2p retained target-site nicking but no
85   RT or TPRT activity (Fig. 1d and Extended Data Fig. 1d-e). These controls validate our direct
86   readout of robust L1 ORF2p-mediated TPRT activity, bypassing the PCR-based amplification
87   required previously[10].
88
89   **Structure of template RNA bound L1 ORF2p**
90   We sought to capture the structure of L1 ORF2p. While our initial attempts at cryo-EM
91   reconstruction of L1 ORF2p without nucleic acids were unsuccessful, we were able to capture L1
92   ORF2p engaged with RNA. We imaged L1 ORF2p bound to Alu AJh RNA with a poly-thymidine
93   (polyT) primer base-paired to its 3' end to mimic the initiation of cDNA synthesis (Fig. 1e). In the
94   resulting 4.4 Å resolution density map, we could place the predicted AlphaFold model of human
95   L1 ORF2p[36] and further identify extra density consistent with the Alu RNA stem-loop bound on
96   one side of the protein and its 3' tail in the L1 ORF2p RT core in an orientation that is topologically
97   compatible with the co-binding of the Alu RNA partner, the SRP9/14 heterodimer[33] (Fig. 1e and
98   Extended Data Fig. 2a, c). However, the Alu RNP map suffered from preferred orientation issues
99   and did not have the resolution to visualize amino acid side chains (Extended Data Fig. 2).
100  We improved the quality and resolution of our density map when we used an L1 ORF2p
101  complex with a synthetic RNA template mimicking Alu RNA features (Fig. 1f, right), harboring a
102  5' stem-loop and a 3' single-stranded region of sufficient length to span the distance between the
103  Alu RNA stem-loop position and the active site of L1 ORF2p seen in our 4.4 Å RNP map. Because
104  cellular assays concur that L1 templates require a 3' polyA tract[37,38], we used adenosine (A) in the
105  single-stranded region. We halted elongation after 5 base-pairs (bp) of cDNA synthesis with
106  dideoxyguanosine triphosphate (ddGTP) replacing dGTP (Fig. 1f, right). Using this sample, we
107  obtained the cryo-EM structure of the RNP in a paused elongation state at an overall resolution of
108  3.2 Å (Fig. 1f, Extended Data Figs. 3-4 and Extended Data Table 1). This resolution allowed us to
109  model the entire protein chain and the individual nucleotides, including a dTTP bound as
110  nucleotide substrate but unable to join the cDNA 3' end (Fig. 1g and Extended Data Fig. 5a-b).
111  Only 8 nt of template RNA near the loop and 3 bp of RNA-DNA duplex farthest from the active
112  site could not be modeled.
113  The L1 ORF2p RT core consists of the palm and fingers (altogether, RT domain) in the right-
114  hand architecture shared by many polymerases, followed by the Thumb domain and preceded by
115  an N-terminal extension (NTE) domain previously noted in L1 ORF2p as "Z-domain"[39,40], all
116  shared with prokaryotic and eukaryotic retrotransposon RTs[41] (Fig. 1a, f, g). The RT and Thumb
117  domains cradle the RNA-DNA duplex emerging from the active site. Preceding the NTE, L1
118  ORF2p has an N-terminal apurinic/apyrimidinic (AP) EN domain fold[42], connected to the rest of
119  the protein through a folded domain incorporating the "cryptic motif"[39] and hereafter designated
120  EN linker, which packs against an adjacent portion of the NTE. The 209 amino acid L1 C-terminal
121  segment (CTS), together with the NTE and EN linker domains, create an extended surface of
122  contacts with the polyA tract of the template RNA proximal to the active site (Fig. 1f, g;
123  summarized in Fig. 2a). The region between the CTS and Thumb, which we labeled as a previously
124  unidentified RNA binding domain (RBD, Fig. 1a), contacts both RT-bound template RNA and

125 peripheral RNA stem-loop (Fig. 1f, g and Extended Data Fig. 5c; summarized in Fig. 2a). The
126 array of protein-RNA interactions bends the template RNA to follow an L-shaped architecture
127 (Fig. 1f, g). Overall, our structure reveals a previously unknown topology and indicates
128 biochemical roles for the different L1 ORF2p domains.
129
## Features of the catalytic core
131 L1 ORF2p RT activity is supported by numerous side chain interactions with nucleic acids. Of the
132 traceable 11 bp of RNA-DNA duplex, 9 bp are almost fully enclosed, predominantly by
133 interactions with the RT, Thumb, and RBD domains (Fig. 2a and Extended Data Fig. 5). The
134 incoming dTTP and ddG-13 at the primer 3' end are positioned by the canonical FADD active site
135 motif and by the conserved aromatic residues Phe566 and Phe605 (Extended Data Fig. 5a). The
136 incoming dTTP hydrogen bonds with three RT domain residues, including the Arg531 side chain
137 (Extended Data Fig. 5b). These contacts parallel the configuration of a Group II intron RT active
138 site[43–46]. The RNA strand of the heteroduplex exiting the active site contacts residues in the NTE
139 and RT domains, and it also contacts the RBD domain not shared with Group II intron RTs
140 (Extended Data Fig. 5c). The cDNA strand has fewer contacts: an electrostatic interaction between
141 the DNA backbone and the side chain of Arg375 in the NTE domain, and several hydrophobic
142 contacts with sugars by Thumb and RT domain residues (Extended Data Fig. 5d). All contacts to
143 nucleic acids in the RT core are sequence non-specific.
144
## Single-stranded RNA recognition
146 Side chains across several domains in the protein define the surface for recognition of the 15-nt
147 single-stranded polyA tract template (Fig. 2a-b). The EN linker, NTE, RT and Thumb domains
148 engage the polyA tract proximal to the active site, whereas the CTS domain interacts with the
149 polyA tract predominantly adjacent to the stem-loop (Fig. 2a-b). This architecture suggests a
150 "threshold" model in which a substantial length of 3' polyA would be required for template binding
151 and threading into the active site. To define the polyA length for optimal TPRT, we designed and
152 purified AJh RNAs with variable polyA tail length and 3' tail sequences (Supplementary Table 1
153 and Extended Data Fig. 5e) and used them as templates for TPRT by L1 ORF2p. Templates with
154 75A, 50A, 25A and 20A were used efficiently, whereas shorter A-tails of 15A, 10A and 5A
155 produced much less or no TPRT product (Fig. 2c). These results agree with our structure-based
156 prediction: a template with 20A, allowing 5A for base-pairing with the nicked primer and at least
157 15 nt of single-stranded polyA, can be efficiently used for TPRT initiation, while for a template
158 with 25A product synthesis reaches the same level obtained using templates with longer A tracts
159 (Fig. 2c, bar graph). Notably, AJh RNA with either 75A or 50A produced a heterogeneous size
160 distribution of TPRT products, with 75A displaying a distinct skew toward lower length of cDNA
161 product than the expected 200-nt full-length cDNA (Fig. 2c). This heterogeneity suggests that the
162 longer polyA tracts exceed the length of single-stranded RNA recognized by L1 ORF2p. Overall,
163 our findings agree with the polyA tail length shown to be required for *in vivo* mobility of L1[37] or
164 Alu SINEs[38].
165 Strikingly, we observed side chain interactions with A-bases distributed across the entire
166 length of polyA tract (Fig. 2a), including contacts that specifically recognize the adenine base (Fig.
167 2d-f). The A-60 base forms adenine-specific hydrogen bonds with Arg385 and Asn388 of the NTE
168 domain, as well as a hydrophobic contact with Ile517 from the RT domain (Fig. 2d). The A-57
169 base forms a hydrogen-bond with Lys1236 from the CTS domain and stacks against the Trp365
170 side chain from the EN linker domain (Fig. 2e). The A-55 base forms hydrogen-bonds with Asn371

5

171 and Cys804 from the NTE and Thumb domains, respectively, and is caged in a hydrophobic pocket
172 formed by leucines from the NTE domain and Phe366 from the EN linker (Fig. 2f). The CTS
173 domain also contributes to adenine-specific recognition (see below). To investigate the
174 dependence of TPRT on the single-stranded polyA sequence, we generated AJh-based RNA
175 templates terminating in 25N (a 25 nt sequence with mixed base composition) or 20N with 3' 5A
176 to retain template-primer base pairing (see Fig. 2 legend, Supplementary Table 1, and Extended
177 Data Fig. 5e). Neither template supported L1 ORF2p's TPRT activity (Fig. 2c). Further intrigued
178 by the large number of hydrogen bonds with the polyA tract, we created mutant L1 ORF2p with
179 alanine substitutions for all 8 side chains that make base contacts to single-stranded RNA (Fig.
180 2a). When assayed for TPRT activity, the L1 ORF2p mutant for single-stranded RNA base
181 interactions (Δss) showed distinctly reduced TPRT while retaining significant EN activity and RT
182 activity when assayed by primer extension (Fig. 2g and Extended Data Fig. 1d-e). We suggest that
183 these contacts contribute to a conformation of L1 ORF2p poised for cDNA synthesis.
184
185 **Novel roles for the C-terminal segment**
186 The template RNA stem-loop and polyA region distal to the RT active site are predominantly
187 engaged by the CTS domain (Fig. 2a-b). Adjacent to the stem-loop, the A-49 base makes adenine-
188 specific hydrogen bonds with Lys1107 and His1113 in the CTS domain (Fig. 3a). However, CTS
189 domain interactions with the RNA are predominantly hydrophobic, without much sequence
190 specificity, in agreement with previous work[47] (Fig. 2a). This involves aromatic side chains of
191 Trp1208, Trp1131 and His1113 that present stacking opportunities for the RNA bases (Fig. 3b).
192 Strikingly, our structure captures the CTS domain forcing apart the RNA stem-loop strands at the
193 base of the stem through the steric barrier defined by an alpha helix (hereafter, termed "insertion
194 helix"), which forks the RNA stem. In the structure, the first three stem base-pairs are splayed
195 apart (Fig. 3c) concurrent with Ile1121 and Ile1122 of the insertion helix forming hydrophobic
196 interactions with the splayed bases G-1 and C-46 (Ile1121 and Ile1122), G-45, and U-44 (Ile1122)
197 (Fig. 3c). These interactions induce a distortion in RNA conformation away from the canonical A-
198 form helix at the base of the stem (Extended Data Fig. 5f).
199 To investigate the role of the insertion helix and the entire CTS domain overall, we generated
200 mutant L1ORF2 proteins with the entire CTS domain deleted (ΔCTS) or with the insertion helix
201 replaced by negatively charged residues (ΔIH). Both showed notably reduced TPRT activity, and
202 the ΔCTS protein was further compromised for target-site nicking activity (Fig. 3d), suggesting a
203 role of the CTS domain beyond interacting with and unwinding the template RNA. To validate the
204 structural integrity of L1 ORF2p mutants, particularly without the CTS domain, we verified that
205 the mutant proteins had similar or greater than wild-type RT activity in our primer extension assays
206 (Extended Data Fig. 1d-e).
207 The entire template RNA stem is nestled into a positively charged surface composed of the
208 CTS, RBD, and Thumb domains (Fig. 3e and Extended Data Fig. 5g), which engage but do not
209 contort the RNA stem aside from the stem's base (Extended Data Fig. 5g). To investigate the
210 significance of the RNA stem-loop for L1 ORF2p's TPRT activity, we generated AJh template
211 RNA variants that differ from native stem structure by increased (AJhm) or decreased (AJh-uf)
212 base-pairing (Supplementary Table 1). While removing mismatches did not increase TPRT
213 product, unpairing the stem-loop with mismatches resulted in a modest decrease in TPRT
214 efficiency and a dramatic increase in the heterogeneity of TPRT product lengths, in which shorter
215 than full-length cDNA products were generated (Fig. 3f). These results suggest that the stem-loop
216 could contribute to defining where TPRT initiates within the template RNA.

6

217    To explore if other RT families share a CTS-like domain with a similar function, we searched
218  for a homologous structure across the evolutionary tree. Our structure-based search revealed a
219  distant relationship to nucleic acid-interacting motifs in the *Bombyx mori* R2 retrotransposon
220  protein[48,49] and in the human telomerase catalytic core[50] (Extended Data Fig. 6a-b). However, it
221  remains to be determined if these partial CTS-like motifs share the same function as the CTS
222  domain in L1 ORF2p. In contrast, primary sequence comparison found homology only within the
223  L1 family. L1 enzymes from fish to human show conservation of the overall hydrophobic content
224  of the CTS domain insertion helix, with L1 ORF2p Ile1122 being replaced only by another
225  hydrophobic residue (Extended Data Fig. 6c).
226
227  **Target site architecture for TPRT**
228  To investigate what structural features may influence recognition and cleavage of target DNA, we
229  superimposed the structure of the L1 EN domain co-crystallized with DNA duplex[51] onto our L1
230  ORF2p RNP structure (Fig. 4a). We observed that the consensus cleavage site (TTTTT/AA) is
231  accessible to the EN domain when located close to the 5' end of the DNA duplex (Fig. 4a, top
232  panel). Surprisingly, adding extra DNA base-pairs upstream (5' of TTTTT) of the consensus
233  cleavage site introduced a steric clash with the L1 ORF2p CTS domain (Fig. 4a, bottom panel).
234  We predicted that as little as ~10 upstream bp could severely inhibit EN domain engagement with
235  the target site. To test this structure-based prediction, we designed DNA duplexes with the
236  consensus cleavage site positioned at different distances from the edge of the base-paired duplex.
237  TPRT assays revealed a drastic inhibition of EN nicking activity and subsequent TPRT from an
238  upstream duplex region as short as 11 bp, with optimal EN nicking and TPRT for an upstream
239  duplex of ~7-9 bp (Fig. 4b-c). Off-target EN nicking (not at the consensus site) was common for
240  non-optimal target-site duplexes and occurred between pyrimidine and purine nt, in agreement
241  with non-consensus cleavage in cells[8,30] (Extended Data Fig. 7). Consistent with what would be
242  expected from the structure, deletion of L1 ORF2p's CTS domain (ΔCTS mutant) enabled nicking
243  of DNA substrates with an upstream duplex region greater than 13 bp (Extended Data Fig. 8).
244  Nonetheless, the ΔCTS mutant did not nick all target sites equally (Extended Data Fig. 8),
245  indicating that there are other determinants of efficient nicking beyond the minimal consensus
246  TTTTT/AA.
247    L1-mediated TPRT in cells is coupled with DNA replication, with preferential EN nicking
248  of the lagging-strand template[30,52]. We therefore hypothesized that an optimal target site could
249  have a single-stranded 5' overhang upstream of the duplex region containing the EN consensus
250  sequence, a design that mimics the lagging strand template with an Okazaki fragment primer. To
251  test this possibility, we compared EN nicking and TPRT activity using DNA duplexes with
252  different 5' overhang lengths upstream of the consensus target site. We found that the presence of
253  an overhang was strongly stimulatory, with some influence from the overhang nucleotide
254  composition (Fig. 4d). Remarkably, increasing the upstream overhang length from 9 to 27 nt gave
255  a tremendous stimulation of nicking efficiency, with two-thirds of the target DNA harboring the
256  longest overhang converted into on-target nicked product (Fig. 4e). Consequently, a 6-fold
257  increase in the TPRT product was also observed with increasing overhang length from 9 to 27 nt
258  (Fig. 4e). We conclude that L1 target sites are partial duplex structures with a long single-stranded
259  5' overhang, with the EN cleavage site is positioned on duplex DNA near the single-strand/duplex
260  transition (Fig. 4f). This unanticipated structure of optimal target-site DNA architecture supports
261  efficient TPRT by L1 ORF2p (Fig. 4f) and explains why previous reconstitutions resulted in low
262  TPRT efficiency[32]. Our results have profound implications for the understanding of L1 and Alu

7

263 mobility in the human genome.

264

265 **Discussion**

266

267 **Adaptation for nucleic acid recognition**

268 Phylogenetic characterization suggests that a prokaryotic mobile Group IIB intron protein gave
269 rise to eukaryotic single-ORF retrotransposons with a domain architecture like the R2
270 retrotransposon, which in turn spawned two-ORF retrotransposons like those in the L1 family[41].
271 We compared the L1 ORF2p structure and substrate engagement to that of its ancestral Group IIB
272 intron from *Thermosynechococcus elongatus*[45], and with the recently reported cryo-EM structure
273 of non-LTR retrotransposon R2 from *Bombyx mori* (R2Bm)[48,49]. Template RNA binds to L1
274 ORF2p with similar topology to that of Group IIB intron RT binding to intron RNA and that of
275 R2Bm binding to target-site DNA upstream of the nick site (Extended Data Fig. 9). However, and
276 despite their evolutionary relationship, our work highlights major differences between the TPRT
277 strategies of L1 ORF2p and R2Bm proteins. First, while the CTS-like domain of R2Bm melts
278 duplex DNA (Extended Data Figs. 9b), the analogous L1 ORF2p CTS domain can bind and
279 facilitate unwinding of RNA. Second, the EN domains are in distinct positions relative to their RT
280 cores (Extended Data Fig. 9b). Third, whereas R2Bm engages a long duplex DNA with sequence-
281 specific DNA binding domains, L1 ORF2p has a largely sequence-independent target-site
282 association that relies on limited duplex length 5' of the target site and a single-stranded DNA
283 overhang.

284

285 **Implications for L1 and SINE lifecycles**

286 Together, our structural and biochemical studies reveal novel insights into the retrotransposition
287 of L1 and SINEs and offer mechanistic rationale for the observed biological properties of L1-
288 mediated genomic insertions (Fig. 4f). First, the extensive surface of L1 ORF2p dedicated to
289 binding single-stranded polyA with adenine-specific contacts explains the loss of transposition of
290 RNAs that lack the long polyA tract of L1 RNAs or the genome-encoded polyA tract of
291 SINEs[27,37,38]. Second, template anchoring to L1 ORF2p by a stem-loop structure can explain why
292 Alu RNAs outcompete the L1 3'UTR for L1 ORF2p binding[24–27], even if both are associated to
293 the same ribosome, because the L1 3'UTR lacks a distinct stem-loop structure. Third, the long
294 single-stranded DNA upstream of the EN cleavage site required for L1 ORF2p's activity helps
295 explain the preference for nicking the lagging DNA strand template at replication forks, and why
296 L1 ORF2p's chromatin engagement and transposition are coupled with DNA replication[30,52].

297 A complete L1 retrotransposition cycle is assumed to require nicking of the second strand
298 of a target site prior to the second-strand synthesis that generates a double-stranded copy of L1 or
299 SINE. The L1 ORF2p target-site architecture, where first-strand cleavage occurs at a limited length
300 of duplex away from a single-strand/duplex transition on the 5'-overhang strand, produces a nick
301 only ~10 bp away from the 5' overhang. This ~10 bp of duplex is prone to dissociation, eliminating
302 the need for second-strand nicking (Fig. 4f). The target-site DNA architecture also accounts for
303 sequence duplication surrounding the new L1 insertion, although the observed target-site
304 duplication lengths[8,29] would also depend on other factors, e.g. the extent of unpairing of upstream
305 duplex by Replication Protein A from the adjacent single-stranded DNA. L1 ORF2p interaction
306 with factors such as PCNA could facilitate target-site selection[18,19]. The predicted PCNA-
307 interacting protein (PIP) box motif in L1 ORF2p[18] is located on a highly accessible α-helix of the
308 NTE domain (Extended Data Fig. 10a), and we found that addition of PCNA gives a modest

8

309  increase in TPRT activity in our biochemical assays, despite the short linear duplex (Extended
310  Data Fig. 10b). Overall, the combination of L1 ORF2p's target-site structure specificity and its
311  interaction with PCNA can explain preferential insertion into the lagging-strand template behind
312  a replication fork, where there would be an intact leading-strand duplex to support DNA break
313  repair.
314

329  **Author contributions**
330  A.T., E.N. and K.C. conceived the project. A.T. collected and analyzed the electron microscopy
331  data, performed manual model building and refinements, analyzed the structures and performed
332  biochemical assays. A.J.F.A. obtained the first L1 ORF2p structure and provided advice on single
333  particle analysis and model building. A.T. wrote the paper with input and revisions from all
334  authors.
335

336  **Data Availability**
337  The 3.2 Å cryo-EM map reported in this work is deposited under EMD-42637 in the Electron
338  Microscopy Data Bank and the corresponding atomic model under PDB 8UW3 on the Protein
339  Data Bank. All other datasets generated and analyzed during the current study are available from
340  the corresponding authors on request.

**References**

341
342　1.　Kazazian, H. H. & Moran, J. V. Mobile DNA in Health and Disease. *N Engl J Med* **377**,
343　　　361–370 (2017).
344　2.　Han, J. S. Non-long terminal repeat (non-LTR) retrotransposons: mechanisms, recent
345　　　developments, and unanswered questions. *Mob DNA* **1**, 15 (2010).
346　3.　Brouha, B. *et al.* Hot L1s account for the bulk of retrotransposition in the human population.
347　　　*Proc Natl Acad Sci U S A* **100**, 5280–5285 (2003).
348　4.　Baillie, J. K. *et al.* Somatic retrotransposition alters the genetic landscape of the human brain.
349　　　*Nature* **479**, 534–537 (2011).
350　5.　Hancks, D. C. & Kazazian, H. H. Roles for retrotransposon insertions in human disease. *Mob*
351　　　*DNA* **7**, 9 (2016).
352　6.　Burns, K. H. Repetitive DNA in disease. *Science* **376**, 353–354 (2022).
353　7.　Martin, S. L. & Bushman, F. D. Nucleic acid chaperone activity of the ORF1 protein from
354　　　the mouse LINE-1 retrotransposon. *Mol Cell Biol* **21**, 467–475 (2001).
355　8.　Feng, Q., Moran, J. V., Kazazian, H. H. & Boeke, J. D. Human L1 retrotransposon encodes a
356　　　conserved endonuclease required for retrotransposition. *Cell* **87**, 905–916 (1996).
357　9.　Mathias, S. L., Scott, A. F., Kazazian, H. H., Boeke, J. D. & Gabriel, A. Reverse
358　　　transcriptase encoded by a human transposable element. *Science* **254**, 1808–1810 (1991).
359　10.　Cost, G. J., Feng, Q., Jacquier, A. & Boeke, J. D. Human L1 element target-primed reverse
360　　　transcription in vitro. *EMBO J* **21**, 5899–5910 (2002).
361　11.　Moran, J. V. *et al.* High frequency retrotransposition in cultured mammalian cells. *Cell* **87**,
362　　　917–927 (1996).
363　12.　Lander, E. S. *et al.* Initial sequencing and analysis of the human genome. *Nature* **409**, 860–
364　　　921 (2001).
365　13.　Zhang, X., Zhang, R. & Yu, J. New Understanding of the Relevant Role of LINE-1
366　　　Retrotransposition in Human Disease and Immune Modulation. *Front Cell Dev Biol* **8**, 657
367　　　(2020).
368　14.　Chen, P. J. & Liu, D. R. Prime editing for precise and highly versatile genome manipulation.
369　　　*Nat Rev Genet* **24**, 161–177 (2023).
370　15.　Manoj, F., Tai, L. W., Wang, K. S. M. & Kuhlman, T. E. Targeted insertion of large genetic
371　　　payloads using cas directed LINE-1 reverse transcriptase. *Sci Rep* **11**, 23625 (2021).
372　16.　Zhao, B., Chen, S.-A. A., Lee, J. & Fraser, H. B. Bacterial Retrons Enable Precise Gene
373　　　Editing in Human Cells. *CRISPR J* **5**, 31–39 (2022).
374　17.　Lopez, S. C., Crawford, K. D., Lear, S. K., Bhattarai-Kline, S. & Shipman, S. L. Precise
375　　　genome editing across kingdoms of life using retron-derived DNA. *Nat Chem Biol* **18**, 199–
376　　　206 (2022).
377　18.　Taylor, M. S. *et al.* Affinity proteomics reveals human host factors implicated in discrete
378　　　stages of LINE-1 retrotransposition. *Cell* **155**, 1034–1048 (2013).
379　19.　Taylor, M. S. *et al.* Dissection of affinity captured LINE-1 macromolecular complexes. *Elife*
380　　　**7**, (2018).
381　20.　Goodier, J. L., Cheung, L. E. & Kazazian, H. H. Mapping the LINE1 ORF1 protein
382　　　interactome reveals associated inhibitors of human retrotransposition. *Nucleic Acids Res* **41**,
383　　　7401–7419 (2013).
384　21.　Moldovan, J. B. & Moran, J. V. The Zinc-Finger Antiviral Protein ZAP Inhibits LINE and
385　　　Alu Retrotransposition. *PLoS Genet* **11**, e1005121 (2015).

386  22. Kulpa, D. A. & Moran, J. V. Cis-preferential LINE-1 reverse transcriptase activity in
387      ribonucleoprotein particles. *Nat Struct Mol Biol* **13**, 655–660 (2006).
388  23. Viollet, S., Doucet, A. J. & Cristofari, G. Biochemical Approaches to Study LINE-1 Reverse
389      Transcriptase Activity In Vitro. *Methods Mol Biol* **1400**, 357–376 (2016).
390  24. Boeke, J. D. LINEs and Alus--the polyA connection. *Nat Genet* **16**, 6–7 (1997).
391  25. Wei, W. *et al.* Human L1 retrotransposition: cis preference versus trans complementation.
392      *Mol Cell Biol* **21**, 1429–1439 (2001).
393  26. Esnault, C., Maestre, J. & Heidmann, T. Human LINE retrotransposons generate processed
394      pseudogenes. *Nat Genet* **24**, 363–367 (2000).
395  27. Dewannieux, M., Esnault, C. & Heidmann, T. LINE-mediated retrotransposition of marked
396      Alu sequences. *Nat Genet* **35**, 41–48 (2003).
397  28. Deininger, P. Alu elements: know the SINEs. *Genome Biol* **12**, 236 (2011).
398  29. Jurka, J. Sequence patterns indicate an enzymatic involvement in integration of mammalian
399      retroposons. *Proc Natl Acad Sci U S A* **94**, 1872–1877 (1997).
400  30. Flasch, D. A. *et al.* Genome-wide de novo L1 Retrotransposition Connects Endonuclease
401      Activity with Replication. *Cell* **177**, 837-851.e28 (2019).
402  31. Sultana, T. *et al.* The Landscape of L1 Retrotransposons in the Human Genome Is Shaped by
403      Pre-insertion Sequence Biases and Post-insertion Selection. *Mol Cell* **74**, 555-570.e7 (2019).
404  32. Cost, G. J. & Boeke, J. D. Targeting of human retrotransposon integration is directed by the
405      specificity of the L1 endonuclease for regions of unusual DNA structure. *Biochemistry* **37**,
406      18081–18093 (1998).
407  33. Ahl, V., Keller, H., Schmidt, S. & Weichenrieder, O. Retrotransposition and Crystal
408      Structure of an Alu RNP in the Ribosome-Stalling Conformation. *Mol Cell* **60**, 715–727
409      (2015).
410  34. Bennett, E. A. *et al.* Active Alu retrotransposons in the human genome. *Genome Res* **18**,
411      1875–1883 (2008).
412  35. Sahakyan, A. B., Murat, P., Mayer, C. & Balasubramanian, S. G-quadruplex structures
413      within the 3' UTR of LINE-1 elements stimulate retrotransposition. *Nat Struct Mol Biol* **24**,
414      243–247 (2017).
415  36. Jumper, J. *et al.* Highly accurate protein structure prediction with AlphaFold. *Nature* **596**,
416      583–589 (2021).
417  37. Doucet, A. J., Wilusz, J. E., Miyoshi, T., Liu, Y. & Moran, J. V. A 3' Poly(A) Tract Is
418      Required for LINE-1 Retrotransposition. *Mol Cell* **60**, 728–741 (2015).
419  38. Dewannieux, M. & Heidmann, T. Role of poly(A) tail length in Alu retrotransposition.
420      *Genomics* **86**, 378–381 (2005).
421  39. Adney, E. M. *et al.* Comprehensive Scanning Mutagenesis of Human Retrotransposon LINE-
422      1 Identifies Motifs Essential for Function. *Genetics* **213**, 1401–1414 (2019).
423  40. Clements, A. P. & Singer, M. F. The human LINE-1 reverse transcriptase:effect of deletions
424      outside the common reverse transcriptase domain. *Nucleic Acids Res* **26**, 3528–3535 (1998).
425  41. Xiong, Y. & Eickbush, T. H. Similarity of reverse transcriptase-like sequences of viruses,
426      transposable elements, and mitochondrial introns. *Mol Biol Evol* **5**, 675–690 (1988).
427  42. Weichenrieder, O., Repanas, K. & Perrakis, A. Crystal structure of the targeting
428      endonuclease of the human LINE-1 retrotransposon. *Structure* **12**, 975–986 (2004).
429  43. Stamos, J. L., Lentzsch, A. M. & Lambowitz, A. M. Structure of a Thermostable Group II
430      Intron Reverse Transcriptase with Template-Primer and Its Functional and Evolutionary
431      Implications. *Mol Cell* **68**, 926-939.e4 (2017).

432   44. Zhao, C. & Pyle, A. M. Crystal structures of a group II intron maturase reveal a missing link
433        in spliceosome evolution. *Nat Struct Mol Biol* **23**, 558–565 (2016).
434   45. Haack, D. B. *et al.* Cryo-EM Structures of a Group II Intron Reverse Splicing into DNA.
435        *Cell* **178**, 612-623.e12 (2019).
436   46. Chung, K. *et al.* Structures of a mobile intron retroelement poised to attack its structured
437        DNA target. *Science* **378**, 627–634 (2022).
438   47. Piskareva, O., Ernst, C., Higgins, N. & Schmatchenko, V. The carboxy-terminal segment of
439        the human LINE-1 ORF2 protein is involved in RNA binding. *FEBS Open Bio* **3**, 433–437
440        (2013).
441   48. Wilkinson, M. E., Frangieh, C. J., Macrae, R. K. & Zhang, F. Structure of the R2 non-LTR
442        retrotransposon initiating target-primed reverse transcription. *Science* eadg7883 (2023)
443        doi:10.1126/science.adg7883.
444   49. Deng, P. *et al.* Structural RNA components supervise the sequential DNA cleavage in R2
445        retrotransposon. *Cell* **186**, 2865-2879.e20 (2023).
446   50. Ghanim, G. E. *et al.* Structure of human telomerase holoenzyme with bound telomeric DNA.
447        *Nature* **593**, 449–453 (2021).
448   51. Miller, I. *et al.* Structural dissection of sequence recognition and catalytic mechanism of
449        human LINE-1 endonuclease. *Nucleic Acids Res* **49**, 11350–11366 (2021).
450   52. Mita, P. *et al.* LINE-1 protein localization and functional dynamics during the cell cycle.
451        *Elife* **7**, e30058 (2018).
452

12

453     **Methods**
454     
455     **Protein Expression and Purification**
456     Full length human LINE1 ORF2 was synthesized (Genscript) and cloned into pFastbac1 vector
457     with His and ZZ-tags. The L1 ORF2p mutation and truncation constructs consisted of the
458     following residues: RT mutant (D702A, D703A), EN mutant (D145A, Y226K[51]), ssRNA (Δss)
459     binding mutant (N371A, R385A, N388A, C804A, R855A, K1107A, H1113A, K1236A),
460     ΔInsertion helix (from V1117 to K1124 mutated to EDDDDDE), ΔCTS (missing residues 1067-
461     1275). All constructs were fully sequenced. The plasmids were transformed into DH10Bac *E. coli*
462     strain to produce bacmids and transfected into Sf9 cells using Bac-to-Bac system (Invitrogen).
463     Three rounds of baculoviral expansion were performed and used for infection of Sf9 cells and used
464     for infection of Sf9 or High Five cells. The insect cells were lysed with sonication and the lysate
465     was clarified by centrifugation at 40,000 rpm in Ti45 rotor (Beckman Coulter) for 30-45 minutes.
466     The proteins were purified with the IgG Sepharose resin (Cytiva), eluted by cleavage with TEV
467     protease, followed by a Heparin column (Cytiva) and finally via gel filtration with a Superdex 200
468     10/300 column (Cytiva). Peak elution fractions were analyzed on SDS PAGE, concentrated, flash
469     frozen in liquid and stored in -80°C. Protein concentrations were determined by analyzing with
470     Bradford reagent (Biorad) against a known Bovine Serum Albumin standard. Mass spectrometry
471     was performed to verify that the full length L1 ORF2p protein was obtained.
472             Human PCNA with N-terminal His-tag was expressed in E. coli (Rosetta2 strain) and
473     purified with NiNTA-affinity (Qiagen), followed by HiTrapQ column (Cytiva) and finally via gel
474     filtration on Superdex 200 10/300 column (Cytiva). Peak elution fractions were analyzed on SDS
475     PAGE, concentrated, flash frozen in liquid nitrogen and stored in -80°C. Protein concentrations
476     were determined by analyzing with Bradford reagent (Biorad) against a known Bovine Serum
477     Albumin standard.
478     
479     **RNA Transcription and Purification**
480     The sequence of the youngest SINE element, AluY, was PCR amplified from a parent vector[53] to
481     include the T7 RNA polymerase promoter followed by a 25A sequence. Full-length AluJ SINE
482     element sequence[33] was synthesized (IDT) and PCR amplified to include the T7 RNA polymerase
483     promoter followed by 25A tail. AluJ half SINE RNA was PCR amplified to isolate the 5' folded
484     Alu domain followed by variable polyA tail from 75A to 5A, non-A tail or ending in 23A-GC for
485     cryo-EM template. L1 3'UTR sequence of youngest L1 family, L1.3 (Genbank L19088.1) was
486     synthesized (IDT). Full-length L1 3'UTR or a truncation lacking 1-78 nt containing a G-
487     quadruplex were PCR amplified to include the T7 RNA polymerase promoter followed by 25A
488     sequence as the 3' end. RNA for *in vitro* reverse transcription assay was designed to result in
489     minimal secondary structure features; transcription templates were synthesized (IDT) and PCR
490     amplified. All RNAs were transcribed with T7 RNA polymerase in 40-100 μl reactions with
491     HiScribe T7 High Yield RNA Synthesis Kit (NEB). For high-resolution structure determination,
492     a synthetic template RNA was generated harboring a GC-rich hairpin, 15A sequence followed by
493     CAATA sequence for L1 ORF2p to polymerize and trap with a dideoxy-G, and 8 nucleotides
494     (TCGGCGCG) sequence complementary to the DNA primer (Supplementary Table 1). The DNA
495     template for this RNAs was synthesized as complementary oligonucleotides (IDT) to include the
496     T7 RNA polymerase promoter, sense and antisense strands were annealed by heating to 95°C and
497     slow cooling to 4°C and then transcribed using T7 RNA polymerase as described above. The *in*
498     *vitro* transcription reaction was performed for 5 hours at 37°C. The template DNA was removed

13

499  with DNase RQ1 (Promega), and the transcribed RNA was separated on a 6-9% denaturing
500  polyacrylamide gel. The RNA band was excised and eluted with RNA elution buffer (300 mM
501  NaCl, 10 mM Tris pH 8, 0.5% SDS, 5 mM EDTA) overnight at 4°C. The RNA was supplemented
502  with 25 µg glycogen and 300mM $NH_4OAc$ and further precipitated with ethanol, centrifuged, and
503  washed with 70% ethanol. The precipitated RNA was air dried before being dissolved in RNase-free
504  $H_2O$ and supplemented with Ribolock (ThermoFisher) for long-term storage in -20°C.
505
506  **Cryo-EM Sample Preparation and Data Collection**
507  Preparation of graphene oxide grids was adapted from our previously developed protocol[54].
508  Briefly, Quantifoil Au/Cu R1.2/1.3 grids 200-mesh (Quantifoil, Micro Tools GmbH, Germany)
509  were cleaned by applying two drops of chloroform, then glow discharged. 4 µl of 1mg/ml
510  polyethylenimine HCl MAX Linear Mw 40k (PEI, Polysciences) in 25mM K-HEPES pH 7.5 was
511  applied to the grids, incubated for 2 minutes, blotted away, washed twice with $H_2O$, and dried for
512  15 minutes on Whatman paper. Graphene oxide (Sigma, 763705) was diluted to 0.2 mg/ml in $H_2O$,
513  vortexed for 30 seconds, and precipitated at 1,200 xg for 60 seconds. 4 µl of supernatant was
514  applied to the PEI treated grids, incubated for 2 minutes, blotted away, washed twice with 4 µl
515  $H_2O$ each, and dried for 15 minutes on Whatman paper before using for grid preparation.
516      AluJ half SINE RNA for EM (141 nt) was diluted to 10 µM, then refolded in RNase-free
517  $H_2O$ by heating to 70°C for 5 minutes and slow cooling to 4°C for 2 hours. A 7 nt DNA primer
518  was added to refolded RNA at 1.5:1 primer:RNA molar ratio and annealed by heating to 30°C for
519  3 minutes and slow cooling to 4°C to assemble RD duplex. Synthetic template RNA (74nt) was
520  diluted to 10 µM, then refolded in RNase-free $H_2O$ by heating to 90°C for 3 minutes and snap
521  cooling to 4°C. An 8 nt DNA primer was added to refolded RNA at 1.5:1 primer: RNA molar ratio
522  and annealed by heating to 45°C for 3 minutes and snap cooling to 4°C to assemble RD duplex.
523  The cryo-EM sample was prepared by diluting wild-type L1 ORF2p to 600 nM concentration in
524  cryo-EM buffer (30 mM K-HEPES pH 7.9, 150 mM KCl, 10 mM $MgCl_2$, 5 mM EGTA, 1 mM
525  DTT). Assembled RD duplex was added to L1 ORF2p at 2:1 RD duplex: protein molar ratio. For
526  synthetic template RNA, dNTPs were added to the reaction to a final concentration of 1 mM dTTP,
527  1 mM dATP and 1 mM dideoxyGTP (ddGTP) to trap the L1 ORF2p-mediated reverse
528  transcription reaction. For SINE RNA, 1 mM dideoxyTTP (ddTTP) was added. The assembled
529  reaction was incubated at 37°C for 30 s to allow nucleic acid binding and complementary DNA
530  synthesis. 4 mM BS3 (ThermoFisher) was added to the reaction to crosslink the sample on ice for
531  5 minutes. 4 µl of sample was applied to the graphene oxide coated grid, incubated for 90 s at room
532  temperature, then washed with cryo-EM buffer. The grid was then blotted for 6 s with a blot force
533  of 5 at 20°C in 100% humidity and vitrified by plunging into liquid ethane using a Vitrobot Mark
534  IV (ThermoFisher).
535      For the L1 ORF2p-Alu RNP, micrographs were collected on a Titan Krios microscope
536  (ThermoFisher) operated at 300 keV and equipped with a K3 Summit direct electron detector
537  (Gatan). 23,878 movies were recorded using the program SerialEM at a nominal magnification of
538  105,000x in super-resolution mode (super-resolution pixel size of 0.405 Å/pixel) and with a
539  defocus range of -1.5 µm to -2.5 µm. The electron exposure was about 50 e⁻/Å². Each movie stack
540  contained 50 frames. For the L1 ORF2p-synthetic template RNP, the initial reconstruction was
541  obtained from datasets collected on a Talos Arctica microscope. 11,711 movies were recorded at
542  a nominal magnification of 45,000x in super-resolution mode (super-resolution pixel size of
543  0.4495 Å/pixel) and with a defocus range of -1.2 µm to -2.5 µm. The electron exposure was about
544  50 e⁻/Å². Each movie stack contained 50 frames. For the final reconstruction of the L1 ORF2p-

545 synthetic template RNP, we collected a large dataset on a Titan Krios G3i (ThermoFisher) operated
546 at 300 keV and equipped with a K3 Summit direct electron detector (Gatan) and an energy filter
547 with a slit width of 20 eV. A total of 23,874 movies were recorded at a nominal magnification of
548 105,000x in super-resolution mode (super-resolution pixel size of 0.405 Å/pixel), with a defocus
549 range of -1.0 μm to -2.5 μm. The electron exposure was about 50 e$^-$/Å$^2$. Each movie stack contained
550 50 frames.
551
552 **Cryo-EM Data Processing**
553 Cryo-EM data processing workflows are outlined in Extended Data Figs. 2 and 3. All movie frames
554 were motion corrected using MotionCor2[55] in RELION 3.1.1 and the corresponding super-
555 resolution pixels size was binned 2x during this process. Contrast transfer function (CTF)
556 parameters for each micrograph were estimated using CTFFIND4.1[57]. For the L1 ORF2p-synthetic
557 template RNP, a subset of micrographs were selected and around 2000 particles were manually
558 picked and inspected to train a Cryolo model using Cryolo v1.7.6[58]. The trained models were used
559 to predict particle locations on the entire dataset, for both the initial dataset acquired with a Talos
560 Arctica and the final dataset acquired with a Titan Krios. The particle picks from the Talos Arctica
561 session were imported to cryoSPARC v.3 to sort particles by 2D classification. 238,798 particles
562 from the initial dataset acquired with the Talos Arctica were imported back to RELION and a 3D
563 initial model was generated. After 3D classification of this dataset, class 1, containing 89,150
564 particles with apparent RNA density, were further processed to produce a 4.2 Å reconstruction.
565 For the final Titan Krios dataset, 786,013 particles, obtained after Cryolo picking and 2D
566 classification with cryoSPARC v.3 , were imported back to RELION and binned by 2. The 4.2 Å
567 reconstruction from the Talos Arctica dataset was filtered to 25 Å and used as the initial model for
568 a first round of 3D classification. A subset of 222,012 particles displaying a clearer RNA density
569 was selected, re-extracted with no binning, and refined to 3.3 Å. RNA-Focused 3D classification
570 without alignment was then performed and one class that displays the most complete RNA density,
571 containing 120,397 particles, was selected. Particle polishing and CTF refinement was performed
572 on this sub-set, followed by focused classification without alignment on the polyA tract RNA. The
573 final reconstruction was obtained at 3.2 Å nominal resolution from 111,564 particles. The cryo-
574 EM map was sharpened with post-processing in RELION for model building and display in the
575 figures.
576        For the L1 ORF2p-Alu RNP complex, the motion-corrected micrographs were imported to
577 cryoSPARC, 13 million particles were picked with a blob picker and sorted with 2D classification
578 down to 399,535 particles, which were then imported to RELION 3.1.1 for further processing. A
579 subset of these particles was used to generate an initial 3D model. 3D classification was performed
580 with the entire set of particles into 3 classes. A subset of 155,822 particles displaying a clear
581 density of the endonuclease domain and Alu RNA stem-loop and 5' fold was selected and refined
582 to 4.4Å.
583
584 **Model Building and Refinement**
585 Model building was initiated by rigid-body fitting the AlphaFold[36] model of human L1 ORF2p
586 into the final 3.3 Å cryo-EM density map using UCSF ChimeraX[60]. The endonuclease domain was
587 removed at this point due to the lower resolution in that part of the density map. The L1 ORF2p
588 protein was first manually inspected in COOT[61] to correct the amino acid sequence and then
589 subjected to real space refinement in PHENIX[62]. Amino acid side chains were manually inspected
590 in COOT and modified when needed before another round of real space refinement in PHENIX.

15

591 Nucleic acid was built using a difference density map generated from the cryo-EM density map
592 with the protein density subtracted. Core RNA-DNA duplex from a yeast RNA Pol III structure
593 (PDB 5FJ8) and dsRNA from a Drosophila Dicer-2 structure (PDB 7W0C) were first manually
594 docked into the cryo-EM map using UCSF ChimeraX. The L1ORF2-RNP was then manually
595 rebuilt in COOT using the nucleic acid difference map and the correct RNA and DNA sequences
596 bound to the protein core and the dsRNA sequence bound to L1 ORF2p. The single stranded RNA
597 was built de novo in COOT using the nucleic acid difference map. The model was corrected to
598 include the dideoxy-guanosine in the terminating DNA polymer obtained from PDB 1QSS, and
599 the following unincorporated dTTP obtained from PDB 1CR1. Both were docked into the density
600 map using UCSF Chimera and manually rebuilt with the corresponding DNA chain in COOT. The
601 model was subjected to global refinement using iterative rounds of real-space refinements in
602 PHENIX with rotamer and Ramachandran restraints. For dideoxy-guanosine, ligand restraints
603 were generated in PHENIX using the eLBOW tool. For the dTTP, ligand restraints were obtained
604 from PDB. PHENIX refinements were performed with these input restraints. At this point, the
605 endonuclease domain from the AlphaFold model of human L1 ORF2p was manually docked in
606 UCSF Chimera and merged into the model with COOT. The complete model was subjected to a
607 final real-space refinement and validation in PHENIX. Model building and validation statistics are
608 listed in Extended Data Table 1.
609

610 ***In vitro* Reverse Transcriptase Reactions**
611 For RT assays, the DNA primer was 5'-labeled with $^{32}$P γ-ATP (Perkin Elmer) using T4 PNK
612 (NEB). Unlabeled nucleotide was removed by spin column (Cytiva). Primer was annealed the RT
613 template RNA at 1:1 concentration by heating to 75°C for 3 minutes and slow cooling to 4°C for
614 1 hour. RT reactions were assembled on ice in 20 μl volume with final concentrations of 25 mM
615 Tris-HCl pH 7.5, 75 mM KCl, 35 mM NaCl, 5 mM MgCl$_2$, 10 mM DTT, 2% PEG-6K, 100 nM
616 RNA-DNA duplex, 0.1 units/μl M-MLV RT (Promega) or 100 nM L1 ORF2p wild-type or mutant
617 protein, 1 mM dNTPs. RT reactions were incubated at 37°C. 4.5 μl reaction was withdrawn at 0,
618 1, 5 and 20 minutes and mixed with 100 μl of stop solution (50 mM Tris-HCl pH 7.5, 20 mM
619 EDTA, 0.2 % SDS). Nucleic acid was purified with 1 volume (100 μl) of
620 phenol/chloroform/isoamyl alcohol and precipitated with 3 volumes of ethanol. Samples were then
621 pelleted at ~18,000 x g for 20 minutes at room temperature, washed with 7 volumes of 70% ethanol
622 and pelleted again at ~18,000 x g for 3 minutes. The pellet was air-dried resuspended in 5 μl water
623 and supplemented with 7 μl formamide loading dye (95% deionized formamide, 0.025% w/v
624 bromophenol blue, 0.025% w/v xylene cyanol, 5 mM EDTA pH 8.0). The sample was heated to
625 95°C for 3 minutes then placed on ice before loading the sample on a 7-8% Urea-PAGE gel. After
626 electrophoresis, the gel was dried, exposed to a phosphoimaging screen, and imaged by Typhoon
627 Trio (Cytiva). To quantitatively compare the RT activity of enzymes, we measured the gel intensity
628 of the full-length cDNA band for all enzymes used at various time points with ImageJ. The reaction
629 product generated by M-MLV RT at 5 minutes was used to normalize each intensity measurement
630 prior to combining data points from three separate repetitions of the RT assay. The mean intensity
631 and its standard deviation are plotted for each enzyme at each time point in Extended Data Figure
632 1e.
633

634 ***In vitro* Target Primed Reverse Transcription Reactions**
635 The target DNA site was synthesized (IDT) to have 3' phosphorylation modification on both the
636 top and bottom strands to block direct extension of the 3' ends by L1 ORF2p. The target DNA

637 strands were gel purified with denaturing Urea-PAGE (Supplementary Table 1), with the top
638 strand containing the cleavage (TTTTTAA) sequence. The top strand was 5'-labeled with $^{32}$P γ-
639 ATP (Perkin Elmer) using T4 PNK (NEB). Unlabeled nucleotide was removed by spin column
640 (Cytiva). The two strands were annealed at equimolar ratio by heating to 95°C and slow cooling
641 to 4°C over 1.5 hours. The template RNA was independently refolded by melting at 70°C for 5
642 minutes and snap cooling to 4°C prior to assembling the reaction. Target primed reverse
643 transcription (TPRT) reactions were assembled in 10 μl volume with final concentrations of 25
644 mM Tris-HCl pH 7.5, 75 mM KCl, 35 mM NaCl, 5 mM MgCl$_2$, 10 mM DTT, 2% PEG-6K, 1 mM
645 dNTPs, 50 nM annealed DNA duplex, 50 nM template RNA, 0.4 units/μl M-MLV RT (Promega),
646 200 nM L1 ORF2p wild-type or mutant proteins. Buffer or 200 nM PCNA was added in addition
647 to L1 ORF2p at 1:1 molar ratio in Extended Data Fig. 10b. TPRT reactions were incubated at 37°C
648 for 30 minutes and mixed with 90 μl of stop solution (50 mM Tris-HCl pH 7.5, 20 mM EDTA,
649 0.2% SDS). Nucleic acid was purified with 1 volume (100 μl) of phenol/chloroform/ isoamyl
650 alcohol and precipitated with 3 volumes of ethanol. Samples were then pelleted at ~18,000 x g for
651 15 mins at room temperature, washed with 7 volumes of 70% ethanol and pelleted again at ~18,000
652 x g for 3 minutes. The pellet was air-dried resuspended in 5 μl water and supplemented with 7 μl
653 formamide loading dye (95% deionized formamide, 0.025% w/v bromophenol blue, 0.025% w/v
654 xylene cyanol, 5 mM EDTA pH 8.0). The sample was heated to 95°C for 3 minutes then placed
655 on ice before loading the sample on a 9% Urea-PAGE gel. After electrophoresis, the gel was dried,
656 exposed to a phosphoimaging screen, and imaged by Typhoon Trio (Cytiva). To quantitatively
657 compare the EN nicking and TPRT activity across distinct target sites (Fig. 4b-e), distinct template
658 RNAs (Figs. 2c, 3f), protein mutations (Fig. 2g, 3d), or with addition of co-factors (Extended Data
659 Fig. 10b), we measured the gel intensity of the full-length TPRT product with ImageJ. The relative
660 TPRT product was measured by dividing the total TPRT product generated with each template
661 RNA, target site, or protein mutation by the total product for the condition used for the
662 normalization, highlighted in each figure legend. The relative EN nicking activity was measured
663 by dividing the total nicked target generated with each protein site by the total nicked target for
664 the condition used for the normalization, highlighted in each figure legend. The experiment and
665 analyses were repeated three independent times and the resulting average, and its standard
666 deviation is plotted in the bar graphs below each gel.
667

668 **Bioinformatics analysis**
669 Structure-based search for L1 ORF2p C-terminal segment homologs was performed by isolating
670 the coordinates for the C-terminal segment and comparing against 3D structures with the DALI
671 server[63]. Two hits for RTs included the insect non-LTR retroelement (PDB 8GH6) and human
672 TERT (PDB 7BG9). The C-terminal segment was aligned with these coordinates using the
673 MatchMaker tool in ChimeraX and displayed in Extended Data Fig. 6.
674 L1 ORF2p family of protein sequences were collected from a recent work[18] and by
675 searching for similar proteins in the UniProt database[64]. In total 14 full-length sequences were
676 aligned using Multiple Sequence Comparison by Log-Expectation (MUSCLE) tool in SnapGene
677 6.0 software (www.snapgene.com). Local alignments near the region of interest are displayed in
678 Extended Data Fig. 6c and the corresponding Genbank accession number or UniProt ID for each
679 sequence is listed.
680

681 **Comparison with R2 RT and Group II Intron RT**

17

682 *Bombyx mori* R2 RT (PDB 8GH6) and the *Thermosynechococcus elongatus* Group IIB intron RT
683 (PDB 6ME0) were aligned with human L1 ORF2p protein chain using the MatchMaker tool in
684 UCSF ChimeraX.
685
686

**Methods References**
688 53. Kroutter, E. N., Belancio, V. P., Wagstaff, B. J. & Roy-Engel, A. M. The RNA polymerase
689     dictates ORF1 requirement and timing of LINE and SINE retrotransposition. *PLoS Genet* **5**,
690     e1000458 (2009).
691 54. Patel, A., Toso, D., Litvak, A. & Nogales, E. *Efficient graphene oxide coating improves*
692     *cryo-EM sample preparation and data collection from tilted grids.*
693     http://biorxiv.org/lookup/doi/10.1101/2021.03.08.434344 (2021)
694     doi:10.1101/2021.03.08.434344.
695 55. Zheng, S. Q. *et al.* MotionCor2: anisotropic correction of beam-induced motion for improved
696     cryo-electron microscopy. *Nat Methods* **14**, 331–332 (2017).
697 56. Zivanov, J. *et al.* New tools for automated high-resolution cryo-EM structure determination
698     in RELION-3. *Elife* **7**, e42166 (2018).
699 57. Rohou, A. & Grigorieff, N. CTFFIND4: Fast and accurate defocus estimation from electron
700     micrographs. *J Struct Biol* **192**, 216–221 (2015).
701 58. Wagner, T. *et al.* SPHIRE-crYOLO is a fast and accurate fully automated particle picker for
702     cryo-EM. *Commun Biol* **2**, 218 (2019).
703 59. Punjani, A., Rubinstein, J. L., Fleet, D. J. & Brubaker, M. A. cryoSPARC: algorithms for
704     rapid unsupervised cryo-EM structure determination. *Nat. Methods* **14**, 290–296 (2017).
705 60. Pettersen, E. F. *et al.* UCSF ChimeraX: Structure visualization for researchers, educators,
706     and developers. *Protein Sci* **30**, 70–82 (2021).
707 61. Emsley, P., Lohkamp, B., Scott, W. G. & Cowtan, K. Features and development of Coot.
708     *Acta Crystallogr D Biol Crystallogr* **66**, 486–501 (2010).
709 62. Adams, P. D. *et al.* PHENIX: a comprehensive Python-based system for macromolecular
710     structure solution. *Acta Crystallogr D Biol Crystallogr* **66**, 213–221 (2010).
711 63. Holm, L., Laiho, A., Törönen, P. & Salgado, M. DALI shines a light on remote homologs:
712     One hundred discoveries. *Protein Sci* **32**, e4519 (2023).
713 64. UniProt Consortium. UniProt: the Universal Protein Knowledgebase in 2023. *Nucleic Acids*
714     *Res* **51**, D523–D531 (2023).
715

18

716 **Figure legends**

717

718 **Fig. 1. In vitro TPRT activity and cryo-EM structures of human L1 ORF2p-RNPs.** (a)
719 Domains of human L1 ORF2p. EN, Apurinic/apyrimidinic endonuclease; NTE, N-terminal
720 extension; RT, reverse transcriptase; RBD, RNA binding domain; CTS, C-terminal segment. (b)
721 Schematic of L1 ORF2p-mediated TPRT. (c) Denaturing gel analysis of TPRT reaction products.
722 Yellow square represents the $^{32}$P-labeled 5' end of the target DNA strand. Triangles indicate
723 expected TPRT product for full-length template (blue), incomplete cDNA synthesis (magenta) and
724 possible internal initiation (mustard). Wild-type L1 ORF2p was assayed with different template
725 RNAs with 25A 3' end: AY, AluY SINE (307nt); AJ, AluJ SINE (306nt); AJh, AluJ half-SINE
726 (141nt); L1, L1 3'UTR (231nt); L1Δ, L1 3'UTR ΔG-quadruplex (149nt). Here and for all
727 subsequent gels, DNA ladder length in nt is indicated on the left. The experiment was replicated
728 three times. Full-length cDNA product was quantified, normalized by the full-length cDNA
729 product with AJh RNA, and its mean±s.d. for n=3 biologically independent replicates are
730 displayed below. Here, and for all quantifications, black dots depict individual data points. (d)
731 Denaturing gel analysis of TPRT reaction products with wild-type L1 ORF2p, EN-dead (ΔEN)
732 and RT-dead (ΔRT) mutants. AJh 25A (141nt) was used as the template. The experiment was
733 replicated three times. (e) Cryo-EM density of L1 ORF2p in complex with AJh RNA:polyT primer
734 in an elongation state, segmented and colored by domains. (f) Cryo-EM density of L1 ORF2p in
735 an elongation state with synthetic template RNA and primer extended with cDNA, segmented and
736 colored by domains. Schematic of synthetic template RNA and cDNA primer used to obtain the
737 high-resolution cryo-EM structure is shown below, with the cDNA 3' ddG in yellow and the
738 incoming dTTP unable to join the cDNA also depicted. (g) Ribbon diagram of the L1 ORF2p-
739 RNP structure derived from (f) colored by domains and shown.

740

741 **Fig. 2. Recognition of the template RNA and its polyA tract.** (a) Schematic of direct interactions
742 between L1 ORF2p and the template RNA. Black lines denote hydrogen bonds, while mustard
743 lines denote hydrophobic contacts. Dashed lines represent direct contacts with the nucleobases or
744 ribonucleobases. (b) Recognition of the polyA tract by EN linker, NTE, RT, Thumb, and CTS
745 domains. (c) Denaturing gel analysis of TPRT reaction products with AJh template RNAs differing
746 in the 3' polyA tail length, including 75A (191nt), 50A (166nt), 25A (141nt), 20A (136nt), 15A
747 (131nt), 10A (126nt), 5A (121nt) or with a 3' tail of 25 non-A nt (25N, 141nt) or 20 non-A and
748 5A nt (20N5A, 141nt). 25N sequence is ggtaacgagaactgtcatgcacccc and 20NA5 sequence is
749 ggtaacgagaactgtcatgcaaaaa (Supplementary Table 1). The experiment was replicated three times.
750 Full-length cDNA product was quantified, normalized by the full-length cDNA product with AJh
751 75A, and its mean±s.d. for n=3 biologically independent replicates are displayed below. (d)
752 Adenine-specific hydrogen bonds between template A-60 and side chains in the NTE and Thumb
753 domains, alongside a hydrophobic contact with the RT domain. (e-f) Hydrogen bonds and
754 hydrophobic interactions between the template A-57 base and side chains in CTS and EN linker
755 domains, and between the A-55 base and EN linker, Thumb and NTE domain residues. Heteroatom
756 representation (oxygens in red and nitrogens in blue) is displayed. (g) Denaturing gel analysis of
757 TPRT reaction products with wild-type or Δss (single-stranded RNA binding) mutant L1 ORF2p
758 using AJh 25A template RNA. The experiment was replicated three times. The full-length cDNA
759 was quantified as the TPRT product, and the nicked product at the expected size was quantified
760 independently. Relative EN nicking and TPRT in the + RNA lanes were normalized by the wild-
761 type L1 ORF2p, and its mean±s.d. for n=3 biologically independent replicates are displayed below.

19

762

763 **Fig. 3. Engagement and unwinding of the template RNA by L1 ORF2p CTS.** (a) Base-reading
764 hydrogen bonds between the duplex-proximal polyA tract and residues in the CTS. Heteroatom
765 representation (oxygens in red and nitrogens in blue) is depicted. (b) Aromatic side chains from
766 the CTS domain near the 5' end of the polyA tract. (c) Isoleucine side chains from the CTS
767 insertion helix oblige unwinding of the RNA stem. (d) Denaturing gel analysis of TPRT reaction
768 products with wild-type L1 ORF2p, ΔIH (Insertion helix mutant) and ΔCTS mutants performed
769 with AJh 25A template RNA (141nt). The experiment was replicated three independent times. The
770 full-length cDNA was quantified as the TPRT product, and the nicked product at the expected size
771 was quantified independently. Relative EN nicking and TPRT in the + RNA lanes were normalized
772 by the wild-type L1 ORF2p, and its mean ± s.d. for n=3 biologically independent replicates are
773 displayed in the bar graph below. (e) Electrostatic rendering of the surface of L1 ORF2p engaging
774 the RNA stem-loop. Blue corresponds to positively charged and red negatively charged surface.
775 (f) Denaturing gel analysis of TPRT reaction products with wild-type L1 ORF2p and template
776 RNAs of variable stem-loop structures: AJh 25A, AluJ half-SINE (141nt); AJhm 25A, AluJ half-
777 SINE 25A with reduced stem bulges (142nt); AJh-uf 25A, AluJ half-SINE 25A with disrupted
778 stem base-pairing (141nt). The experiment was replicated three independent times. Full-length
779 cDNA product was quantified as the relative TPRT product, normalized by the full-length cDNA
780 product with AJh 25A, and its mean ± s.d. for n=3 biologically independent replicates are displayed
781 in the bar graph below.

782

783 **Fig. 4. Target site position and upstream single-stranded DNA determine efficiency of**
784 **nicking and TPRT.** (a) The full-length L1 ORF2p RNP structure (this study) superposed with a
785 structure of the EN domain: duplex DNA complex (PDB 78NS). For a cleavage site near the 5'
786 end of the duplex DNA, there is no steric clash with L1 ORF2p (upper panel). A modeled, longer
787 DNA duplex engaged with the EN domain, illustrates the steric clash of upstream duplex DNA
788 with L1 ORF2p (lower panel). (b-c) Denaturing gel analysis of TPRT reaction products using
789 target DNA with varying cleavage site position from 7 to 26 bp from the 5' end of the duplex DNA
790 in (b) and 5 to 13 bp in (a). (d-e) Denaturing gel analysis of TPRT reaction products using target
791 DNA with varying length and sequence of upstream single-stranded DNA. Blunt duplex end and
792 a short overhang with T-rich sequences were used in (d), while longer overhang lengths from 9 to
793 27 nt in (e). Red arrowhead in (b-e) denotes the expected nicked product size from cleavage at the
794 consensus target site. The experiments in (b-e) were replicated three independent times. Relative
795 amount of full-length cDNA was quantified as the TPRT product, and its mean ± s.d. for n=3
796 replicates are displayed in the bar graphs. The purple bar in (b-e) indicates the common DNA
797 target site across all assays, which was used for normalization of the relative TPRT product. The
798 rightmost lanes in (d) use a T-rich overhang of alternative sequence. AJh 25A (141 nt) was used
799 as the template RNA across (b-e). (f) Model for the initial stages of template engagement, target
800 site identification and first strand synthesis by L1 ORF2p. The overhang single-stranded DNA is
801 drawn near the CTS domain for illustration purposes only.

802

803 **Extended Data Fig. 1. Purification, electron microscopy and reverse transcriptase activity of**
804 **human L1 ORF2p and mutants.** (a) Size exclusion chromatogram (top) and SDS-PAGE of the
805 final step of L1 ORF2p purification stained with Coomassie dye (bottom). The experiment was
806 replicated more than 10 independent times. (b) Cryo-electron micrograph of L1 ORF2p-RNP
807 complex. The experiment was replicated more than 5 independent times. (c) Denaturing gel

20

808  analysis of TPRT reaction products with M-MLV RT (negative control) and wild-type L1 ORF2p
809  using AJh 25A as the template RNA (141nt). The experiment was replicated 3 independent times.
810  (d) Denaturing gel analysis of the amount of reverse transcribed product with RT template RNA
811  (129nt), base-paired at its 3' end to a 9 nt primer, after 0, 1, 5 and 20 minutes by M-MLV RT,
812  wild-type L1 ORF2p and L1 ORF2p mutants. RT mutant is RT-dead, and EN mutant is EN-dead.
813  (e) Intensity of full-length cDNA product was quantified and plotted across time for all proteins.
814  The experiment in (d) was replicated 3 independent times, cDNA product was normalized by the
815  cDNA product generated by M-MLV RT at 5 minutes, and the mean and standard deviation across
816  three repeats are plotted.

817

818  **Extended Data Fig. 2. Cryo-EM of L1 ORF2p RNP with Alu RNA.** (a) Secondary structure
819  schematic of AluJ half-SINE (AJh-EM) RNA and the DNA primer extended by the addition of
820  dideoxy-TTP used for cryo-EM. Bold nts denote the RNA and DNA bases visible in the cryo-EM
821  density map. RNA regions that bind SRP9/14 and L1 ORF2p are denoted with gray and blue
822  shading, respectively. (b) Cryo-EM data processing pipeline for L1 ORF2p in complex with the
823  Alu RNA and base-paired primer. A final cryo-EM density map at 4.4 angstrom resolution and the
824  corresponding FSC curve are displayed. (c) SRP9/14 bound to AJh RNA (PDB 5AOX[33]) was
825  superimposed with L1 ORF2p-Alu RNA structure using the common RNA stem between the two
826  structures, showing engagement of the SRP proteins to a distinct Alu RNA domain.

827

828  **Extended Data Fig. 3. Cryo-EM data processing for L1 ORF2p RNP complex bound to**
829  **synthetic template RNA.** Summary of single particle analysis pipeline leading to the
830  reconstruction of the L1 ORF2p RNP engaged with the synthetic template RNA described in
831  Figures 1-4 of this paper.

832

833  **Extended Data Fig. 4. Resolution estimation.** (a) Gold-standard FSC curve for the L1 ORF2p
834  RNP map, and map versus model FSC obtained from the final model after validation in Phenix (b)
835  Unsharpened density map obtained from analysis in Extended Data Fig. 4 was colored by local
836  resolution as estimated by Relion 3.1. (c) Particle orientation distribution in the final
837  reconstruction. (d) Representative map densities with atomic models for regions of RNA-DNA
838  duplex and protein.

839

840  **Extended Data Fig. 5. Active site conformation and supporting data for investigation of the**
841  **polyA tract and stem-loop engagement.** (a) RT active site residues involved in hydrophobic
842  interactions with the DNA and incoming dTTP are shown relative to the metal-binding aspartic
843  acid side chains of the active site (D702 and D703). (b) Hydrogen bonding interactions with the
844  incoming dTTP. Density of the EM map for the dTTP is displayed. (c) Interactions of RT and RBD
845  domain side chains with the duplex region of the template RNA. (d) NTE, Thumb and RT domain
846  residues interacting with the DNA primer and the cDNA, including a cDNA hydrogen bond with
847  Arg375 of the NTE. (e) Denaturing gel and SYBR-Gold staining of purified RNAs used in the
848  TPRT assays in Fig. 1c and Fig. 2c, showing their integrity and migration as expected. (f)
849  Deviation of the RNA stem-loop from a canonical A-form helix at the end contacted by the CTS
850  insertion helix. (g) Schematic of positively charged L1 ORF2p residues surrounding the RNA
851  stem-loop. Cα positions of all lysines and arginines near the RNA stem-loop from the CTS, RBD
852  and Thumb domain are displayed (bottom).

853

854 **Extended Data Fig. 6. Structure and sequence-based bioinformatics analysis on L1 ORF2p**
855 **CTS domain.** (a-b) Comparison of the L1 ORF2p CTS domain structure and CTS-like structures
856 in the *Bombyx mori* R2 enzyme in (a) and in human telomerase reverse transcriptase in (b). (c)
857 Sequence conservation for the insertion helix showing that Ile1122 is highly conserved across
858 predicted proteins from the L1 family. The IDs gi*xxx* represent Genbank accession codes and
859 A*xxx* and P*xxx* are Uniprot IDs.
860
861 **Extended Data Fig. 7. Analysis of off-target cleavage by L1 ORF2p.** The target DNA sequences
862 from Fig. 4b-c are indicated, and L1 ORF2p cleavage products are matched to sequence using
863 different colors of arrowhead. The green, yellow, and blue shades represent off-target cleavage
864 products, while red represent on-target cleavage. The annotated off-target cleavage products are
865 consistent with the L1 ORF2p cleavage site analysis described in a previous work[30].
866
867 **Extended Data Fig. 8. Analysis of target cleavage by ΔCTS L1 ORF2p.** Denaturing gel analysis
868 of EN cleavage products using target DNA with varying position of the cleavage site varied
869 between 7 and 26 bp from the 5' end of the duplex DNA, as denoted in the schematics above each
870 set of lanes. Expected nicked product size from cleavage at the consensus target site is denoted
871 with a red arrowhead. The experiment was replicated three times.
872
873 **Extended Data Fig. 9. Comparison between the L1 ORF2p RNP and related structures.** (a)
874 Comparison to target DNA-engaged Group IIB intron RNP structure with the RT protein bound
875 to intron RNA (PDB 6ME0)[45]. The RT domains are colored to directly compare with L1 ORF2p.
876 The DNA is colored gray, the intron RNA is colored red. DBD, DNA binding domain, colored
877 yellow. For clarity, only the regions of intron RNA, DNA and the RT protein that have an
878 equivalent in the L1 ORF2p RNP structure are displayed. (b) Comparison with the R2Bm TPRT
879 complex (PDB 8GH6)[48]. R2Bm domains are colored to directly compare with human L1 ORF2p.
880 The downstream DNA was removed for clarity, while the upstream DNA is colored gray, and the
881 RNA colored red for comparison with L1 ORF2p. RL, restriction enzyme like; ZnF, zinc finger
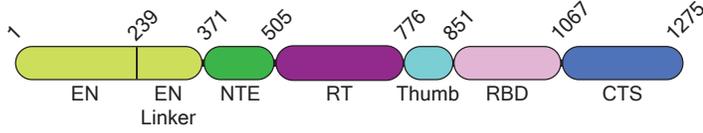882 domain for DNA binding; Myb, Myb domain for DNA binding.
883
884 **Extended Data Fig. 10. Proposed configuration of PCNA interaction with L1 ORF2p.** (a) Top
885 panel: the predicted PCNA interaction domain (PIP box)[18] within the NTE domain is highlighted.
886 Bottom panel: putative orientation of the PCNA trimer and L1 ORF2p based on existing structures
887 of PCNA with PIP-box containing protein complexes (PDB 7NV0). Based on the superposition of
888 the PIP box, PCNA would be expected to interact near the face of L1 ORF2p for entry of nucleic
889 acids (DNA and RNA), not near the exit channel of the product duplex. (b) Denaturing gel analysis
890 of TPRT reactions with L1 ORF2p in the presence of equimolar PCNA with AJh 25A (141nt) as
891 the template RNA. The experiment was replicated three times. Full-length cDNA product was
892 quantified as the relative TPRT product, normalized by the full-length cDNA product without
893 PCNA, and its mean and standard deviation of error across three replicates are displayed below.
894
895 **Extended Data Table 1. Cryo-EM data collection, refinement and validation statistics**
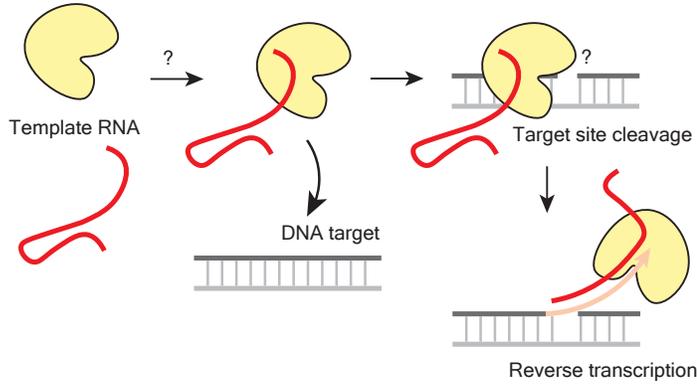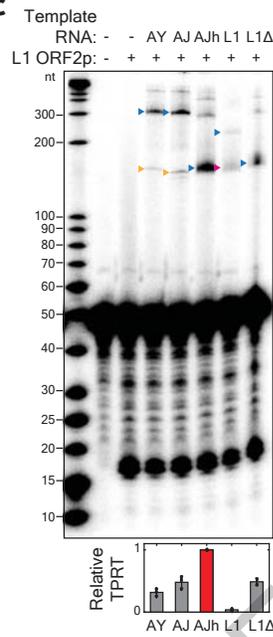896
897

# Figure 1

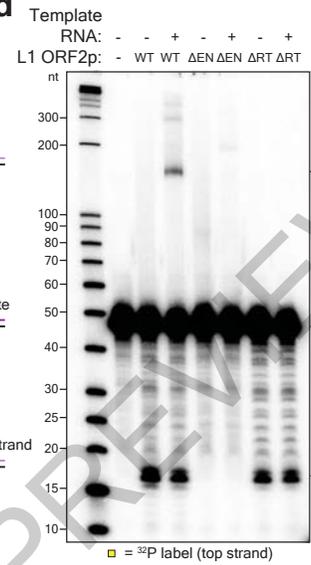**a** Human LINE-1 ORF2 protein



**b** L1 ORF2p



Template RNA
DNA target
Target site cleavage
Reverse transcription
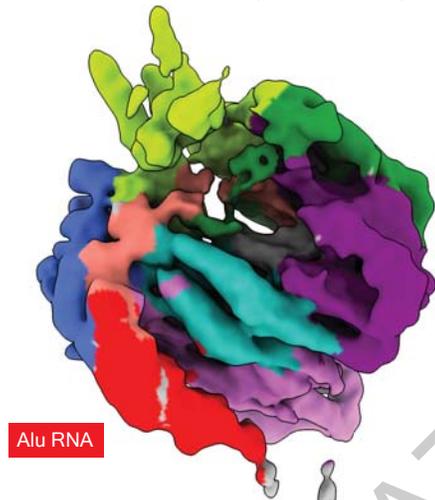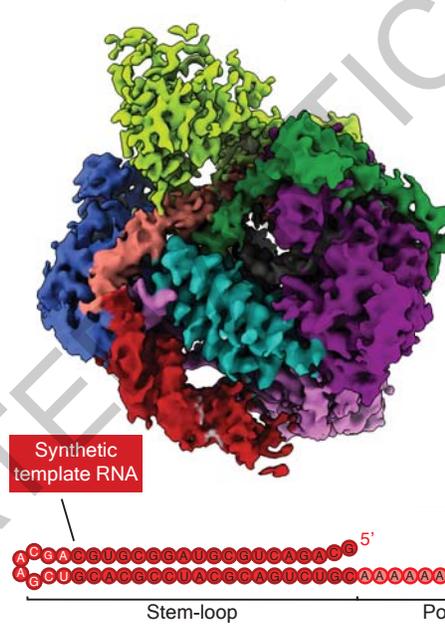
**c**



**d**



= $^{32}P$ label (top strand)

**e** cryo-EM density
4.4 Å resolution (FSC = 0.143)



Alu RNA

**f** cryo-EM density
3.2 Å resolution (FSC = 0.143)



Synthetic template RNA

Stem-loop    PolyA tract    cDNA    Primer

**g**



EN
NTE
CTS
Thumb
Synthetic template RNA
RT
RBD
cDNA

# Figure 2

# Figure 3



**a**

CTS
H1113
K1107
A-49
PolyA tract

**b**

PolyA tract
A-49
A-48
A-47
H1113
W1131
W1208
CTS

**c**

Insertion helix
I1122
I1121
G-1
C-2
A-3
5'
3'
C-46
G-45
U-44
RNA stem-loop

**d**

Template
RNA: − − + − + − +
L1 ORF2p: − WT WT ΔIH ΔIH ΔCTS ΔCTS
nt
300
200
100
90
80
70
60
50
40
30
25
20
15
10

Relative Nicking / Relative TPRT
WT  ΔIH  ΔCTS

**e**

Electrostatic surface color

**f**

Template
RNA: − AJh AJhm AJh-uf
L1 ORF2p: − + + +
nt
300
200
100
90
80
70
60
50
40
30
25
20
15
10

Relative TPRT
AJh  AJhm  AJh-uf

# Figure 4

# Extended Data Figure 1

**a** Purification of full-length LINE-1 ORF2 protein
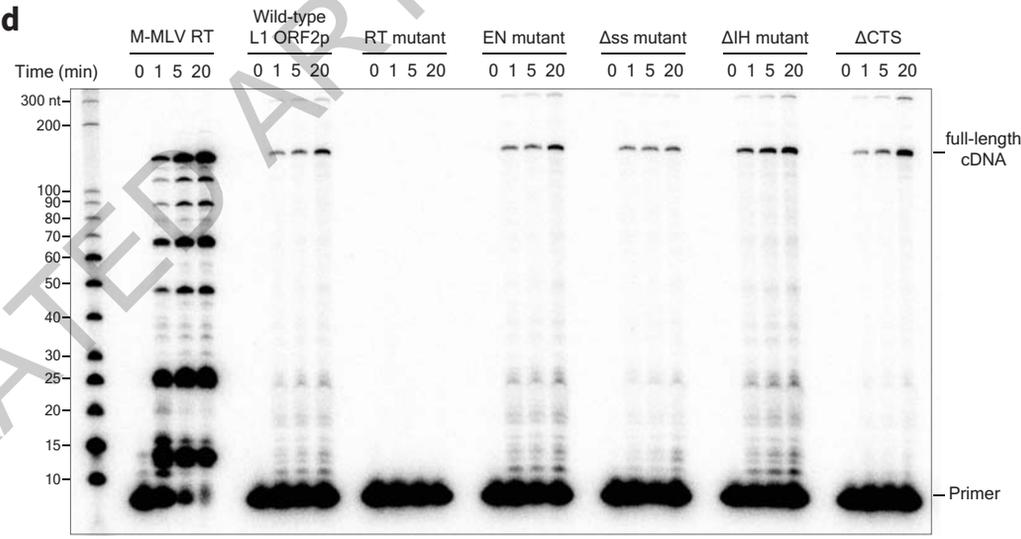
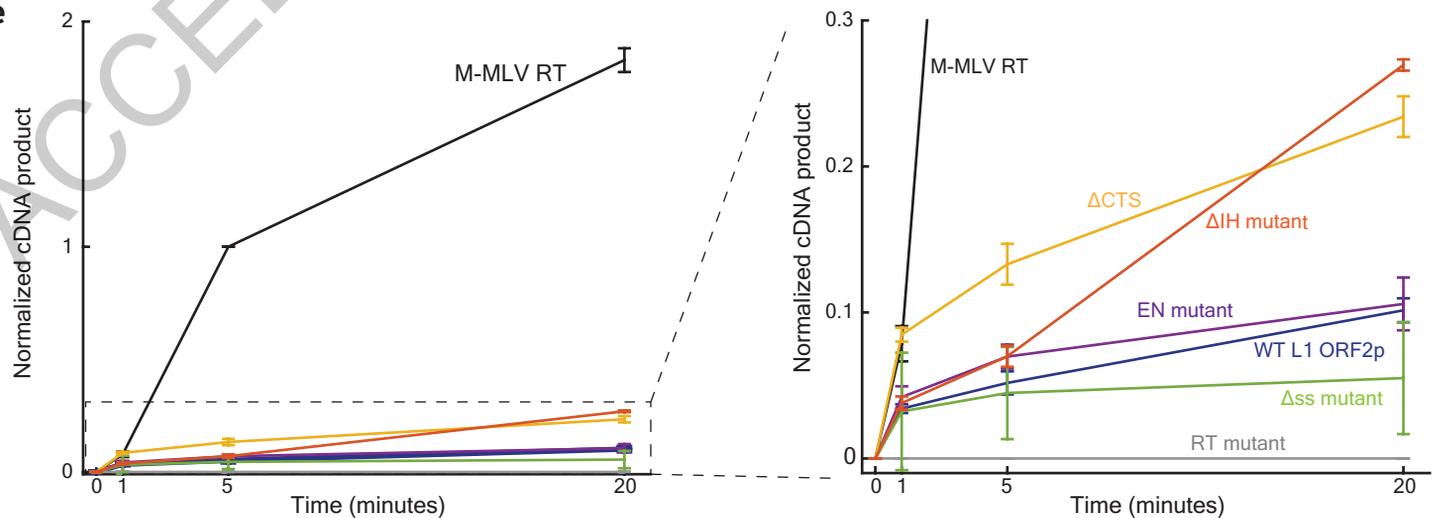**b** Cryo-electron micrograph of L1 ORF2p RNP complex



50 nm

**c**



**d**



**e**

# Extended Data Figure 2



**a** Template-primer duplex for cryo-EM

SRP 9/14 binding

L1 ORF2p binding

AJh-EM RNA

ddTTP

**b**

Titan Krios Dataset (AJh-EM RNA)

23,878 movies

13,000,000 p

2D classification
399,535 p

10 nm

Ab initio model and refinement

3D classification

75,902 p

167,811 p
(EN domain not visible)

155,822 p

Refinement
4.4Å Resolution

90°

(sharpened with B-factor -70)

Fourier Shell Correlation

— map vs. map

FSC = 0.143

4.4 Å

Resolution (1/Å)

**c**

L1 ORF2p

Alu RNA 3' tail

Alu stem loop (superposition with PDB 5aox)

SRP9/14 (PDB 5aox)

120°

L1 ORF2p

Alu stem loop (superposition with PDB 5aox)

Alu RNA 5' end (PDB 5aox)

SRP9/14 (PDB 5aox)

# Extended Data Figure 3



Talos Arctica Data
11,711 movies

1,082,827 p

Cryosparc Initial model
Flip hand

2D classification
237,798 p

~6 Å

Import to Relion
238,798 p

3D classification

89,150 p    102,383 p    47,265 p

Refinement
4.2Å Resolution

80º

Initial model
(filter to 25Å)

Titan Krios Data
23,874 movies

5,363,791 p

2D classification
786,013 p

Import to Relion
(bin 2)

10 nm

3D classification

222,012 p    317,237 p    246,764 p

Refinement

Re-extract unbinned particles

Refinement

RNA hairpin mask

3D focused classification w/o
alignment (N=3, T=10)

60,417 p    120,397 p    41,198 p

Refinement (3.53Å)

Polishing

CTF refinement

3D focused classification
on polyA tract (N=3, T=10)

6,494 p    2,339 p    111,564 p

Refinement
3.18Å Resolution

140º

(sharpened with B-factor -50)

# Extended Data Figure 4

**a**  Fourier Shell Correlation



**b**  Local resolution



Local resolution (Å)

High map threshold

Low map threshold

**c**  Angular distribution



**d**  Model fitting

# Extended Data Figure 5



**a** RT active site

D702, D703, A701, F700, F566, F605, dTTP, ddG-13, 3', 5'

**b** RT

R531, K603, A604, A-61, dTTP, 5', ddG-13, 3', cDNA

**c** Template, 5', S534, A-61, G660, K541, Q552, G-69, N865, G569, K1047, S1051, RBD, Y878, RT

**d** ddG-13, 3', F700, cDNA, Y823, P819, M815, Thumb, R375, NTE

**e** Purified template RNAs

**f** Distortion of RNA hairpin

C-46, G-45, U-44, A-3, G-1, C-2, 120°

L1 ORF2p bound RNA stem-loop
Ideal A-form RNA

**g** Positively charged residues near RNA stem-loop

CTS, Thumb, RBD, Lysine & Arginine (Cα), G-1, 5'

# Extended Data Figure 6

**a** Comparison with Bombyx mori R2 RT
(PDB 8GH6)

**b** Comparison with human Telomerase RT
(PDB 7BG9)

**c** Conservation of Insertion helix



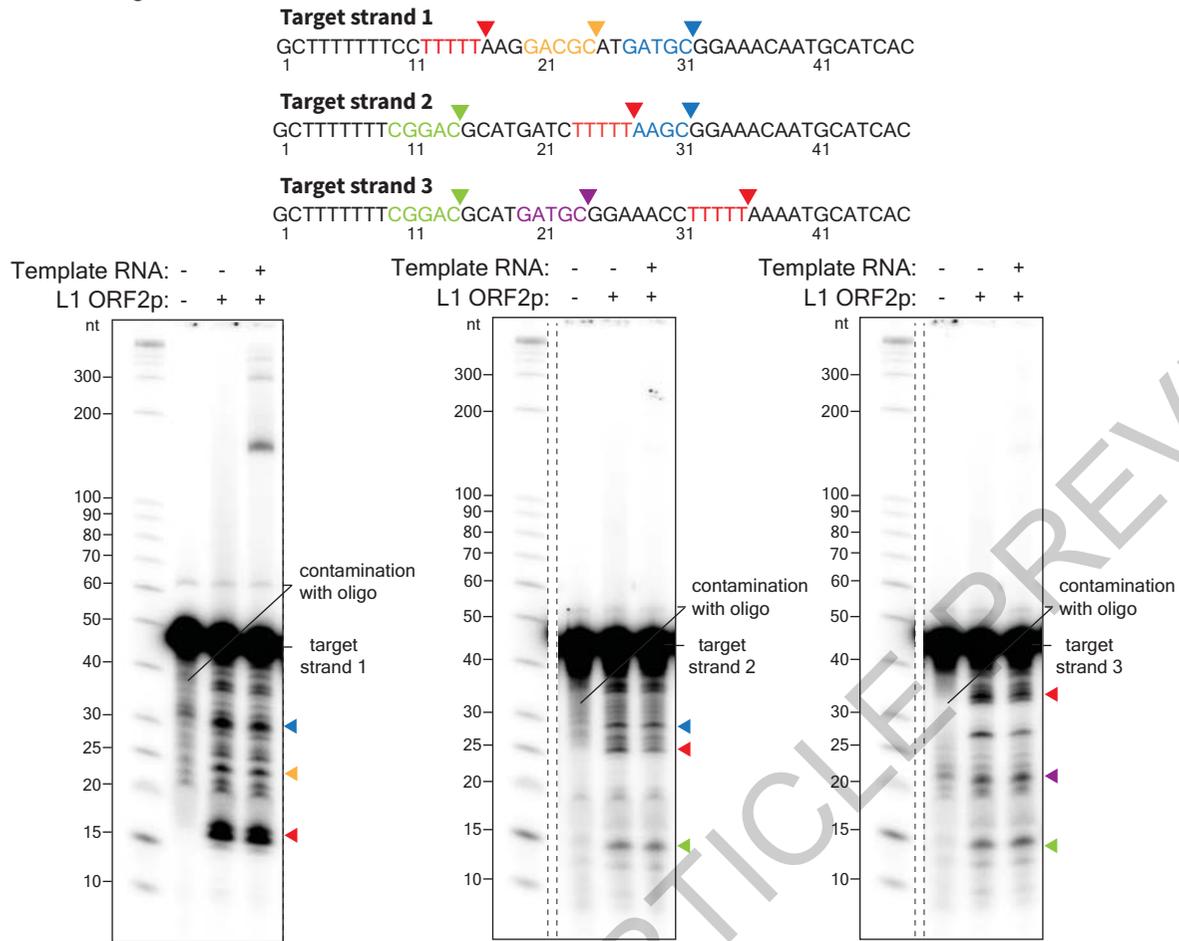| | | | 1122 |
|---|---|---|---|
| gi307098 | Human | IKTTMRYHLTP | VRMAIIKK |
| A0A8I5MVJ6 | Baboon | IKTTMRYHLTP | VRMAIIKK |
| A0A3Q1LG01 | Bovine | IKTTMRYHFTP | VRMAAIQK |
| A0A8D0MHJ9 | Pig | IKTTMRYHLTP | ARMAIIQK |
| P09548 | Slow loris | IKTTLRYHLTP | VRVAHITK |
| gi2981631 | Dog | IKTTMRYHLTP | VRMGKINK |
| gi3599320 | Mouse | IKTTLRFHLTP | VRMAKIKN |
| gi1791243 | Rat | IKTTLRFHLTP | VRMAKIKN |
| gi34392550 | Pufferfish | WRILHGAVAMN | IFISRMNP |
| gi34392555 | Pufferfish | WRLVYGVLAVN | KFVSILSL |
| gi34392557 | Pufferfish | WRVLHGIFPVN | SFVSTINQ |
| gi34392560 | Zebrafish | WRILHGAIAVN | AFVSIINP |
| gi34392563 | Zebrafish | WRIVHGIIATN | RHRAHIDP |
| gi34392575 | Zebrafish | FMLRHNCIMTE | IIFKKIGV |

Insertion helix
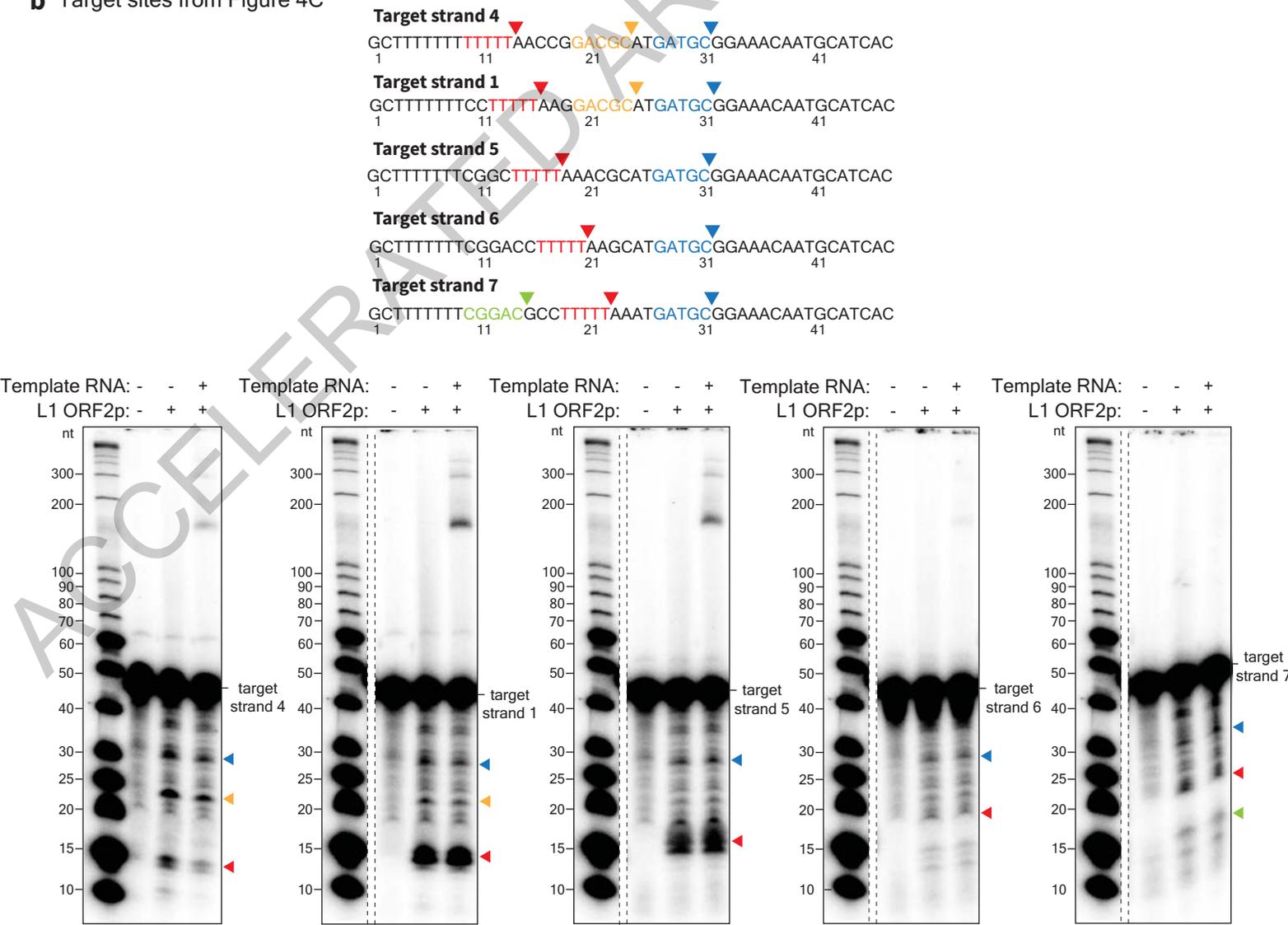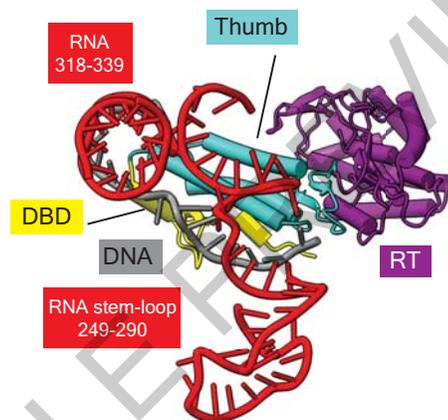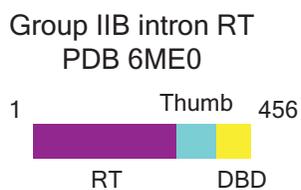
# Extended Data Figure 7

**a** Target sites from Figure 4B

**Target strand 1**

GCTTTTTTTCC TTTTT AAG GACGC AT GATGC GGAAACAATGCATCAC

1       11       21       31       41

**Target strand 2**

GCTTTTTTT CGGAC GCATGATC TTTTT AAGC GGAAACAATGCATCAC

1       11       21       31       41

**Target strand 3**

GCTTTTTTT CGGAC GCAT GATGC GGAAACC TTTTT AAAATGCATCAC

1       11       21       31       41



**b** Target sites from Figure 4C

**Target strand 4**

GCTTTTTTT TTTTT AACCG GACGC AT GATGC GGAAACAATGCATCAC

1       11       21       31       41

**Target strand 1**

GCTTTTTTTCC TTTTT AAG GACGC AT GATGC GGAAACAATGCATCAC

1       11       21       31       41

**Target strand 5**

GCTTTTTTT CGGC TTTTT AAACGCAT GATGC GGAAACAATGCATCAC

1       11       21       31       41

**Target strand 6**

GCTTTTTTT CGGACC TTTTT AAGCAT GATGC GGAAACAATGCATCAC

1       11       21       31       41

**Target strand 7**

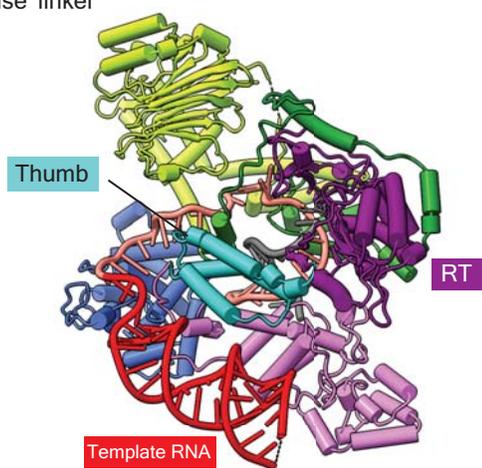GCTTTTTTT CGGAC GCC TTTTT AAAT GATGC GGAAACAATGCATCAC
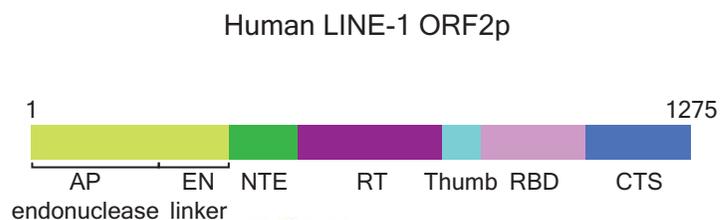
1       11       21       31       41

# Extended Data Figure 8

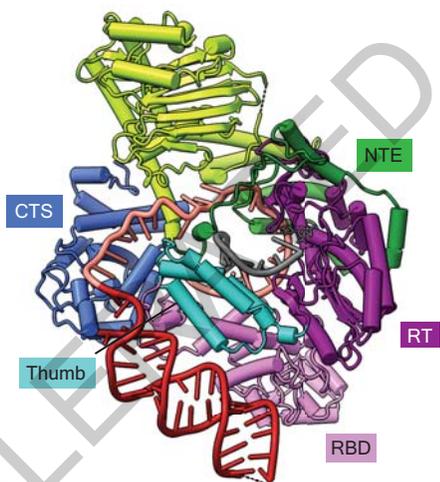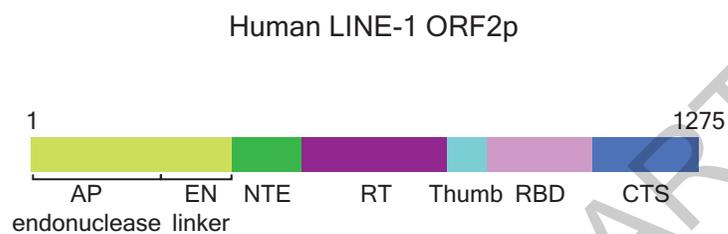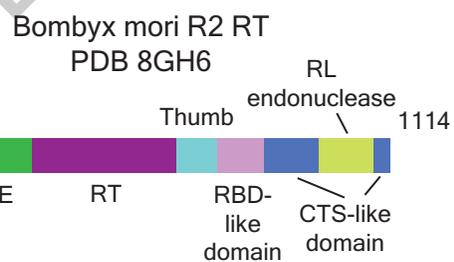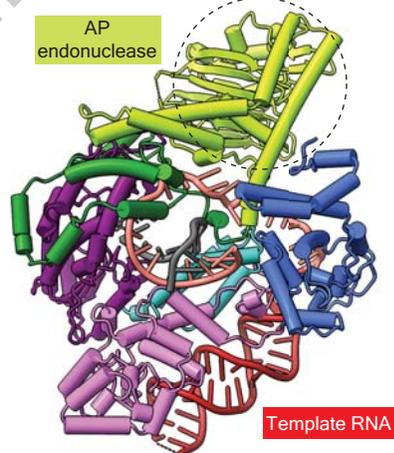# Extended Data Figure 9



**a**

**Human LINE-1 ORF2p**

1 — AP endonuclease / EN linker — NTE — RT — Thumb — RBD — CTS — 1275

**Group IIB intron RT PDB 6ME0**

1 — RT — Thumb — DBD — 456

**b**

**Human LINE-1 ORF2p**

1 — AP endonuclease / EN linker — NTE — RT — Thumb — RBD — CTS — 1275

**Bombyx mori R2 RT PDB 8GH6**

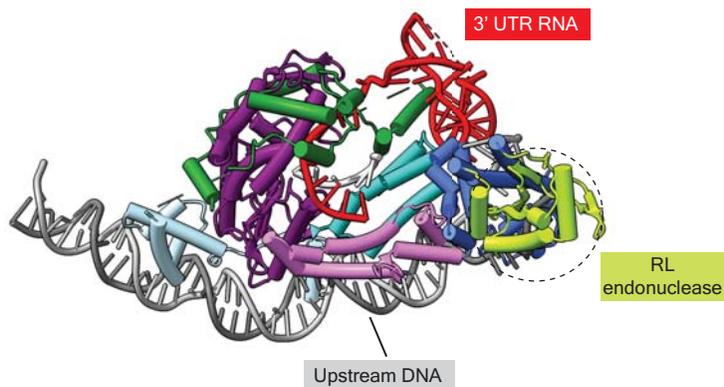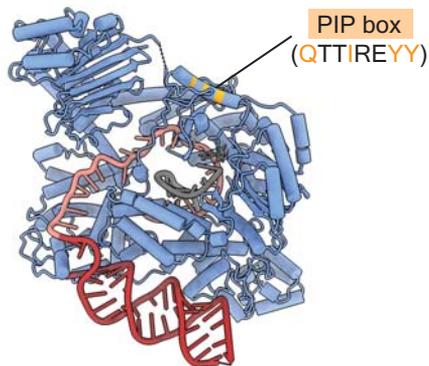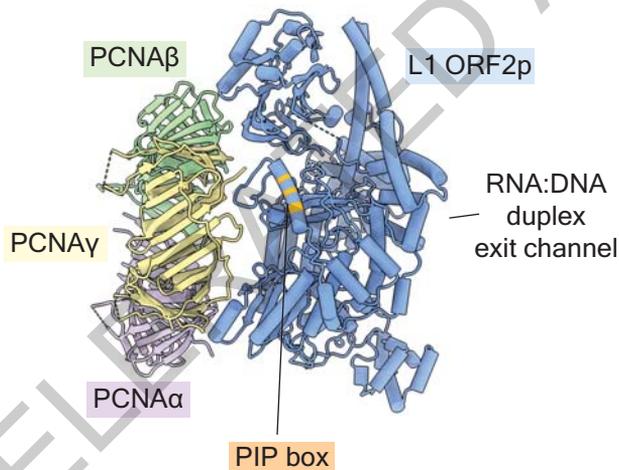111 — ZnF Myb — NTE — RT — Thumb — RBD-like domain — RL endonuclease — CTS-like domain — 1114

# Extended Data Figure 10

**a**   PCNA interaction domain



Orientation of PCNA relative to L1 ORF2p (modeled)

**b**

|  | Hs L1ORF2p RNP (EMD-42637) (PDB 8UW3) |
|---|---|
| **Data collection and processing** | |
| Magnification | 105,000 |
| Voltage (kV) | 300 |
| Electron exposure (e–/Å$^2$) | 50 |
| Defocus range (μm) | -1.0 to -2.5 |
| Pixel size (Å) | 0.81 |
| Symmetry imposed | *C1* |
| Initial particle images (no.) | 786,083 |
| Final particle images (no.) | 111,564 |
| Map resolution (Å) | 3.2 |
| FSC threshold | 0.143 |
| Map resolution range (Å) | 3.0 to 6.6 |
| | |
| **Refinement** | |
| Initial model used (PDB code) | none (generated in AlphaFold) |
| Model resolution (Å) | 3.3 |
| FSC threshold | 0.5 |
| Map sharpening *B* factor (Å$^2$) | -50 |
| Model composition | |
| Non-hydrogen atoms | 12,012 |
| Protein residues | 1,265 |
| Nucleic acid atoms | 73 |
| Ligands | 1 (dTTP) |
| *B* factors (Å$^2$) | |
| Protein | 162.81 |
| Nucleotide | 90.65 |
| Ligand | 80.32 |
| R.m.s. deviations | |
| Bond lengths (Å) | 0.004 |
| Bond angles (°) | 0.596 |
| Validation | |
| MolProbity score | 2.14 |
| Clashscore | 8.53 |
| Poor rotamers (%) | 2.51 |
| Ramachandran plot | |
| Favored (%) | 94.53 |
| Allowed (%) | 5.39 |
| Disallowed (%) | 0.08 |

**Extended Data Table 1**

Corresponding author(s): Eva Nogales, Kathleen Collins

Last updated by author(s): Nov 7, 2023

# Reporting Summary

Nature Portfolio wishes to improve the reproducibility of the work that we publish. This form provides structure for consistency and transparency in reporting. For further information on Nature Portfolio policies, see our Editorial Policies and the Editorial Policy Checklist.

## Statistics

For all statistical analyses, confirm that the following items are present in the figure legend, table legend, main text, or Methods section.

| n/a | Confirmed | |
|---|---|---|
| ☐ | ☒ | The exact sample size (*n*) for each experimental group/condition, given as a discrete number and unit of measurement |
| ☐ | ☒ | A statement on whether measurements were taken from distinct samples or whether the same sample was measured repeatedly |
| ☒ | ☐ | The statistical test(s) used AND whether they are one- or two-sided *Only common tests should be described solely by name; describe more complex techniques in the Methods section.* |
| ☒ | ☐ | A description of all covariates tested |
| ☒ | ☐ | A description of any assumptions or corrections, such as tests of normality and adjustment for multiple comparisons |
| ☐ | ☒ | A full description of the statistical parameters including central tendency (e.g. means) or other basic estimates (e.g. regression coefficient) AND variation (e.g. standard deviation) or associated estimates of uncertainty (e.g. confidence intervals) |
| ☒ | ☐ | For null hypothesis testing, the test statistic (e.g. *F*, *t*, *r*) with confidence intervals, effect sizes, degrees of freedom and *P* value noted *Give P values as exact values whenever suitable.* |
| ☒ | ☐ | For Bayesian analysis, information on the choice of priors and Markov chain Monte Carlo settings |
| ☒ | ☐ | For hierarchical and complex designs, identification of the appropriate level for tests and full reporting of outcomes |
| ☒ | ☐ | Estimates of effect sizes (e.g. Cohen's *d*, Pearson's *r*), indicating how they were calculated |

*Our web collection on statistics for biologists contains articles on many of the points above.*

## Software and code

Policy information about availability of computer code

| Data collection | Serial EM 4-0-20 for cryo-EM data collection |
|---|---|
| Data analysis | The EM softwares used: Relion 3.1.1, cryoSPARC v.3, cryoSPARC v.4, Cryolo 1.7.6. Structures were built using Coot 0.8.9, Chimera 1.14, ChimeraX 1.3, Phenix 1.20, Pymol 2.5.4. Gels were analyzed using ImageJ (Fiji) 2.1.0 |

For manuscripts utilizing custom algorithms or software that are central to the research but not yet described in published literature, software must be made available to editors and reviewers. We strongly encourage code deposition in a community repository (e.g. GitHub). See the Nature Portfolio guidelines for submitting code & software for further information.

## Data

Policy information about availability of data

All manuscripts must include a data availability statement. This statement should provide the following information, where applicable:

- Accession codes, unique identifiers, or web links for publicly available datasets
- A description of any restrictions on data availability
- For clinical datasets or third party data, please ensure that the statement adheres to our policy

The described cryo-EM maps and coordinate files were deposited in the Electron Microscopy Data Bank (EMDB) with accession code EMD-42637 and in Protein Data Bank (PDB) with accession code PDB 8UW3. All other datasets, reagents or resources generated during this study are available upon request from the corresponding authors.

# Research involving human participants, their data, or biological material

Policy information about studies with [human participants or human data](). See also policy information about [sex, gender (identity/presentation), and sexual orientation]() and [race, ethnicity and racism]().

| | |
|---|---|
| Reporting on sex and gender | n/a |
| Reporting on race, ethnicity, or other socially relevant groupings | n/a |
| Population characteristics | n/a |
| Recruitment | n/a |
| Ethics oversight | n/a |

Note that full information on the approval of the study protocol must also be provided in the manuscript.

# Field-specific reporting

Please select the one below that is the best fit for your research. If you are not sure, read the appropriate sections before making your selection.

☒ Life sciences   ☐ Behavioural & social sciences   ☐ Ecological, evolutionary & environmental sciences

For a reference copy of the document with all sections, see [nature.com/documents/nr-reporting-summary-flat.pdf](http://nature.com/documents/nr-reporting-summary-flat.pdf)

# Life sciences study design

All studies must disclose on these points even when the disclosure is negative.

| | |
|---|---|
| Sample size | In total 23,874 microscope raw movies collected from two different grid preparations were used for data processing of the highest resolution structure, sufficient to provide a high resolution structure. This data size was determined in order to reconstruct a high-resolution cryo-EM map for structure determination were obtained (at around 3 angstrom resolution). For low resolution, Alu structure 23,878 microscope raw movies collected from one grid preparations were used for data processing to yield a structure where protein and RNA densities could be clearly fitted at around 4 angstrom resolution. For biochemical assays, at least three independent biological replicates were performed, as recommended and as is the standard in similar works. |
| Data exclusions | Poor resolution data was excluded from cryo-EM analysis through 2D classifications and 3D classifications. This is standard step in single-particle cryo-EM analysis workflow and necessary to obtain highest resolution structures. |
| Replication | All biochemical experiments were repeated in three or more independent replicates, specified within the figure legends for individual experiments. All replicates which showed similar results. Data from all replicates were pooled for quantification and reported in bar graphs |
| Randomization | In the Fourier shell correlation (FSC) measurement step of the Relion 3.1 data processing pipeline, data were randomly divided into two halves resulting in two independently determined 3D volumes that were used for the FSC calculation. |
| Blinding | Data division in the FSC calculation step is a computer-based, unbiased process. Individual processing of different datasets collected from different human heart samples gave rise to the same 3D structures. |

# Reporting for specific materials, systems and methods

We require information from authors about some types of materials, experimental systems and methods used in many studies. Here, indicate whether each material, system or method listed is relevant to your study. If you are not sure if a list item applies to your research, read the appropriate section before selecting a response.

## Materials & experimental systems

| n/a | Involved in the study |
|---|---|
| ☒ ☐ | Antibodies |
| ☐ ☒ | Eukaryotic cell lines |
| ☒ ☐ | Palaeontology and archaeology |
| ☒ ☐ | Animals and other organisms |
| ☒ ☐ | Clinical data |
| ☒ ☐ | Dual use research of concern |
| ☒ ☐ | Plants |

## Methods

| n/a | Involved in the study |
|---|---|
| ☒ ☐ | ChIP-seq |
| ☒ ☐ | Flow cytometry |
| ☒ ☐ | MRI-based neuroimaging |

# Eukaryotic cell lines

| | |
|---|---|
| Cell line source(s) | SF9 cell line for baculovirus generation. SF9 and High5 cells lines for protein production. The SF9 and High5 cells were obtained from Invitrogen, ThermoFisher. |
| Authentication | No authentication of cell lines was performed as they were purchased from reliable commercial sources. |
| Mycoplasma contamination | Cells were tested for mycoplasma contamination and were found to be negative. Cell lines were monitored for doubling time and correct morphology. |
| Commonly misidentified lines (See ICLAC register) | No misidentified cell lines were used in this work. |

# Plants

| | |
|---|---|
| Seed stocks | n/a |
| Novel plant genotypes | n/a |
| Authentication | n/a |