

Recording gene expression order in DNA by CRISPR addition of retron barcodes

Santi Bhattarai-Kline¹, Sierra K. Lear^{1,2}, Chloe B. Fishman¹, Santiago C. Lopez^{1,2}, Elana R. Lockshin³, Max G. Schubert^{4,5}, Jeff Nivala⁶, George Church^{4,5}, Seth L. Shipman^{1,7*}

¹Gladstone Institute of Data Science and Biotechnology, San Francisco, CA, USA

²Graduate Program in Bioengineering, University of California, San Francisco and Berkeley, CA, USA

³Department of Neurobiology, Duke University Medical Center, Durham, NC, USA

⁴Department of Genetics, Harvard Medical School, Boston, MA, USA

⁵Wyss Institute for Biologically Inspired Engineering, Harvard University, Boston, MA, USA

⁶Paul G. Allen School of Computer Science & Engineering, University of Washington, Seattle, WA, USA

⁷Department of Bioengineering and Therapeutic Sciences, University of California, San Francisco, CA, USA

*Correspondence to: seth.shipman@gladstone.ucsf.edu

1 **ABSTRACT**

2 **Biological processes depend on the differential expression of genes over time, but**
3 **methods to make physical recordings of these processes are limited. Here we report a**
4 **molecular system for making time-ordered recordings of transcriptional events into**
5 **living genomes. We do this via engineered RNA barcodes, based on prokaryotic retrons¹,**
6 **which are reverse-transcribed into DNA and integrated into the genome using the**
7 **CRISPR-Cas system². The unidirectional integration of barcodes by CRISPR integrases**
8 **enables reconstruction of transcriptional event timing based on a physical record via**
9 **simple, logical rules rather than relying on pre-trained classifiers or post-hoc inferential**
10 **methods. For disambiguation in the field, we will refer to this system as a Retro-**
11 **Cascorder.**

12 **INTRODUCTION**

13 DNA is the universal storage medium for cellular life. In recent years, an emerging field
14 of biotechnology has begun repurposing DNA to store data that has no cellular function. The
15 same qualities of DNA that are beneficial in a biological context – high information density, ease
16 of copying, and durability – also enable flexible storage of text, images, and sound³⁻⁶. Extending
17 this general concept, researchers have developed data storage systems contained within living
18 organisms that allow the recording of biological signals into DNA, such as endogenous
19 transcription and environmental stimuli. One particular avenue of interest for such systems is in
20 the longitudinal recording of biological processes within cells⁷⁻⁹.

21 These recordings address a fundamental limitation in standard methods to interrogate
22 complex biological processes that require the destruction of cells and, thus, can only provide
23 measurements at single points in time (e.g. RNA-Seq). Because biological processes are not
24 perfectly synchronized at the cellular level, any individual cell collected in the middle of a
25 biological process could be either ahead or behind in the progression of events relative to any

26 other cell collected at that same time. This cellular heterochronicity makes it impossible to
27 definitively reconstruct time-dependent processes from the destructive measurement of parallel
28 samples. Indeed, cell-to-cell heterochronicity has actually been exploited in computational
29 methods to infer position in a biological process among cells within a single sample (e.g. single-
30 cell RNA-Seq pseudotime)¹⁰. However, these methods of inference make assumptions about
31 the relationship between cells that are not explicitly known, and often require user-imposed
32 constraints or the incorporation of prior biological knowledge¹¹.

33 An approach known as molecular recording provides an alternative to statistical
34 inference. Molecular recorders are biological devices that continuously record cellular
35 processes, storing a physical record of the data permanently in cellular DNA, so that it may be
36 retrieved at the very end of an experiment or process. Approaches to build molecular recorders
37 have relied on different methods of modifying DNA, including site-specific recombinases and
38 CRISPR-Cas nucleases^{7,12,13}. Another approach to molecular recording, which we have worked
39 to develop, leverages CRISPR-Cas integrases¹⁴. CRISPR-Cas integrases have been previously
40 used to encode information into CRISPR arrays through the delivery of chemically synthesized
41 oligos^{4,14} or by modulating the copy number of a reporter plasmid in response to a biological
42 stimulus^{5,8}. However, the ability to record the temporal order of more than one different
43 biological signals into the CRISPR array of a single cell has not yet been demonstrated.

44 Here, we demonstrate successful recording of temporal relationships by adding a new
45 molecular component to the system: a retroelement called a retron. The compact size,
46 specificity, and flexibility of retrons to produce customizable DNA *in vivo* make them an
47 attractive tool for biotechnology. Previously, retrons have been used in applications such as
48 genome editing in several host systems¹⁵⁻¹⁸ and early analog molecular recorders¹⁹. By
49 combining the functions of retrons and CRISPR-Cas integrases, we have built a system to make
50 temporal recordings of transcriptional events.

51 To record transcriptional events, we engineered retrons to produce a set of compact,
52 specific molecular tags, which can be placed under the control of multiple promoters of interest
53 inside a single cell. When a tagged promoter is active, the tag sequence is transcribed into
54 RNA, and reverse transcribed by the retron RT to generate a DNA ‘receipt’ of transcription. That
55 DNA ‘receipt’ is then bound by Cas1-Cas2 and integrated into the cell’s CRISPR array, creating
56 a permanent record of transcription. If another tagged promoter subsequently becomes active, a
57 different DNA ‘receipt’ can be generated and integrated into the CRISPR array following the first
58 spacer. By producing a linear record of these ‘receipts’ in the genome, we have built a biological
59 device, called a Retro-Cascorder, that records the temporal history of specific gene expression
60 events into the CRISPR arrays of individual cells (Fig. 1a).

61 **RESULTS**

62 **Cas1-Cas2 integrates retron RT-DNA**

63 CRISPR-Cas systems function as adaptive immune systems in bacteria and archaea.
64 During the first phase of the immune response to infection by phage or mobile genetic elements,
65 called adaptation, the CRISPR proteins Cas1 and Cas2 integrate a piece of foreign DNA into a
66 genomic CRISPR array. The CRISPR array consists of a leader sequence followed by unique
67 spacer sequences derived from foreign DNA, which are all separated by identical sequences
68 called repeats. The sequence information stored in the spacers serves as an immunological
69 memory of previous infection. This machinery, comprised of the CRISPR array, Cas1, and
70 Cas2, is a ready-made storage device. When the Cas1-Cas2 complex integrates a spacer into
71 the CRISPR array, it is added next to the leader sequence and the previous spacers are shifted
72 away from the leader^{20,21}. Thus, spacers which are further away from the leader sequence were
73 acquired further in the past, and those closer to the leader acquired more recently.

74 The first challenge in building a temporal recorder of gene expression was to generate
75 specific DNA barcodes following a transcriptional event, which can be permanently stored in a

76 cell's genome via integration by Cas1-Cas2. For integration, Cas1-Cas2 require DNA of at least
77 35 bases from end-to-end, with a 23 base complementary core region, and a protospacer-
78 adjacent motif (PAM)²². To generate these acquirable pieces of DNA on-demand in cells, we
79 used retrons. Recently determined to function in bacteria as a defense system against phage
80 infection²³, a typical retron consists of a single operon that controls the expression of: (1) a
81 small, highly structured noncoding RNA (retron ncRNA), (2) a retron reverse transcriptase (RT)
82 that specifically recognizes and reverse-transcribes part of its cognate ncRNA, and (3) one or
83 more effector proteins which are implicated in downstream functions^{23,24}. We designed variant
84 ncRNA sequences of a native *E. coli* retron, Eco1^{1,25} (Ext. Data Fig. 1a), for integration into the
85 genome by the type I-E CRISPR system of *E. coli* BL21-AI cells²⁰ after they are reverse
86 transcribed (Fig. 1b).

87 We tested multiple variants of Eco1 ncRNA for both reverse-transcription functionality
88 and the ability of their RT-DNA to be acquired by the CRISPR adaptation machinery, and
89 identified two that accomplish these aims. When overexpressed in *E. coli*, variants v32 and v35
90 (Ext. Data Fig. 1b-c) produced robust levels of RT-DNA that could be easily visualized on a
91 PAGE gel (Fig. 1c), had perfect 3'-TTC PAM sequences, and could theoretically hybridize to
92 create a 23-base core. Rather than a single copy of the retron RT-DNA hairpin forming the
93 prespacer for acquisition, both v32 and v35 are designed such that two copies of the RT-DNA
94 can form a duplex, which we hypothesized would be efficiently integrated into the CRISPR array
95 (Fig. 1d, Ext. Data Fig. 1b-c). To measure the ability of variant retrons to be acquired, we
96 overexpressed the variant ncRNA, Eco1 RT, and Cas1-Cas2 in BL21-AI cells which harbor a
97 single CRISPR array in their genome. We then sequenced the CRISPR arrays of these cells to
98 quantify integrations. In both cases, we found new spacers in these cells that matched the
99 sequence of the retron RT-DNA that was expressed (Fig. 1e). Critically, arrays containing
100 retron-derived spacers were only seen when cells also harbored a plasmid coding for Eco1 RT
101 (Fig. 1e), indicating that the retron-derived spacers were indeed a result of the production of RT-

102 DNA, rather than being derived exclusively from plasmid DNA. Retron v35 was acquired at a
103 higher rate than v32 (Fig. 1e), and was selected for use in subsequent work.

104 We further modified v35 by extending the length of the non-hairpin duplex region
105 referred to as the a1/a2 region (Fig. 1f). We have previously shown that this modification to
106 retrons both increases production of RT-DNA in bacteria and yeast and increases the efficiency
107 of genome editing methods which rely on retrons¹⁸. Consistent with our previous findings,
108 extending the a1/a2 region of retron v35 resulted in an increase in the percentage of arrays
109 which contained retron-derived spacers (Fig. 1f). This suggests that, like RT-DNA-templated
110 genome editing, the rate of acquisition of retron-derived spacers is dependent on the
111 abundance of RT-DNA. To take advantage of this improved acquisition efficiency, we
112 incorporated this modification into all future Eco1 constructs.

113 To better characterize the acquisition of spacers by Cas1-Cas2 over time, we expressed
114 retron v35 and Cas1-Cas2 for 24 hours and sampled arrays at regular intervals throughout (Fig.
115 1g-i). This showed that the number of arrays that contain retron-derived spacers increased
116 regularly over time (Fig. 1g). As retron-derived spacers accumulated, they were accompanied
117 by spacers derived from the cell's genome and from plasmids, as previously described²⁰. These
118 non-retron-derived spacers also increased in arrays over time (Fig. 1h). The proportion of new
119 retron-derived spacers remained relatively stable over time, making up between 1-10% of new
120 spacer acquisitions (Fig. 1i). Thus, the abundance of retron-derived spacers can be used as a
121 proxy for the duration of a transcriptional event. This result demonstrates a new implementation
122 of analog molecular recording, similar in function to those previously described¹⁹, but based on
123 the marriage of retrons and CRISPR-Cas integrases.

124

125 **Diversification of retron-based barcodes**

126 A crucial advantage of retron-based molecular recording is the ability to follow multiple
127 transcripts of interest by capturing distinct events within a single genomic CRISPR array. This

128 enables the recording of gene expression timing within genetically identical cells, rather than
129 relying on a mixed population of cells, each harboring different sensors. The specificity of
130 retrons also enables more focused recordings compared to promiscuous RTs, which cannot be
131 made to selectively reverse-transcribe individual transcripts^{9,26}. Additionally, in contrast to
132 recombinase-based molecular recording systems, the retron-based approach should enable a
133 much larger set of sensors to coexist within a population of genetically identical cells. This is
134 because the set of barcoded retrons is only limited by DNA sequence, rather than by the
135 comparatively small number of well-characterized recombinases^{7,27}. To construct a set of unique
136 retron tags, we chose to use the loop in retron v35's RT-DNA hairpin as a six-base barcode
137 (Fig. 2a, Ext. Data Fig. 2). This barcoding strategy allows multiple otherwise identical ncRNAs to
138 be reverse transcribed by the same RT, but remain easily distinguishable by sequence in
139 CRISPR arrays. We synthesized a set of barcoded retrons, expressed them in cells along with
140 Cas1-Cas2, and analyzed how efficiently they were acquired by sequencing CRISPR arrays
141 (Fig. 2b). We compared these barcoded variants to the original v35 retron, and included a dead-
142 RT version of the v35 retron as a negative control. Overall, we observed differences in the rate
143 at which different barcoded retrons were acquired, ranging from no significant difference up to a
144 ~70% reduction in acquisitions (Fig. 2b). The differences in acquisition efficiency of the different
145 barcoded retrons may come from changes in the efficiency of RT-DNA production, which we
146 have observed to occur when changing the stem and loop region of retron ncRNAs¹⁸, and
147 changes in the efficiency of acquisition, which we have observed to occur with different
148 prespacer sequences¹⁴.

149 To test our ability to discriminate between the barcoded spacers derived from this set of
150 retrons, we searched the sequence data from each sample expressing one barcoded retron for
151 all of the other barcodes in the set. In our computational pipeline, we specify a tolerance of up to
152 3 bases of mismatches or indels (out of a 23 base search sequence) when determining the
153 identity of a retron-derived spacer. This is to compensate for minor differences which may be

154 found in mature spacers compared to their hypothetical sequence. As such, if our retron
155 barcodes are faithfully preserved through all steps of the recording process (DNA coding
156 sequence → RNA → RT-DNA → CRISPR array), then we should be able to effectively
157 distinguish between barcodes which differ by 4 bases or more. This proved to be true when we
158 examined our original set of 9 barcodes for orthogonality *in-silico*. Barcodes which differed by
159 less than 4 bases could not be differentiated and barcodes which differed by 4 bases or more
160 could be distinguished from each other with perfect accuracy, forming a set of 6 mutually
161 orthogonal barcodes (Fig. 2c-d). This demonstrated that barcode sequences in retron-based
162 transcriptional tags are faithfully preserved throughout the process of molecular recording,
163 allowing for the facile construction of sets of mutually orthogonal tags.

164

165 **Mechanism of RT-DNA spacer acquisition**

166 While it has been demonstrated that Cas1-Cas2 can integrate prespacers consisting of
167 two complementary strands of DNA into the CRISPR array¹⁴, recent evidence suggests that
168 Cas1-Cas2 are capable of binding ssDNA and may in fact bind the two strands of a prespacer
169 separately, in a stepwise fashion²⁸. To date, all experimentally characterized retrons have been
170 shown to use the 2'-OH from a conserved guanosine to initiate reverse-transcription^{1,24}, leaving
171 a 2'-5' RNA-DNA linkage. We hypothesized that this unique feature of RT-DNA prespacers
172 might allow us to further interrogate the mechanism of prespacer loading and spacer acquisition
173 by Cas1-Cas2.

174 Unlike many prespacers examined in prior work, which generally form perfect DNA
175 duplexes, the duplex which we believe is formed by our retron has three characteristic regions
176 where the prespacer should contain mismatches. The first of these regions is a stretch of five
177 bases which, after integration, is located closest to the leader sequence. We will refer to this
178 region as the leader-proximal, or LP, region (Fig. 3a). Next, there is a single base mismatch
179 which falls near the middle of the mature spacer. We will refer to this region as the middle, or M,

180 region (Fig. 3a). Finally, the last of the mismatched regions is found, in the mature spacer, in the
181 five bases furthest from the leader sequence. We will refer to this as the leader-distal, or LD,
182 region (Fig. 3a). We found that in retron-derived spacers, the sequence of these mismatched
183 regions either corresponded to one strand of our hypothesized prespacer duplex or the other
184 (Fig. 3b). In this analysis, we will refer to the two strands of the hypothetical prespacer as the
185 (+) and (-) strands. In Eco1-derived spacers, the sequence in the LP region overwhelmingly
186 corresponded to the (-) strand. This (-) strand contains the PAM-proximal 3'-end, which
187 determines directionality¹⁴ and has been shown to be integrated second in the spacer
188 integration process^{28,29}. This pattern of preserving the PAM-derived 3'-end sequence in the LP
189 region was also seen when cells were electroporated with a synthetic oligonucleotide version of
190 the retron RT-DNA (Fig. 3b).

191 At the opposite end of the spacer, however, retron-derived and oligo-derived spacers
192 were not identical. In the LD region, oligo-derived spacers overwhelmingly mapped to the (+)
193 strand of our hypothesized duplex, whereas the LD regions of retron RT-DNA-derived spacers
194 predominantly mapped to the (-) strand (Fig. 3b). Because the *in vivo*-produced RT-DNA
195 contains a 2'-5' linkage and the oligo does not, we suspected that the 2'-5' linkage present in the
196 Eco1 RT-DNA may interfere with the CRISPR adaptation process. To test this, we treated
197 purified Eco1 RT-DNA with the eukaryotic debranching enzyme DBR1, which natively
198 processes RNA lariats by cleaving 2'-5' bonds in RNA³⁰. Treatment of Eco1 RT-DNA with DBR1
199 *in vitro* resulted in a characteristic downward shift in the size of Eco1 RT-DNA from the loss of a
200 small number of ribonucleotides remaining at the branch point. DBR1 treatment also rendered
201 Eco1 RT-DNA sensitive to the 5'-exonuclease recJ (Fig. 3c). This indicates that DBR1 is able to
202 remove the 2'-5' linkage and produce Eco1 RT-DNA with an unbranched 5'-end. When purified
203 Eco1 RT-DNA was treated with DBR1 and electroporated back into cells expressing Cas1-
204 Cas2, the LD sequences of retron-derived spacers closely resembled those of retron-derived
205 spacers after oligo electroporation (Fig. 3d), indicating that the presence of the 2'-5' linkage in

206 Eco1 RT-DNA is responsible for its unique pattern of spacer sequences. One potential
207 explanation for the spacer pattern observed in Figure 3b is that, in addition to duplexed retron
208 RT-DNA, the integrases may also bind and integrate prespacers consisting of one molecule of
209 RT-DNA as the (-) strand and one molecule of plasmid-derived ssDNA (the retron coding
210 sequence) as the (+) strand.

211 Beyond the apparent difference in prespacer processing due to the 2'-5' linkage, we
212 were curious to see whether the efficiency of acquisition would increase if the 2'-5' linkage was
213 removed. We approached this question by electroporating cells with three different prespacer
214 types: purified RT-DNA, purified and debranched RT-DNA, and a synthetic oligo version of the
215 RT-DNA. Debranched RT-DNA and oligos tended to be acquired more efficiently than the
216 natively-branched RT-DNA, but this trend did not reach statistical significance (Fig. 3e).

217 In some retrons, processing naturally occurs following reverse transcription to remove
218 the 2'-5' linkage³¹, so we tested such a retron to see whether this processing would change the
219 pattern or efficiency of retron-derived acquisitions. While the biosynthesis of the retron Eco4 RT-
220 DNA still depends on priming from the 2'-hydroxyl of a conserved guanosine, its RT-DNA is
221 cleaved 4 bases away from the 5' branch point by an ExoVII exonuclease complex-dependent
222 mechanism^{31,32}, leaving a mature RT-DNA lacking a 2'-5' linkage (Ext. Data Fig. 3a)³¹. We
223 expressed wildtype Eco4 ncRNA, Eco4 RT, and Cas1-Cas2 in cells and then sequenced their
224 CRISPR arrays to measure acquisitions. Notably, unlike the variant Eco1 retron, acquisitions
225 from Eco4 occurred in two different orientations (Fig. 3f, Ext. Data Fig. 3b-c). Although the
226 wildtype Eco4 RT-DNA does not have any perfect PAM sites (3'-TTC), both orientations
227 observed in Eco4-derived spacers had a near-perfect PAM (3'-GTC) which proved sufficient for
228 integration. We found no evidence as to whether these Eco4-derived spacers were derived from
229 single hairpins or duplexes. We next analyzed the mismatched regions of the Eco4-derived
230 spacers. As expected, almost all the LP regions mapped to the (-) strand, but unlike with variant
231 Eco1, the LD region of Eco4-derived spacers almost entirely mapped to the (+) strand (Fig. 3g).

232 Oligo-derived Eco4 spacers produced similar patterns of acquisition (Fig. 3g), indicating that
233 retron Eco4, and likely other unbranched RT-DNAs, avoid the peculiarities caused by using a
234 branched RT-DNA as a prespacer.

235 To confirm that Eco4 RT-DNA is debranched *in vivo*, we treated purified Eco4 RT-DNA
236 with DBR1 and did not observe a size shift that would indicate removal of ribonucleotides (Fig.
237 3h). In addition, the RT-DNA was not recJ sensitive because there were fewer than 6 bases of
238 single stranded DNA on the 5' end, which recJ requires for exonuclease activity³³.

239 The final test for Eco4 was to determine the overall efficiency of acquisition. We
240 observed that retron-derived spacers from Eco4 were dependent on the presence of Eco4 RT,
241 but their frequency was ultimately lower than Eco1-derived spacers (Fig. 3i). Based on these
242 baseline efficiencies, we have focused our efforts on engineering Eco1 for the purpose of
243 molecular recording. However, these results demonstrate that other retrons can also be used for
244 molecular recording and, as is the case with Eco4, may possess unique qualities which affect
245 their function in these applications.

246

247 **Temporal recordings of gene expression**

248 Having built and characterized the requisite tools, we set out to make a temporal
249 recording of gene expression using retron-based tags. We first constructed a signal plasmid and
250 a recording plasmid. The recording plasmid contained the coding sequence for retron Eco1 RT,
251 expressed from the constitutive promoter J23115, and the coding sequences for Cas1 and
252 Cas2, both under the control of a T7/lac promoter. The signal plasmid, pSBK.134, harbored two
253 copies of the Eco1 v35 ncRNA with different barcodes in the loop, which we will refer to as “A”
254 and “B”, under different inducible promoters. “A” was under the control of the
255 anhydrotetracycline-inducible promoter, pTet*, and ncRNA “B” was under the control of the
256 choline chloride-inducible promoter, pBetI (Fig. 4a)³⁴. We tested both the pTet* and pBetI
257 promoters individually using YFP fluorescence and confirmed that both are responsive to their

258 respective inducers with a similar maximum fluorescence, although the pBetI is 'leakier' with
259 higher uninduced fluorescence (Ext. Data Fig. 4). We found no effect on growth of harboring
260 these plasmids, and no effect on growth of adding inducers to cells that harbor the plasmids, but
261 a small effect of inducing the pTet* promoter versus cells that harbor no plasmids (Ext. Data Fig.
262 5). The recorded responses to induction of pTet* and pBetI were well matched, with 24 hours of
263 induction of each promoter yielding similar numbers of "A" and "B" derived spacers (Fig. 4b).

264 To record a time-ordered biological event, we transformed *E. coli* BL21-AI cells with both
265 the signal and recording plasmids, and grew them under two different experimental conditions
266 for a total of 48 hours. In the first temporal recording condition, cells were grown for 24 hours
267 with inducers driving the expression of Eco1 RT, Cas1-2, and ncRNA "A". The cells were then
268 grown for another 24 hours while expressing ncRNA "B", along with the Eco1 RT and Cas1-
269 Cas2 (Fig. 4c). In the second condition, the order of expression of ncRNA "A" and "B" was
270 reversed (ncRNA "B" was expressed for the first day and ncRNA "A" for the second) (Fig. 4d).
271 Samples were taken at 24 and 48 hours. Examination of the expanded arrays revealed a
272 significant increase in the percentage of cells that received a retron-derived spacer in the 24
273 hours where its chemical inducer was present, compared to the 24 hours where it was absent.
274 This held true for both ncRNAs "A" and "B" under both the "A"-before-"B" and "B"-before-"A"
275 expression schemes (Fig. 4c-d). The number of non-retron-derived spacers also increased
276 consistently over 48 hours (Fig. 4e).

277 To further test the generalizability of the system, we made a recording of a different set
278 of promoters driving the same retron ncRNAs. For this second arrangement, the recording
279 plasmid remained the same, but in the signal plasmid, pSBK.136, ncRNA "A" was placed under
280 the control of the sodium salicylate-inducible promoter, pSal, and "B" under the control of pTet*
281 (Fig. 4f). We also validated the pSal promoter individually using YFP fluorescence and found it
282 to be responsive to its inducer, with a higher maximum fluorescence than the pTet* promoter
283 (Ext. Data Fig. 4). Consistent with this difference, 24 hours of induction of each promoter

284 resulted in a much higher rate of acquisitions from retron “A” driven by pSal than acquisitions of
285 retron “B” from pTet* (Fig. 4g). In this case, induction of the pSal promoter did result in a
286 negative effect on population growth (Ext. Data Fig. 5). Notably, in one biological replicate, the
287 recording system appeared to break, resulting in nearly non-existent acquisitions; this sample
288 was excluded from further analysis following its identification as an outlier by Grubbs’ test^{35,36}
289 (Fig. 4g). Next, we tested two experimental conditions: “A”-before-“B” and “B”-before-“A” (Fig.
290 4h-i). Despite the mismatched promoter strengths, when arrays were examined at the 24- and
291 48-hour timepoints, more arrays were expanded with retron-derived spacers in the presence of
292 their respective inducers than in their absence (Fig. 4h-i). In addition, the numbers of non-
293 retron-derived spacers again increased over 48 hours (Fig. 4j).

294 This analysis of spacer acquisitions from the signal plasmid was enabled by a timepoint
295 sampling in the middle of the overall transcriptional sequence. However, the aim of this work is
296 to reconstruct the timing of transcriptional events using only data acquired at an endpoint.
297 Therefore, we defined logical rules that should govern the ordering of spacers in the CRISPR
298 arrays, and allow us to reconstruct the order of transcription of separate ncRNAs. Because
299 spacers are acquired unidirectionally, with newer spacers closer to the leader sequence, we
300 postulated that if transcript “A” is expressed before transcript “B”, arrays of the form “A” → “B” →
301 Leader should be more numerous than “B” → “A” → Leader. Accordingly, if “B” is expressed
302 before “A”, then the opposite should be true: the number of “B” → “A” → Leader arrays should
303 be greater than the number of “A” → “B” → Leader arrays.

304 Another feature of using CRISPR arrays for recording is that Cas1-2 also acquire
305 spacers derived from the plasmid and genome²⁰. These untargeted acquisitions can also be
306 used to interpret temporal information^{5,8}. If we assume that these non-retron-derived spacers
307 (denoted “N”) are acquired at a constant rate throughout the experiment, we can define a set of
308 rules that govern the order of “N” → “A” → Leader versus “A” → “N” → Leader arrays and of “N”
309 → “B” → Leader versus “B” → “N” → Leader arrays. In the “A”-before-“B” case, since “A” is

310 expressed in the first half of an experiment, arrays of the form “A” → “N” → Leader should be
311 more numerous than “N” → “A” → Leader. And since “B” is expressed in the second half of the
312 experiment, “N” → “B” → Leader arrays should be more numerous than “B” → “N” → Leader
313 arrays. Likewise, in the “B”-before-“A” condition, “N” → “A” → Leader arrays should be more
314 numerous than “A” → “N” → Leader arrays and “B” → “N” → Leader arrays should be more
315 numerous than “N” → “B” → Leader arrays. Restating these as mathematical statements, we
316 can take the difference between possible array types (e.g. “A” → “B” → Leader minus “B” → “A”
317 → Leader) as the numerator and the sum of the two possibilities (e.g. “A” → “B” → Leader plus
318 “B” → “A” → Leader) as the denominator (Fig. 4k) to yield a number between -1 and 1 for each
319 ordering rule (A/B, A/N, and B/N). By the convention of our ordering rules, positive values would
320 indicate that “A” was present before “B”, and a negative output would indicate that “B” was
321 present before “A”. The magnitude of the output ($0 \leq |x| \leq 1$) is a measure of how strongly the
322 rule is satisfied in a given direction, or in other words, how complete is the separation of the two
323 signals in time.

324 To test these predictions, we sequenced the CRISPR arrays of all of our samples at the
325 48-hour endpoint. Across 6 biological replicates of samples with signal plasmid pSBK.134, the
326 samples in which “A” was expressed before “B” yielded positive values when subjected to
327 analysis by our ordering rules, correctly identifying the order of expression. Likewise, for
328 samples where “B” was expressed before “A”, the rules yielded negative values, again correctly
329 identifying the order (Fig. 4l). We also calculated a composite score by taking a weighted
330 average of all three rules. This score consists of the average between the A/B rule and the sum
331 of the A/N rule and B/N rule. We devised this formulation based on what the ordering rules
332 represent in an ideal system. By definition, the A/B rule represents the degree of order between
333 A and B and will have a magnitude between 0 and 1. When “A” and “B” are not at all ordered
334 with respect to time the ordering score should be 0, and when “A” and “B” completely separated
335 in time the ordering score should be 1. Likewise, the A/N and B/N rules represent the degree of

336 order between “N” and “A” or “B”, respectively. In an ideal system, where the rate of acquisition
337 of “N” is constant, the magnitude of the A/N and B/N scores should be constrained between 0
338 and 0.5, and the sum of the A/N score and B/N score can be used as a proxy for the order of “A”
339 with respect to “B”. Thus, if we assume that the rate of acquisition of “N” is constant, we can
340 average the sum of the A/N and B/N scores with the A/B score to generate a composite score
341 which integrates all three rules and is representative of the degree of temporal order between A
342 and B. It is important to note that in the *in vivo* recording data, there are samples in which the
343 magnitudes of the A/N and/or B/N scores exceeds the ideal value of 0.5 and, as a result, the
344 composite score exceeds 1. We believe that this could occur due to several reasons. One
345 reason is that our assumption of “N” being a constant signal is not true *in vivo*, and that the
346 strength of signal “N” has some structure in time. Another potential reason is that the recording
347 of these signals is a stochastic process, with randomness and noise introduced at many levels
348 of the system, from RT-DNA synthesis, to spacer acquisition, to cell division, to sampling.

349 When applied to our *in vivo* recording data, this method accurately determined that each
350 experiment yielded directional acquisition of spacers and correctly recalled the order of events
351 for both directions. Critically, this demonstrates our ability to accurately reconstruct the order of
352 two transcriptional events in an endpoint biological sample, using only logical rules derived from
353 first principles. Interestingly, the retron signal driven by the pBetI promoter, which was found to
354 be leakier when uninduced, was not as strongly directional as the pTet^{*}-driven signal in relation
355 to N spacers, as would be expected. When each replicate was examined separately, though all
356 rules were not uniformly satisfied, the order of expression could be consistently determined (Ext.
357 Data Fig. 6a-b).

358 When this analysis was applied to samples with the signal plasmid pSBK.136 (which had
359 mismatched “A” and “B” promoter strengths), we were still able to accurately reconstruct the
360 order of events from endpoint data (Fig. 4m, Ext. Data Fig. 6c-d), demonstrating that the
361 temporal analysis of gene expression can be generalized to different promoters.

362 Finally, Retro-Cascorder data is stably maintained in cells for multiple generations after
363 the completion of a recording. When cells containing retron-derived recordings using signal
364 plasmid pSBK.134 were passaged for multiple days, ordering analysis results remained very
365 stable through roughly 18 generations of cell division (Ext. Data Fig. 7a-d). Only after around 45
366 generations of division did ordering scores begin to experience moderate drift, with severe drift
367 apparent after 81 generations. For reference, the Hayflick limit of human fetal cells in vitro is 40
368 to 60 generations of cell division³⁷. Ultimately, this paradigm enables the reconstruction of
369 temporal histories within genetically-identical populations of cells, based on a physical molecular
370 record.

371

372 **Modeling the Limits of Retron Recording**

373 To better understand the nature of the retron recording system and its behavior in a wide
374 range of conditions, including those which we are unable to recreate in the lab, we developed a
375 computational model of the Retro-Cascorder based on data from our temporal recordings and
376 present understanding of the biology of the system. Using the raw number of acquisitions
377 observed previously from temporal recordings, at the 24- and 48-hour timepoints, we defined a
378 set of rates (of acquisitions per hour) for the different signals recorded by the cells. Based on
379 the overall low rate of acquisitions, we assume that the spacer acquisitions can be modeled
380 faithfully as a Poisson process, wherein the average time between events (here acquisitions) is
381 known but the exact timing of events is random. Using rates estimated from our recordings, we
382 define rates of acquisition for A and B in the presence of their respective inducers, A and B in
383 the absence of their respective inducers (to account for leak from their promoters), and a
384 constant background rate of acquisition of non-retron-derived spacers, or N.

385 To test our model, we first simulated 100 replicates, of 1 million arrays each, of our
386 previous recording experiments using the signal plasmids pSBK.134 and pSBK.136 (Fig. 5a-b).
387 With the results of the simulation appearing to approximate the results from our real recordings,

388 we next sought to understand how changing various parameters of the recording system may
389 affect results. First, we simulated the effect of analyzing different numbers of arrays (Fig. 5c).
390 The simulation suggests that it is important to dedicate a generous number of sequencing reads
391 to a recording experiment in order to properly resolve the process in question. When too few
392 arrays are analyzed from a given sample, the calculated ordering scores will be unreliable, as
393 evidenced by the very wide distribution of composite scores from simulated low-read samples.

394 Next, we simulated making recordings which varied in length (Fig. 5d). The effect that
395 appears here is similar to the effect of varying the number of arrays analyzed. In the range of
396 very short recordings, the system is unable to resolve the order of the signals, but as the length
397 of the recordings increases, the composite scores converge toward a specific value. To check
398 this finding from the model, we made biological recordings of different length using pSBK.134.
399 In these recordings, the trend predicted by the model appeared to hold, with shorter recordings
400 unable to resolve the order of the signals and longer recordings with greater fidelity. Finally, we
401 simulated the effect of varying the rates of acquisition of both retrans. To do this we simulated
402 50 replicates, of 1 million arrays each, across a range of rates of acquisition of both signals A
403 and B. We varied both the induced and uninduced rates of a given signal by the same factor
404 (e.g. A-On and A-Off increased by a factor of 4, B-On and B-Off decreased by a factor of 8).
405 Interestingly, even when acquisition rates are decreased, the mean ordering scores across 50
406 replicates faithfully reflect the order of expression of signals (Fig. 5e-g). Dispersion of the
407 ordering scores among replicates, however, varies dramatically with the rates of acquisition of
408 signals A and B. In short: as the strength of signals A and B increases, we expect to be able to
409 faithfully recall temporal order using fewer replicates, or visa-versa, that if the strength of signals
410 decreases, more replicates will be required to resolve their temporal order.

411 Putting together the pieces above, we believe that we have shown 4 variables that are
412 critical to the design of these recording experiments: (1) signal strength; (2) length of recording;
413 (3) number of reads; and (4) number of replicates. By increasing any of these 4 parameters, the

414 experimenter can expect greater fidelity of their final temporal recording. When these
415 parameters are decreased, one should expect more noise and variability in their recordings. In
416 the laboratory however, there will be practical limits as to how much the experimenter can
417 maximize or alter these parameters. Often, it may not be possible to alter the duration of an
418 experiment or the strength of a transcriptional signal due to the biology of the process of
419 interest, and it may be time- and cost-prohibitive to run large numbers of biological replicates. Of
420 the four parameters then, increasing the number of arrays analyzed (and consequently the
421 number of sequencing reads) from individual samples is likely the cheapest and simplest way to
422 increase the fidelity of temporal recordings and final analyses.

423

424 **DISCUSSION**

425 Here, we have described a system for the recording and reconstruction of transcriptional
426 history in a population of cells, which we call the Retro-Cascorder. We achieved this by
427 engineering an RNA molecular tag, which is specifically reverse-transcribed to produce a DNA
428 ‘receipt’ of transcription that is permanently saved in a CRISPR array. We demonstrated the
429 flexibility and potential for continued development of these tools by making the recording retron
430 more efficient with modifications to the structure of the retron ncRNA, and developed a toolkit of
431 barcoded retrons for future application to more complex systems. Beyond this, we investigated
432 the ability of the CRISPR adaptation system to utilize RT-DNA as a prespacer, and discovered
433 that the retron 2’-5’ linkage causes a marked difference in the type of spacers acquired. Finally,
434 we used this system to record and reconstruct time-ordered biological events in populations of
435 cells, and developed a computational model of retron-recording to more comprehensively
436 explore the limits of the system.

437 One natural aspect of the CRISPR integrases which has proven useful in these
438 recordings is the acquisition of diverse spacers from plasmid and genomic fragments (N
439 spacers). In our temporal recordings of two inducible elements, these N spacers function as a

440 third signal, providing a constant background. Because the integrases are also driven by an
441 inducible promoter, this background signal marks the timing of the recording components. In our
442 recordings with pSBK.134, for instance, the A spacers encode the timing of anhydrotetracycline
443 in the media, the B spacers encode the timing of choline chloride in the media, and the N spacer
444 encode the timing of arabinose and IPTG in the media. The frequency of N spacer acquisitions
445 is unaffected by the retron-derived acquisitions, which we interpret to mean that the acquisition
446 of events in this system is not competitive, but rather additive. Therefore, these N spacers do
447 not interfere with the recording, but rather aid in resolving the temporal order of recorded
448 signals.

449 Here, we recorded two distinct signals within a homogenous population of cells, with the
450 N spacers serving as a constant third signal. The level of complexity of these recordings is
451 similar to previous work using different recombinases to encode events^{7,27,38}. One aspect that is
452 encouraging about the system described here is that the recordings use a common set of
453 protein components, with distinct signals being encoded using variable nucleotides in the retron-
454 derived spacers. Recombinase-based recorders, while robust, are inherently limited in the
455 number of distinct signals to $2^{\text{number of recombinases}}$ ³⁹, which requires identifying and
456 expressing many orthogonal recombinases. In contrast, this approach is limited in the number of
457 distinct signals to $4^{\text{number of nucleotides used in the barcode}}$. This bodes well for the
458 scalability of this approach, but the practicality of scaling will need to be experimentally validated
459 in future work.

460 We believe that this framework of selective tagging and recording of biological signals in
461 an RNA \rightarrow DNA \rightarrow CRISPR direction is a powerful, modular, and extensible method of making
462 temporal recordings in cells. Using only *a priori* ordering rules, we can detect and interpret time-
463 ordered biological signals from a single endpoint sample. Immediate uses of this technology
464 include the construction of living biosensors that sample and record their environment. Here, we
465 record the presence of anhydrotetracycline, choline chloride, sodium salicylate, arabinose, and

466 IPTG. Near future work could modify these systems to record the presence of pollutants,
467 metabolites, or pathogens in an environment. With additional engineering to increase the
468 efficiency of the recordings, we hope that this system will enable recordings of natural gene
469 expression to log transcriptional order during complex cellular events.

470

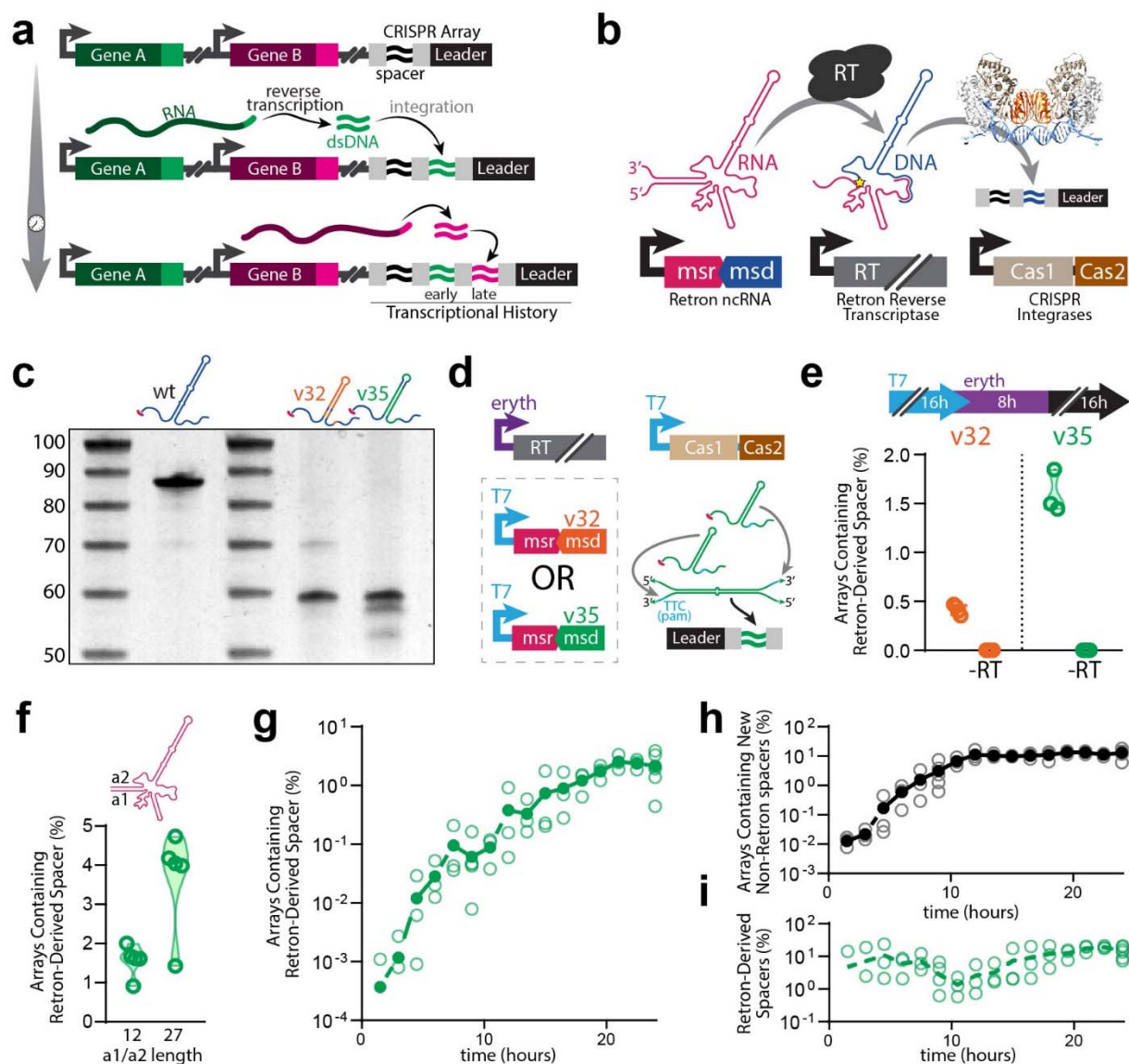
471 **Main References:**

- 472
- 473 1 Simon, A. J., Ellington, A. D. & Finkelstein, I. J. Retrons and their applications in genome
474 engineering. *Nucleic Acids Res* **47**, 11007-11019, doi:10.1093/nar/gkz865 (2019).
- 475 2 Barrangou, R. *et al.* CRISPR provides acquired resistance against viruses in prokaryotes. *Science*
476 **315**, 1709-1712, doi:10.1126/science.1138140 (2007).
- 477 3 Church, G. M., Gao, Y. & Kosuri, S. Next-Generation Digital Information Storage in DNA. *Science*
478 **337**, 1628-1628, doi:10.1126/science.1226355 (2012).
- 479 4 Shipman, S. L., Nivala, J., Macklis, J. D. & Church, G. M. CRISPR–Cas encoding of a digital movie
480 into the genomes of a population of living bacteria. *Nature* **547**, 345-349,
481 doi:10.1038/nature23017 (2017).
- 482 5 Yim, S. S. *et al.* Robust direct digital-to-biological data storage in living cells. *Nat Chem Biol* **17**,
483 246-253, doi:10.1038/s41589-020-00711-4 (2021).
- 484 6 Ceze, L., Nivala, J. & Strauss, K. Molecular digital data storage using DNA. *Nat Rev Genet* **20**, 456-
485 466, doi:10.1038/s41576-019-0125-3 (2019).
- 486 7 Roquet, N., Soleimany, A. P., Ferris, A. C., Aaronson, S. & Lu, T. K. Synthetic recombinase-based
487 state machines in living cells. *Science* **353**, doi:10.1126/science.aad8559 (2016).
- 488 8 Sheth, R. U., Yim, S. S., Wu, F. L. & Wang, H. H. Multiplex recording of cellular events over time
489 on CRISPR biological tape. *Science* **358**, 1457-1461, doi:10.1126/science.aao0958 (2017).
- 490 9 Schmidt, F., Cherepkova, M. Y. & Platt, R. J. Transcriptional recording by CRISPR spacer
491 acquisition from RNA. *Nature* **562**, 380-385, doi:10.1038/s41586-018-0569-1 (2018).
- 492 10 Wagner, D. E. & Klein, A. M. Lineage tracing meets single-cell omics: opportunities and
493 challenges. *Nat Rev Genet* **21**, 410-427, doi:10.1038/s41576-020-0223-2 (2020).
- 494 11 Street, K. *et al.* Slingshot: cell lineage and pseudotime inference for single-cell transcriptomics.
495 *BMC Genomics* **19**, 477, doi:10.1186/s12864-018-4772-0 (2018).
- 496 12 Perli, S. D., Cui, C. H. & Lu, T. K. Continuous genetic recording with self-targeting CRISPR-Cas in
497 human cells. *Science* **353**, aag0511-aag0511, doi:10.1126/science.aag0511 (2016).
- 498 13 Park, J. *et al.* Recording of elapsed time and temporal information about biological events using
499 Cas9. *Cell* **184**, 1047-1063.e1023, doi:10.1016/j.cell.2021.01.014 (2021).
- 500 14 Shipman, S. L., Nivala, J., Macklis, J. D. & Church, G. M. Molecular recordings by directed CRISPR
501 spacer acquisition. *Science* **353**, aaf1175, doi:10.1126/science.aaf1175 (2016).
- 502 15 Simon, A. J., Morrow, B. R. & Ellington, A. D. Retroelement-Based Genome Editing and Evolution.
503 *ACS Synth. Biol.* **7**, 2600-2611, doi:10.1021/acssynbio.8b00273 (2018).
- 504 16 Sharon, E. *et al.* Functional Genetic Variants Revealed by Massively Parallel Precise Genome
505 Editing. *Cell* **175**, 544-557.e516, doi:10.1016/j.cell.2018.08.057 (2018).
- 506 17 Schubert, M. G. *et al.* High-throughput functional variant screens via in vivo production of single-
507 stranded DNA. *PNAS* **118**, e2018181118, doi:10.1073/pnas.2018181118 (2021).
- 508 18 Lopez, S. C., Crawford, K. D., Lear, S. K., Bhattarai-Kline, S. & Shipman, S. L. Precise genome
509 editing across kingdoms of life using retron-derived DNA. *Nat Chem Biol* **18**, 199-206,
510 doi:10.1038/s41589-021-00927-y (2022).
- 511 19 Farzadfard, F. & Lu, T. K. Genomically encoded analog memory with precise in vivo DNA writing
512 in living cell populations. *Science* **346**, 1256272-1256272, doi:10.1126/science.1256272 (2014).
- 513 20 Yosef, I., Goren, M. G. & Qimron, U. Proteins and DNA elements essential for the CRISPR
514 adaptation process in *Escherichia coli*. *Nucleic Acids Res* **40**, 5569-5576, doi:10.1093/nar/gks216
515 (2012).
- 516 21 Nuñez, J. K. *et al.* Cas1–Cas2 complex formation mediates spacer acquisition during CRISPR–Cas
517 adaptive immunity. *Nat Struct Mol Biol* **21**, 528-534, doi:10.1038/nsmb.2820 (2014).

- 518 22 Wang, J. *et al.* Structural and Mechanistic Basis of PAM-Dependent Spacer Acquisition in CRISPR-
519 Cas Systems. *Cell* **163**, 840-853, doi:10.1016/j.cell.2015.10.008 (2015).
- 520 23 Millman, A. *et al.* Bacterial Retrons Function In Anti-Phage Defense. *Cell* **183**, 1551-1561.e1512,
521 doi:10.1016/j.cell.2020.09.065 (2020).
- 522 24 Bobonis, J. *et al.* Bacterial retrons encode tripartite toxin/antitoxin systems. (*Microbiology*,
523 2020).
- 524 25 Lampson, B. C. *et al.* Reverse transcriptase in a clinical strain of *Escherichia coli*: production of
525 branched RNA-linked msDNA. *Science* **243**, 1033-1038, doi:10.1126/science.2466332 (1989).
- 526 26 Silas, S. *et al.* Direct CRISPR spacer acquisition from RNA by a natural reverse transcriptase–Cas1
527 fusion protein. *Science* **351**, aad4234, doi:10.1126/science.aad4234 (2016).
- 528 27 Bonnet, J., Subsoontorn, P. & Endy, D. Rewritable digital data storage in live cells via engineered
529 control of recombination directionality. *PNAS* **109**, 8884-8889, doi:10.1073/pnas.1202344109
530 (2012).
- 531 28 Kim, S. *et al.* Selective loading and processing of prespacers for precise CRISPR adaptation.
532 *Nature* **579**, 141-145, doi:10.1038/s41586-020-2018-1 (2020).
- 533 29 Ramachandran, A., Summerville, L., Learn, B. A., DeBell, L. & Bailey, S. Processing and integration
534 of functionally oriented prespacers in the *Escherichia coli* CRISPR system depends on bacterial
535 host exonucleases. *J. Biol. Chem.* **295**, 3403-3414, doi:10.1074/jbc.RA119.012196 (2020).
- 536 30 Chapman, K. B. & Boeke, J. D. Isolation and characterization of the gene encoding yeast
537 debranching enzyme. *Cell* **65**, 483-492, doi:10.1016/0092-8674(91)90466-C (1991).
- 538 31 Lim, D. Structure and biosynthesis of unbranched multicopy single-stranded DNA by reverse
539 transcriptase in a clinical *Escherichia coli* isolate. *Molecular Microbiology* **6**, 3531-3542,
540 doi:10.1111/j.1365-2958.1992.tb01788.x (1992).
- 541 32 Jung, H., Liang, J., Jung, Y. & Lim, D. Characterization of cell death in *Escherichia coli* mediated by
542 XseA, a large subunit of exonuclease VII. *J Microbiol.* **53**, 820-828, doi:10.1007/s12275-015-
543 5304-0 (2015).
- 544 33 Han, E. S. *et al.* RecJ exonuclease: substrates, products and interaction with SSB. *Nucleic Acids*
545 *Res* **34**, 1084-1091, doi:10.1093/nar/gkj503 (2006).
- 546 34 Meyer, A. J., Segall-Shapiro, T. H., Glassey, E., Zhang, J. & Voigt, C. A. *Escherichia coli*
547 “Marionette” strains with 12 highly optimized small-molecule sensors. *Nat Chem Biol* **15**, 196-
548 204, doi:10.1038/s41589-018-0168-3 (2019).
- 549 35 Grubbs, F. E. Procedures for Detecting Outlying Observations in Samples. *Technometrics* **11**, 1-21
550 (1969).
- 551 36 Stefansky, W. Rejecting Outliers in Factorial Designs. *Technometrics* **14**, 469-479 (1972).
- 552 37 Hayflick, L. & Moorhead, P. S. The serial cultivation of human diploid cell strains. *Experimental*
553 *Cell Research* **25**, 585-621, doi:10.1016/0014-4827(61)90192-6 (1961).
- 554 38 Yang, L. *et al.* Permanent genetic memory with >1-byte capacity. *Nat Methods* **11**, 1261-1266,
555 doi:10.1038/nmeth.3147 (2014).
- 556 39 Yehl, K. & Lu, T. Scaling computation and memory in living cells. *Current Opinion in Biomedical*
557 *Engineering* **4**, 143-151, doi:10.1016/j.cobme.2017.10.003 (2017).

558

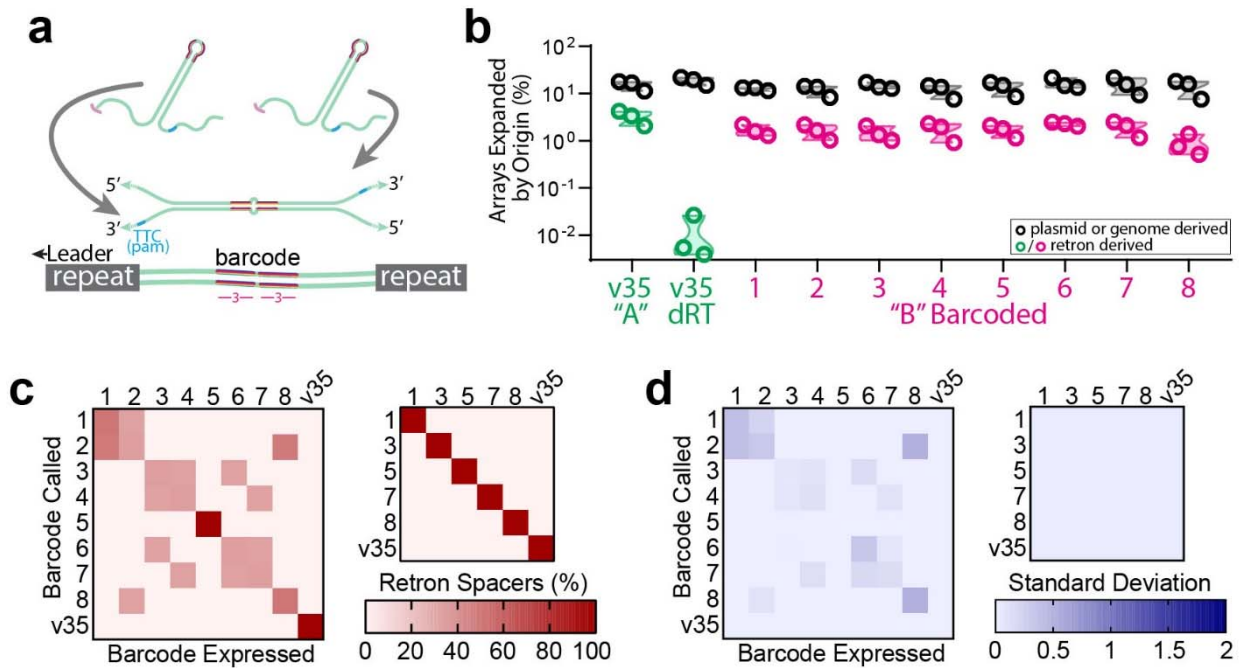
559 **Figures**



560

561

562 **Figure 1. Cas1-Cas2 integrates retron RT-DNA.** **a.** Schematic representation of retroelement-based
 563 transcriptional recording into CRISPR arrays. **b.** Schematic representation of biological components of
 564 the retron-based recorder. **c.** Urea-PAGE visualization of RT-DNA from retron Eco1 ncRNA variants.
 565 From left to right (excluding ladders): wild-type Eco1, Eco1 v32, Eco1 v35. For gel source data, see
 566 Supplementary Figure 1. **d.** Schematic of experimental promoters used to test retron-recorder parts and
 567 cartoon of hypothetical duplex RT-DNA prespacer structure. **e.** Quantification of arrays expanded with
 568 retron-derived spacers using Eco1 variants v32 (orange) and v35 (green). Open circles represent 3
 569 biological replicates. **f.** Quantification of arrays expanded with retron derived spacers with a wild-type (12
 570 bp) and extended (27 bp) a1/a2 region. Open circles represent 5 biological replicates. **g.** Time series of
 571 array expansions from retron-derived spacers. Open circles represent biological replicates, closed circles
 572 are the mean. **h.** Time series of array expansions from non-retron-derived spacers. Open circles
 573 represent biological replicates, closed circles are the mean. **i.** Proportion of total new spacers that are
 574 retron-derived. Open circles represent biological replicates, dashed line is the mean. All statistics in
 575 Supplementary Table 1.



576

577

578 **Figure 2. Diversification of retron-based barcodes.** **a.** Hypothetical structure of duplexed RT-DNA
 579 prespacer with 6-base barcode and retron-derived spacer. **b.** Quantification of array expansions from
 580 barcoded variants of retron Eco1 v35, showing both retron-derived (green/pink) and non-retron derived
 581 (black) spacers for each variant. Open circles represent 3 biological replicates. **c.** Left: Heatmap of *in*
 582 *silico* ability to distinguish between all barcoded Eco1 v35 variants. Right: Heatmap of *in silico* ability to
 583 distinguish between reduced set of barcoded Eco1 v35 variants. **d.** Heatmap of standard deviation
 584 between three separate trials of barcode discrimination test. Left: full set. Right: reduced set. All statistics
 585 in Supplementary Table 1.

586

587

588

589

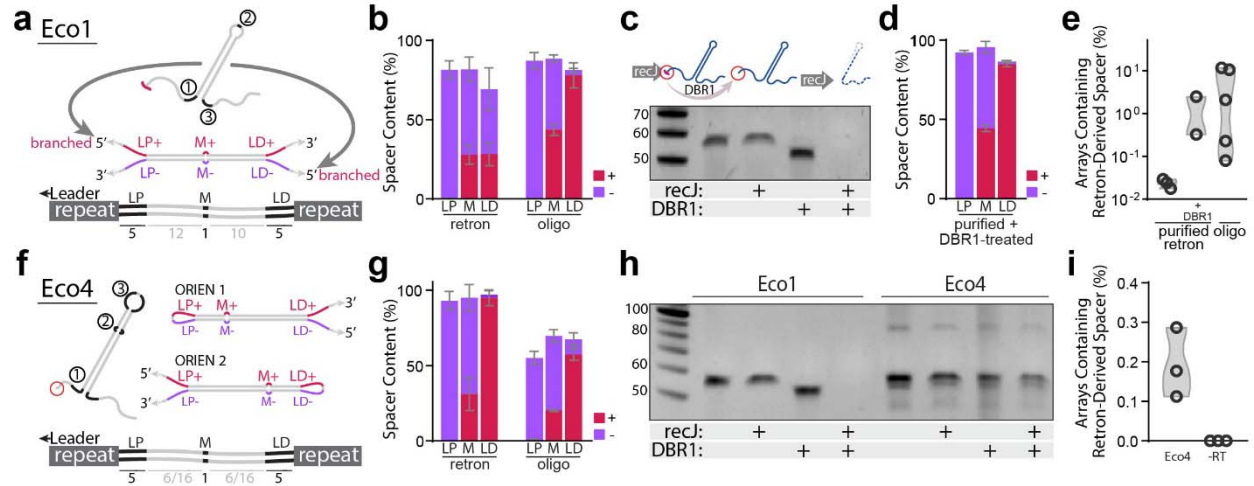
590

591

592

593

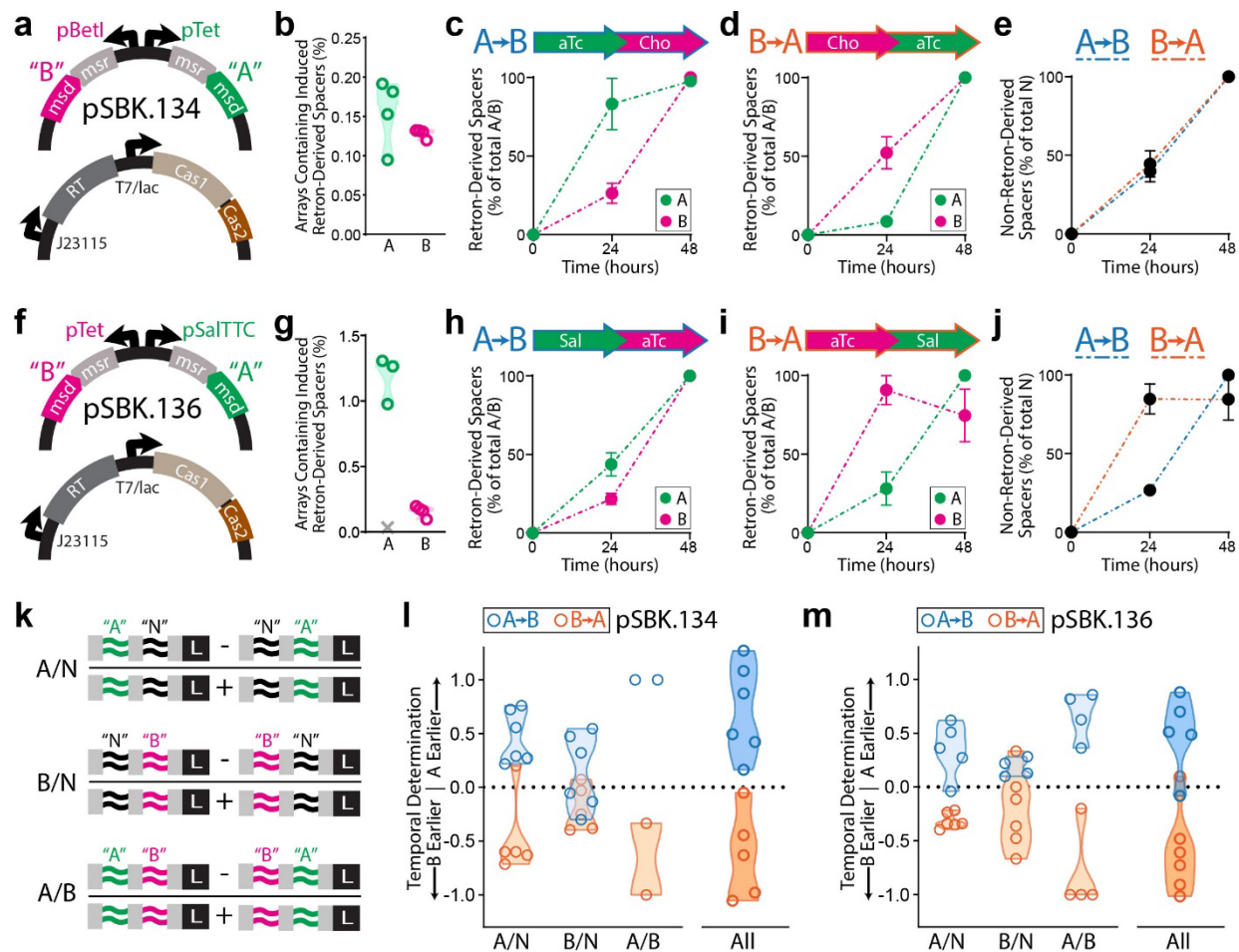
594



595
596

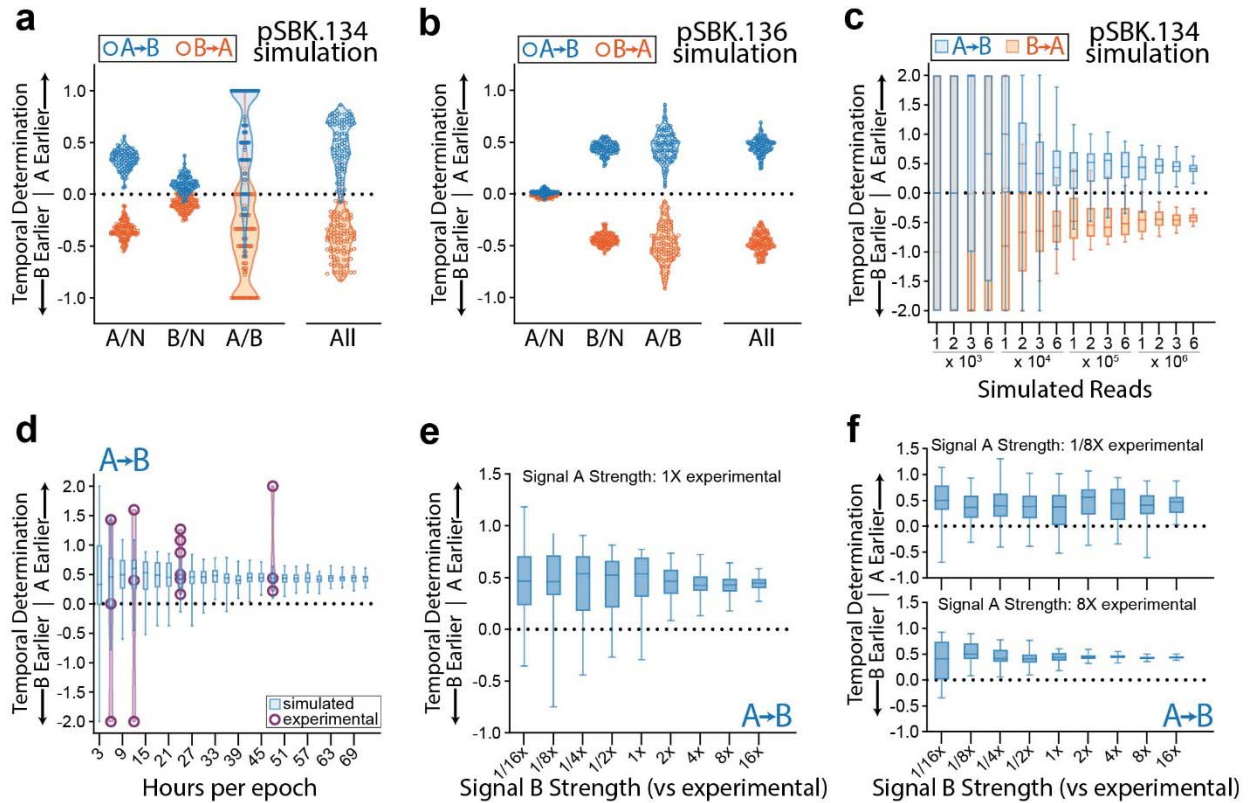
Figure 3. Mechanism of RT-DNA spacer acquisition. **a.** Hypothetical structure of duplexed Eco1 v35 RT-DNA prespacer and retron-derived spacer, with mismatched regions highlighted. **b.** Quantification of mismatch region sequences in spacers from cells expressing Eco1 v35 versus cells electroporated with oligo mimic. Bars represent the mean of 4 and 5 biological replicates for the retron and oligo-derived conditions, respectively (\pm SD). **c.** Urea-PAGE visualization of Eco1 RT-DNA. DBR1 treatment resolves 2'-5' linkage. For gel source data, see Supplementary Figure 1. **d.** Quantification of mismatch region sequences in spacers from cells electroporated with purified, debranched Eco1 v35 RT-DNA. Bars represent the mean of 4 biological replicates (\pm SD). **e.** Quantification of array expansions from different prespacer substrates. Open circles represent 3, 2, and 5 biological replicates (left-right). **f.** Schematic of Eco4 RT-DNA, in both orientations, with mismatch sequences highlighted. **g.** Quantification of mismatch region sequences in cells expressing Eco4 versus cells electroporated with oligo mimic. Bars represent the mean of 3 biological replicates (\pm SD). **h.** Urea-PAGE visualization of Eco4 RT-DNA. DBR1 does not cause size shift of Eco4 RT-DNA. For gel source data, see Supplementary Figure 1. **i.** Quantification of array expansions from retron Eco4. Open circles represent 3 biological replicates (left-right). All statistics in Supplementary Table 1.

611
612
613



614
 615 **Figure 4. Temporal recordings of gene expression.** **a.** Schematic of signal plasmid pSBK.134 used to
 616 express ncRNAs "A" and "B", and recording plasmid used to express Eco1-RT and Cas1 and 2. **b.**
 617 Accumulation of retron-derived spacers from pSBK.134 after 24 hours of induction from their respective
 618 promoters (4 biological replicates). **c.** Retron-derived spacers when ncRNAs were induced in the order
 619 "A" then "B" from pSBK.134. Filled circles represent the mean of four biological replicates (\pm SEM). **d.**
 620 Retron-derived spacers when ncRNAs were induced in the order "B" then "A" from pSBK.134. Filled circles
 621 represent the mean of four biological replicates (\pm SEM). **e.** Non-retron-derived spacers in cells
 622 harboring pSBK.134, in both induction conditions. Filled circles represent the mean of four biological
 623 replicates (\pm SEM). **f.** Schematic of signal plasmid pSBK.136 used to express ncRNAs "A" and "B", and
 624 the recording plasmid. **g.** Accumulation of retron-derived spacers from pSBK.136 after 24 hours of
 625 induction from their respective promoters (4 biological replicates). Outlier sample determined by Grubbs'
 626 test denoted as a grey "X". **h.** Retron-derived spacers when ncRNAs were induced in the order "A" then
 627 "B" from pSBK.136. Filled circles represent the mean of three biological replicates (\pm SEM). **i.** Retron-
 628 derived spacers when ncRNAs were induced in the order "B" then "A" from pSBK.136. Filled circles
 629 represent the mean of four biological replicates (\pm SEM). **j.** Non-retron-derived spacers in cells harboring
 630 pSBK.136, in both induction conditions. Filled circles represent the mean of four biological replicates
 631 (\pm SEM). **k.** Graphical representation of the rules used to determine order of expression from arrays. **l.**
 632 Ordering analysis of recording experiments with signal plasmid pSBK.134. Open circles are 6 biological
 633 replicates. **m.** Ordering analysis of recording experiments with signal plasmid pSBK.136. Open circles are
 634 5 biological replicates. All statistics in Supplementary Table 1.

635
 636
 637



638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656

Figure 5. Modeling the Limits of Retron Recording. **a.** Simulation of 100 replicates each of A-then-B and B-then-A recordings using acquisition rate data from pSBK.134 recordings. Each point represents the calculated ordering score from a single replicate of 1 million arrays. **b.** Simulation of 100 replicates each of A-then-B and B-then-A recordings using acquisition rate data from pSBK.136 recordings. Each point represents the calculated ordering score from a single replicate of 1 million arrays. **c.** Simulation of varying the number of arrays analyzed per sample using acquisition rate data from pSBK.134 recordings. Each box with whiskers represents 100 simulated replicates, with whiskers extending from minimum to maximum. **d.** Simulation of varying the length of each epoch in a retron recording using acquisition rate data from pSBK.134 (blue). Overlaid with real retron recordings of the same length (purple). Each box with whiskers represents 100 simulated replicates of 1 million reads each, with whiskers spanning from minimum to maximum. Each overlaid point is a single biological replicate. Recording experiments with 6, 12, and 48-hour epochs were done in triplicate. Recording experiments with epoch length of 24 hours are the same as in Figure 4l. **e.** Simulation of varying the strength of signal B when signal A remains constant. 1x acquisition rates were obtained from pSBK.134 recordings. Each box with whiskers represents 50 simulated replicates of 1 million arrays each. Whiskers span from minimum to maximum. **f.** Simulation of varying the strength of signal B when signal A is decreased or increased by a factor of 8. 1x acquisition rates were obtained from pSBK.134 recordings. Each box with whiskers represents 50 simulated replicates of 1 million arrays each. Whiskers span from minimum to maximum.

657 **METHODS**

658 All biological replicates were taken from distinct samples, not the same sample measured
659 repeatedly.

660 **Bacterial Strains and Growth Conditions**

661 This work uses the following *E. coli* strains: NEB 5-alpha (NEB C2987, not authenticated),
662 BL21-AI (ThermoFisher C607003, not authenticated), bMS.346, and bSLS.114. bMS.346 was
663 generated from *E. coli* MG1655 by inactivating *exol* and *recJ* genes with early stop codons as in
664 previous work⁴⁰. Additionally, the *araB::T7RNAP-tetA* locus was transferred from BL21-AI by P1
665 phage transduction⁴¹. bSLS.114 (which has been used previously¹⁸) was generated from BL21-
666 AI by deleting the retron Eco1 locus by lambda Red recombinase mediated insertion of an FRT-
667 flanked chloramphenicol resistance cassette. This cassette was amplified from pKD3⁴² with
668 homology arms added to the retron Eco1 locus. This amplicon was electroporated into BL21-AI
669 cells expressing lambda Red genes from pKD46⁴², and clones were isolated by selection on
670 chloramphenicol (10 µg/mL) plates. After genotyping to confirm locus-specific insertion, the
671 chloramphenicol cassette was excised by transient expression of FLP recombinase to leave
672 only an FRT scar. Experimental cultures were grown with shaking in LB broth at 37°C with
673 appropriate inducers and antibiotics. Inducers and antibiotics were used at the following working
674 concentrations: 2 mg/mL L-arabinose (GoldBio A-300), 1 mM IPTG (GoldBio I2481C), 400 µM
675 erythromycin, 100 ng/mL anhydrotetracycline, 100 µM choline chloride, 1 mM sodium salicylate,
676 35 µg/mL kanamycin (GoldBio K-120), 25 µg/mL spectinomycin (GoldBio S-140), 100 µg/mL
677 carbenicillin (GoldBio C-103), 25 µg/mL chloramphenicol (GoldBio C-105; used at 10 µg/mL for
678 selection during recombineering). Additional strain information can be found in Supplementary
679 Table 2.

680 **Plasmid Construction**

681 All cloning steps were performed in *E. coli* NEB 5-alpha. pWUR 1+2, containing Cas1 and Cas2
682 under the expression of a T7lac promoter, was a generous gift from Udi Qimron²⁰. Eco1

683 wildtype ncRNA and Eco1 RT, along with Cas1+2, were cloned into pRSF-DUET (Sigma
684 71341) to generate pSLS.405. Eco1 variant ncRNA sequences v32 and v35 were cloned into
685 pRSF-DUET along with Cas1+2 to generate pSLS.407 and pSLS.408, respectively. Extended
686 a1/a2 v35 ncRNA expression plasmid pSLS.416 was generated from pSLS.408 by site-directed
687 mutagenesis. Retron Eco1 RT and retron Eco4 RT were cloned into pJKR-O-mphR to generate
688 pSLS.402 and pSLS.400, respectively. pJKR-O-mphR was generated previously⁴³ (Addgene
689 plasmid # 62570). Barcoded, extended a1/a2 v35 ncRNA expression plasmids pSBK.009-016
690 were generated from pSLS.416 by site-directed mutagenesis. Wildtype retron Eco4 ncRNA was
691 cloned into pRSF-DUET along with Cas1+2 to generate SLS.419. pSBK.134 and pSBK.136
692 were generated in three steps. First, barcoded, extended a1/a2 v35 ncRNA sequences were
693 cloned into the 'Marionette' plasmids pAJM.717, pAJM.718, and pAJM.771. pAJM.717,
694 pAJM.718, and pAJM.771 were gifts from Christopher Voigt³⁴ (pAJM.717 - Addgene plasmid #
695 108517 // pAJM.718 - Addgene plasmid # 108519 // pAMJ.771 - Addgene plasmid # 108534).
696 Then, in two steps, two ncRNA expression cassettes (for barcoded ncRNAs "A" and "B") from
697 the Marionette plasmids were cloned into pSol-TSF (Lucigen F843213-1) facing in opposite
698 directions. pSBK.079 was generated by cloning the resistance marker AmpR in place of the
699 KanR marker into the plasmid pSLS.425, which was synthesized by Twist biosciences.
700 Additional plasmid information can be found in Supplementary Table 3.

701 **RT-DNA Purification and Visualization**

702 Retron RT-DNA was expressed in *E. coli* BMS.346 and purified in two steps. First, DNA was
703 extracted from cells using a plasmid midiprep kit (Qiagen 12943). This purified DNA was then
704 treated for 30 minutes at 37C with RNase A/T1 mix (ThermoFisher EN0551) and, if required,
705 DBR1 (OriGene TP300024) and/or RecJ_f (NEB M0264). This sample was then used as the
706 input for the Zymo Research ssDNA/RNA Clean & Concentrate kit (Zymo D7011). Samples
707 eluted from the ssDNA kit were resolved using TBE-urea PAGE (ThermoFisher EC6885BOX).

708 Gels were stained with SYBR Gold for imaging (ThermoFisher S11494) and imaged on a Bio-
709 Rad Gel Doc imager.

710 **Retron Acquisition Experiments**

711 Cells were transformed sequentially: first with the RT expression plasmid (pSLS.400 or
712 pSLS.402), and second with the ncRNA and Cas1+2 expression plasmid (eg. pSLS.416). For
713 the -RT condition, cells were only transformed with an ncRNA and Cas1+2 expression plasmid
714 (e.g. pSLS.416). For testing acquisition of retron-derived spacers in figures 1e-f, cells with RT,
715 ncRNA, and Cas1+2 expression plasmids were grown overnight (16 hours) in 3 mL LB with
716 antibiotics and inducers IPTG and arabinose, from individual clones on plates. In the morning,
717 240 uL of overnight culture was diluted into 3 mL fresh media with antibiotics, IPTG, and
718 arabinose and grown for 2 hours. After 2 hours, 320 uL of culture was diluted into 3 mL fresh
719 media with antibiotics and erythromycin (no erythromycin was used in the -RT condition) and
720 grown for 8 hours. After 8 hours, culture was diluted 1:1000 into 3 mL LB with antibiotics and
721 without inducers and grown overnight (16 hours). In the morning, 25 uL of culture was mixed
722 with 25 uL of water, heated to 95C for 5 minutes to lyse cells, cooled, and frozen at -20C for
723 later analysis. For data presented in Figures 2b-d and 3i, cells were grown overnight (16 hours)
724 in 3 mL LB with antibiotics and inducers IPTG and arabinose, from individual clones on plates.
725 In the morning, 240 uL of overnight culture was diluted into 3 mL fresh media with antibiotics,
726 IPTG, and arabinose and grown for 2 hours. After 2 hours, 320 uL of culture was diluted into 3
727 mL fresh media with antibiotics and erythromycin and grown for 2 (rather than 8) hours. At this
728 point, 25 uL of culture was mixed with 25 uL of water, heated to 95C for 5 minutes to lyse cells,
729 cooled, and frozen at -20C for later analysis.

730 For the 24-hour time course experiment, the experiment was broken into two halves: the first 9
731 hours, and the final 15 hours. For the entirety of the time course, cells were grown in media with
732 antibiotics and inducers (arabinose, IPTG, and erythromycin). For the first 9-hour samples,
733 cultures were grown starting from single colonies added to 0.5 mL of media. These cultures

734 were sampled every 1.5 hours until hour 9, with 1 mL of media added at hour 3 and 1.5 mL of
735 media added at hour 6. For the final 15-hour samples, 3 mL of media was inoculated with single
736 colonies from plates and grown for 9 hours. Starting at hour 9, samples were taken every 1.5
737 hours until hour 24. At hour 16.5, 200 uL of culture was diluted into 1.5 mL of fresh media and
738 the experiment continued in the new tube. At hour 21, 1 mL media was added to the culture.

739 **Oligo Prespacer Feeding**

740 For spacer acquisition experiments using exogenous DNA prespacers (purified RT-DNA or
741 synthetic oligos), cells containing pWUR1+2 were grown overnight from individual colonies on
742 plates. In the morning, 100 uL of overnight culture was diluted into 3 mL LB with antibiotics,
743 IPTG, and arabinose. Cells were grown with inducers for 2 hours. For each electroporation, 1
744 mL of culture was pelleted and resuspended in water. Cells were washed a second time by
745 pelleting and resuspension, then pelleted one final time and resuspended in 50 uL of prespacer
746 DNA solution at a concentration of 6.25 uM of single-stranded RT-DNA. All wash steps were
747 done using ice cold water, all centrifugation steps were done in a centrifuge chilled to 4C, and
748 samples kept on ice until electroporation was complete. The cell-DNA mixture was transferred
749 to a 1 mm gap cuvette (Bio-Rad 1652089) and electroporated using a Bio-Rad gene pulser set
750 to 1.8 kV and 25 uF with pulse controller at 200 Ohms. After electroporation, cells were
751 recovered in 3 mL of LB without antibiotics for 2 hours. Then, 25 uL of culture was mixed with
752 25 uL of water, heated to 95C for 5 minutes to lyse cells, cooled, and frozen at -20C for later
753 analysis.

754 **Signal Promoter Strength Measurement**

755 bSLS.114 was transformed with Marionette plasmid³⁴ (pAJM.683, pAJM.011, or pAJM.771) and
756 grown overnight in LB with antibiotic (kanamycin). In the morning, 60ul of overnight culture was
757 added to 2 tubes of 3 ml LB, one with antibiotic and inducer and the other with antibiotic and no
758 inducer. The cultures were grown for 2 hours, 1ml of cell suspension pelleted (8000g for 1 min)
759 and resuspended twice in 1mL PBS, and OD600 (600nm absorbance) and YFP fluorescence

760 (513nm excitation/538nm emission) measurements were taken on a Molecular Devices
761 SpectraMax i3 plate reader using a black, clear bottom 96-well plate and 200ul of resuspended
762 cells. OD600 was measured using a kinetic scan for 2 minutes, taking measurements every 25
763 seconds, with a 1 second shake in between. Fluorescence was measured as a kinetic scan for
764 2 minutes, taking measurements every 20 seconds, with a 1 second shake in between.

765 **Recording System Fitness Measurements**

766 Cells were transformed with plasmids (sequentially in the case of multiple plasmids). Single
767 colonies were picked from plates and grown overnight in 3 mL of LB with antibiotics and without
768 inducers. In the morning, 60 uL of culture was added to 3 mL of LB with antibiotics and left to sit
769 at room temp for ~4 hours. Then, cultures were transferred to a shaking incubator and grown for
770 2 hours. After 2 hours, the OD600 of each culture was measured using a NanoDrop-2000c
771 spectrophotometer and cultures diluted to an OD600 of 0.05 in LB with antibiotics. Next,
772 inducers were added at the appropriate strength and 200 uL of each culture (OD600 = 0.05 with
773 antibiotics and inducers) was transferred to a clear-bottomed 96-well plate. The 96-well plate
774 was loaded onto a Molecular Devices SpectraMax i3 plate reader set at 37C, and OD600
775 measurements taken every 2.5 minutes for the next 15 hours, with a 30 second shake before
776 each reading.

777 **Temporal Recordings**

778 Cells were transformed sequentially, first with pSBK.134 or pSBK.136 and then with pSBK.079.
779 For recording, single colonies were picked from plates and grown overnight in 3 mL of LB with
780 antibiotics and without inducers. In the morning, 150 uL of culture was diluted into 3 mL of LB
781 with antibiotics and appropriate inducers (Fig. 4) and grown for 8 hours. After 8 hours, 60 uL of
782 culture was diluted into 3 mL of LB with appropriate inducers and grown overnight (16 hours). In
783 the morning, 150 uL of culture was diluted into 3 mL of LB with appropriate inducers (for second
784 day of expression) and grown for 8 hours. Samples were collected at this 24-hour timepoint. 25
785 uL of culture was mixed with 25 uL of water, heated to 95C for 5 minutes to lyse cells, cooled,

786 and frozen at -20C for later analysis. After 8 hours, 60 uL of culture was diluted into 3 mL of LB
787 with appropriate inducers and grown overnight (16 hours). In the morning, 25 uL of culture was
788 mixed with 25 uL of water, heated to 95C for 5 minutes to lyse cells, cooled, and frozen at -20C
789 for later analysis.

790 **Computational Model of Retron Recording**

791 A computational model of retron recording was written using Python 3 and the following
792 modules, packages, and libraries: numpy, matplotlib, random, itertools, and xlswriter. For
793 modeling, we assume that spacer acquisition is well approximated as a Poisson process in
794 which acquisitions occur at some average rate over time, but where the precise timing of these
795 events is random and independent of the timing of previous events. We believe this is a fair
796 approximation of spacer acquisition due to the overall low rate of spacer acquisition in the retron
797 recording system (single digit percentages of arrays expanded over 24 hours), the
798 demonstrated ability of CRISPR arrays to be multiply expanded, and the current understanding
799 of CRISPR adaptation indicating that acquisitions occur one at a time. To simulate spacer
800 acquisition in a population of cells, we first define a “Cell” class of which each instance
801 possesses an attribute called an “array”. The user defines the following parameters of the
802 recording experiment: number of cells, rates of acquisition (in units of integrations per hour per
803 cell) of retron-derived signal “A” with and without inducer present, rates of acquisition of retron-
804 derived signal “B” with and without inducer present, rate of acquisition of non-retron-derived
805 signal “N”, time of induction of signal “A”, time of induction of signal “B”, and order of induction
806 (e.g. “A” before “B”). For the first epoch, each “Cell” instance samples three different Poisson
807 distributions (one each for signals “A”, “B”, and “N”) to determine the number of spacers of each
808 type which are added to its “array” during the epoch. The order of these spacers is then
809 randomized and appended to the “array”. For example: when the order of induction is “A” before
810 “B”, the cell is subject to the following rates of acquisition: “A” with inducer, “B” without inducer,
811 and “N”. For each signal (“A”, “B”, and “N”) the cell samples a Poisson distribution defined by

812 the probability mass function $p(k; \lambda) = \frac{\lambda^k e^{-\lambda}}{k!}$ where k is the number of acquisitions of that
813 signal ($k = 0, 1, 2 \dots$) and λ is the expected number of acquisitions of that signal (defined as
814 the rate of acquisition of a given signal times the length of the epoch). It is fair to randomize the
815 order of acquisitions occurring in each epoch, prior to appending them to the array, because the
816 timing of the events is random by definition. For example: given that a cell acquires one “A”
817 spacer and one “N” spacer in an interval with constant rates of acquisition of “A” and “N”, it is
818 equally likely that “A” comes before “N” as it is that “N” comes before “A”. After acquisitions
819 during the first epoch are completed, the process is repeated for the second epoch (using the
820 relevant rates of acquisition for all three signals). At this point, the arrays are complete and
821 ready for analysis using the ordering rules. Recordings, replicates, and ordering rule analysis
822 were simulated using purpose-built scripts to investigate parameters of interest. Relevant data
823 was exported to Excel sheets for further analysis and visualized using GraphPad Prism.

824 **Long-Term Passage for Data Stability**

825 24+24-hour, “A”-then-“B” recordings were made in bSLS.114 cells harboring plasmids
826 pSBK.134 and pSBK.079 as described previously. At hour 48, 25 uL of culture was mixed with
827 25 uL of water, heated to 95C for 5 minutes to lyse cells, cooled, and frozen at -20C for later
828 analysis. 500 uL of culture was combined with 500 uL of 50% glycerol and frozen at -80C for
829 future outgrowth. To begin long-term culture, recording glycerol stocks stored at -80C were
830 thawed at room temperature, 100 ul of thawed cells added to 25 mL of LB with antibiotics, and
831 the culture left to shake at 37C for 24 hours. Every day for the next 14 days, 25 ul of culture was
832 sampled, boiled, and frozen as above, and 50 uL of culture added to 25 mL of fresh LB with
833 antibiotics (ratio of 1:500, yielding roughly 9 generations per day). Samples from days 0, 2, 5,
834 and 9 were sequenced and analyzed.

835 **Analysis of Spacer Acquisition**

836 Analysis of spacer acquisition was conducted by sequencing a library of all CRISPR arrays in
837 an experimental population using an Illumina MiSeq instrument. Libraries were created by
838 amplifying a region of the genomic CRISPR array using PCR, then indexed using custom
839 indexing oligos. Up to 192 conditions were run per flow cell. A list of oligo prespacers and
840 primers can be found in Supplementary Table 4.

841 **Processing and Analysis of MiSeq Data**

842 Sequences were analyzed using custom Python software, which will be available on GitHub
843 upon peer-reviewed publication. In brief, newly acquired spacer sequences were extracted from
844 array sequences based on their position between identifiable repeats and compared to
845 preexisting spacers in the array. In this preliminary analysis, metrics were collected including
846 number of expansions in arrays (unexpanded, single, double, and triple expanded) and
847 proportion of each present in the library. Sequenced arrays were sorted into subcategories
848 based on these characteristics (e.g. doubly expanded with first three repeats identifiable) for
849 further analysis. Next, to determine number of retron-derived spacers and the order of spacers
850 in multiply expanded arrays, two different analyses were used: one strict and one lenient. In the
851 strict analysis (used in figures 1, 2, and 3) a retron-derived spacer is defined to be a spacer
852 which contains the 23-base core region of the hypothetical prespacer structure from a given
853 retron (with three mismatches or indels allowed). In the lenient analysis (used in figures 4 and 5)
854 a retron-derived spacer is defined to be a spacer which contains an 11-base region of the
855 hypothetical prespacer consisting of the 7-base barcode region and 2 bases on either side (with
856 one mismatch or indel allowed). The order of spacers in multiply expanded arrays is then
857 reported (e.g. Leader-NNA) and these data are used to complete the ordering rule analysis.

858 **Data Availability**

859 All data supporting the findings of this study are available within the article and its
860 supplementary information, or will be made available from the authors upon request.
861 Sequencing data associated with this study are available in the NCBI SRA (PRJNA838025).

862 **Code Availability**

863 Custom code to process and analyze data from this study is available on GitHub
864 (<https://github.com/Shipman-Lab/Spacer-Seq>).

865 **Methods References**

- 866
867 40 Mosberg, J. A., Gregg, C. J., Lajoie, M. J., Wang, H. H. & Church, G. M. Improving Lambda Red
868 Genome Engineering in Escherichia coli via Rational Removal of Endogenous Nucleases. *PLOS*
869 *ONE* **7**, e44638, doi:10.1371/journal.pone.0044638 (2012).
870 41 Moore, S. D. in *Strain Engineering: Methods and Protocols Methods in Molecular Biology* (ed
871 James A. Williams) 155-169 (Humana Press, 2011).
872 42 Datsenko, K. A. & Wanner, B. L. One-step inactivation of chromosomal genes in Escherichia coli
873 K-12 using PCR products. *PNAS* **97**, 6640-6645, doi:10.1073/pnas.120163297 (2000).
874 43 Rogers, J. K. *et al.* Synthetic biosensors for precise gene control and real-time monitoring of
875 metabolites. *Nucleic Acids Res* **43**, 7648-7660, doi:10.1093/nar/gkv616 (2015).

ACKNOWLEDGEMENTS

876 Work was supported by funding from the Simons Foundation Autism Research Initiative
877 (SFARI) Bridge to Independence Award Program, the Pew Biomedical Scholars Program, the
878 NIH/NIGMS (1DP2GM140917-01), and the UCSF Program for Breakthrough Biomedical
879 Research. S.L.S. is a Chan Zuckerberg Biohub investigator and acknowledges additional
880 funding support from the L.K. Whittier Foundation. S.K.L. was supported by an NSF Graduate
881 Research Fellowship (2034836). S.C.L. was supported by a Berkeley Fellowship for Graduate
882 Study. We thank Kathryn Claiborn for editorial assistance.

883 **AUTHOR CONTRIBUTIONS**

884 S.L.S. conceived the study with J.N. and G.M.C. contributing. S.B.K. and S.L.S. designed
885 experiments and analyzed all data. Contributions to data collection were made from S.K.L.
886 (Extended Data Figure 4), C.B.F. (Extended Data Figure 5), and S.L.S. (Figure 1c, e; Figure 3b-
887 e). S.B.K. collected all other data. E.L., M.G.S., and S.C.L. performed preliminary experiments
888 not included in the figures. S.B.K. performed the computational modeling of recordings. S.B.K.
889 wrote the manuscript with input from all co-authors.

890 **COMPETING INTERESTS**

891 S.L.S., G.M.C., M.G.S., and J.N. are named inventors on a patent application assigned to
892 Harvard College, Method of Recording Multiplexed Biological Information into a CRISPR Array
893 Using a Retron (US20200115706A1).

894

895 **CORRESPONDING AUTHOR**

896 Correspondence to S.L.S. (seth.shipman@gladstone.ucsf.edu).

897 **Extended Data Figure Legends**

898 **Extended Data Figure 1. Accompaniment to Figure 1.** **a.** Hypothetical Eco1 wild-type ncRNA-linked
899 RT-DNA structure. **b.** Hypothetical Eco1 v32 ncRNA-linked RT-DNA structure and hypothetical duplexed
900 RT-DNA prespacer structure. Nucleotides that are altered from wild-type Eco1 are shown in orange. **c.**
901 Hypothetical Eco1 v35 ncRNA-linked RT-DNA structure and hypothetical duplexed RT-DNA prespacer
902 structure. Nucleotides that are altered from wild-type Eco1 are shown in green.

903

904 **Extended Data Figure 2. Accompaniment to Figure 2.** **a.** Hypothetical barcoded Eco1 v35 ncRNA-
905 linked RT-DNA structure and hypothetical duplexed RT-DNA prespacer structure. Bases used to barcode
906 retons are shown in red.

907

908 **Extended Data Figure 3. Accompaniment to Figure 3.** **a.** Hypothetical wild-type Eco4 ncRNA-linked
909 RT-DNA structure. ExoVII-dependent RT-DNA cleavage site is shown as a red slash. **b.** Eco4-derived
910 spacer sequences and orientations. Bases are colored to match Figure 3f. **c.** Proportion of Eco4-derived
911 spacers in each orientation. Open circles are individual biological replicates.

912

913 **Extended Data Figure 4. Change in YFP fluorescence when expressed using inducible promoters.**
914 The Y-axis shows fluorescence (in arbitrary units) normalized to culture density (OD600).

915

916 **Extended Data Figure 5. Growth curves (upper plot) and max growth rates (lower plot) of E. coli**
917 **with different combinations of retron recording components and inducers.** In growth curve plots the
918 solid line is the mean OD600 of 3 biological replicates, with dotted lines showing the standard deviation.
919 In maximum growth rate plots, each symbol is a single biological replicate. Bars show the mean and
920 standard deviation. Statistically significant differences in maximum growth rate, as calculated by Tukey's
921 multiple comparison's test, are highlighted. **a.** Growth kinetics of E. coli with different combinations of
922 retron recording plasmids, all without inducers. **b.** Growth kinetics of E. coli with recording plasmid
923 pSBK.079, with and without inducers. **c.** Growth kinetics of E. coli with signal plasmid pSBK.134, with and
924 without inducers. Only one biological replicate is present in condition "pSBK.134 + aTc" (pink). **d.** Growth
925 kinetics of E. coli with signal plasmid pSBK.136, with and without inducers. **e.** Growth kinetics of E. coli
926 with signal plasmid pSBK.134 and recording plasmid pSBK.079, with and without inducers. **f.** Growth
927 kinetics of E. coli with signal plasmid pSBK.136 and recording plasmid pSBK.079, with and without
928 inducers.

929

930 **Extended Data Figure 6. Accompaniment to Figure 4.** **a.** Ordering rules for pSBK.134 "A"-before-"B"
931 replicates. The scores for each rule, and the composite score, are shown for each individual replicate. X-
932 containing boxes indicate that no informative arrays, for that particular rule, were present in that replicate.
933 **b.** As in panel (a), ordering rules for pSBK.134 "B"-before-"A" replicates. **c.** As in panel (a), ordering rules
934 for pSBK.136 "A"-before-"B" replicates. **d.** As in panel (a), ordering rules for pSBK.136 "B"-before-"A"
935 replicates.

936

937 **Extended Data Figure 7. Long-term stability of retron-derived recordings in CRISPR arrays.** **a.**
938 Ordering rules for 24+24-hour, "A"-before-"B" recordings during post-recording multiday culture. Individual
939 and composite scores are shown for samples taken on days 0, 2, 5, and 9 of culture. Each open circle
940 represents the score, for that rule, from a single biological replicate. A total of 3 biological replicates are

941 shown here. **b.** Changes in ordering rule scores over time in biological replicate 1. **c.** Changes in ordering
942 rule scores over time in biological replicate 2. **d.** Changes in ordering rule scores over time in biological
943 replicate 3.