

1



2

3 **Main Manuscript for**

4 **Structured sequences emerge from random pool**
5 **when replicated by templated ligation**

6 Patrick W. Kudella¹, Alexei V. Tkachenko², Annalena Salditt¹, Sergei Maslov^{3,4}, Dieter Braun^{*1}

7 ¹Systems Biophysics and Center for NanoScience, Ludwigs-Maximilian-Universität München, 80799
8 Munich, Germany

9 ²Center for Functional Nanomaterials, Brookhaven National Laboratory, Upton, New York 11973, USA

10 ³Department of Bioengineering, University of Illinois at Urbana-Champaign, Urbana, Illinois 61801, USA

11 ⁴Carl R. Woese Institute for Genomic Biology, University of Illinois at Urbana-Champaign, 1206 West
12 Gregory Drive, Urbana Illinois 61801, USA* Paste corresponding author name(s) here.

13 **Email:** dieter.braun@lmu.de

14 Braun: 0000-0001-7751-1448

15 Maslov: 0000-0002-3701-492X

16 Tkachenko: 0000-0003-1291-243X

17 **Classification**

18 Physical Sciences, Biophysics and Computational Biology; Biological Sciences, Evolution

19 **Keywords**

20 Origin of Life, DNA replication, Darwinian evolution, templated ligation, sequence entropy

21 **Author Contributions**

22 P.W.K. performed the experiments, prepared the library for sequencing, performed the
23 demultiplexing, the analysis, programmed the analysis software, analyzed the data, drafted and
24 wrote the manuscript. A.V.T and S.M. performed the theoretical analysis and analyzed the data in
25 context of their already published theoretical work, drafted graphs, drafted and wrote the
26 manuscript. D.B. contrived the experiment, guided the experimental progress, analyzed data and
27 drafted the manuscript.

28 **This PDF file includes:** Main Text and Figures 1 to 5

29

30 **Abstract**

31 The central question in the origin of life is to understand how structure can emerge from
32 randomness. The Eigen theory of replication states for sequences that are copied one base at a
33 time, the replication fidelity has to surpass an error threshold to avoid that replicated specific
34 sequences become random due to the incorporated replication errors [M. Eigen,
35 *Naturwissenschaften* 58(10), 465-523 (1971)]. Here we showed that linking short oligomers from a
36 random sequence pool in a templated ligation reaction reduced the sequences space of product
37 strands. We started from 12mer oligonucleotides with two bases in all possible combinations and
38 triggered enzymatic ligation under temperature cycles. Surprisingly, we found the robust creation
39 of long, highly structured sequences with low entropy. At the ligation site, complementary and
40 alternating sequence patterns developed. However, between the ligation sites, we found either an
41 A-rich or a T-rich sequence within a single oligonucleotide. Our modeling suggests that avoidance
42 of hairpins was the likely cause for these two complementary sequence pools. What emerged was
43 a network of complementary sequences that acted both as templates and substrates of the
44 reaction. This autocatalytic ligation reaction could be restarted by only a few majority sequences.
45 The findings showed that replication by random templated ligation from a random sequence input
46 will lead to a highly structured, long and non-random sequence pool. This is a favorable starting
47 point for a subsequent Darwinian evolution searching for higher catalytic functions in an RNA world
48 scenario.

49 **Significance Statement**

50 The structure of life emerged from randomness. Typically, this is attributed to the selection of
51 function by molecular Darwinian evolution. But can we already find sequence selection before the
52 onset of Darwinian evolution? We experimentally studied a simple model of replication by templated
53 ligation. We did not copy sequences base-by-base, but found long strands emerging by ligation of
54 two substrate strands bound to a third template strand from the same random sequence pool under
55 temperature cycling. We started from a minimal setting of random 12mer sequences that used only
56 two bases. Interestingly, the ligated strands showed highly structured sequences that form a
57 replication network. We discuss theoretical models for how these non-random sequences could
58 emerge. The findings show a likely mode to reduce sequence entropy before the onset of Darwinian
59 evolution.

60 **Main Text**

61 **Introduction**

62 One of the dominant hypotheses to explain the origin of life(1–3) is the concept of RNA world. It is
63 built on the fact that catalytically active RNA molecules can enzymatically promote their own
64 replication(4–6) via active sites in their three dimensional structures(7–9). These so-called
65 ribozymes have a minimal length of 30 to 41 bp(9, 10) and, thus, a sequence space of more than
66 $4^{30} \approx 10^{18}$. The subset of functional, catalytically active sequences in this vast sequence space is
67 vanishingly small(11) making spontaneous assembly of ribozymes from monomers or oligomers all
68 but impossible. Therefore, prebiotic evolution has likely provided some form of selection guiding
69 single nucleotides to form functional sequences and thereby lowering the sequence entropy of this
70 system.

71 The problem of non-enzymatic formation of single base nucleotides and short oligomers in settings
72 reminiscent of the primordial soup has been studied before(12–17). However, the continuation of
73 this evolutionary path towards early replication networks would require a pre-selection mechanism
74 of oligonucleotides (as shown in Figure 1a), lowering the information entropy of the resulting
75 sequence pool(18–22). In principle, such selection modes include optimization for information
76 storage, local oligomer enrichment e.g. in hydrogels or in catalytically functional sites.

77 An important aspect of a selection mechanism is its non-equilibrium driving force. Today's highly
78 evolved cells function through multistep and multicomponent metabolic pathways like glycolysis in

79 the Warburg effect(23) or by specialized enzymes like ATP synthase which provide energy-rich
80 adenosine triphosphate (ATP)(24). In contrast, it is widely assumed(3, 4, 25–28) that selection
81 mechanisms for molecular evolution at the dawn of life must have been much simpler, e.g.
82 mediated by random binding between biomolecules subject to non-equilibrium driving forces such
83 as fluid flow and cyclic changes in temperature.

84 Here, we explored the possibility of a significant reduction of sequence entropy driven by templated
85 ligation(19) and mediated by Watson-Crick base pairing(29). Starting from a random pool of
86 oligonucleotides we observed a gradual formation of longer chains showing reproducible sequence
87 landscape inhibiting self-folding and promoting templated ligation. Here we argue, that base pairing
88 combined with ligation chemistry, can trigger processes that have many features of the Darwinian
89 evolution.

90 As a model oligomer we decided to use DNA instead of RNA since the focus of our study is on
91 base pairing which is very similar for both(30). We start our experiments with a random pool of
92 12mers formed of bases A (adenine) and T (thymine). This binary code facilitates binding between
93 molecules and allows us to sample the whole sequence space in microliter volumes
94 ($2^{12} \ll 10 \mu\text{M} * 20 \mu\text{l} * N_A = 10^{14}$).

95 Formation of progressively longer oligomers from shorter ones requires ligation reactions, a method
96 commonly employed in hairpin-mediated RNA and DNA replication(31, 32). At the origin of life, this
97 might have been achieved by activated oligomers(33, 34) or activation agents(35–37), **while later**
98 **on the formation of simple ribozyme ligases seemed possible(38)**. Our study is focused on inherent
99 properties of self-assembly by base pairing in random pools of oligomers and not on chemical
100 mechanisms of ligation. Hence, we decided to use TAQ DNA ligase - an evolved enzyme for
101 templated ligation of DNA(21), **that is known for its ligation site sequence specificity(39, 40) and**
102 **lack of sequence-dependent ligation rate (compare SI-section 21)**. This allowed for fast turnovers
103 of ligation and enabled the observation of sequence dynamics.

104 **Results**

105 To test templated elongation of polymers in pools of random sequence oligomers, we prepared a
106 10 μM solution of 12mer DNA strands composed with nucleobases A and T (sequence
107 space: 4096) and subjected it to temperature cycling, similar to reference(21) with 20 s at
108 denaturation temperature of 75 °C and 120 s at ligation temperature of 33 °C. Temperatures were
109 selected according to the melting dynamics of the DNA pool; the time steps were prolonged relative
110 to Toyabe and Braun(21) (SI section 5.3) because of a greater sequence space. The larger
111 sequence space of full random 12mers with all four bases did not show any ligation under the same
112 experimental conditions (SI section 5.2). **The sample was split into multiple tubes and exposed to**
113 **200, 400, 600, 800, 1000 temperature cycles and one tube kept at 4 °C for reference, all without**
114 **influx or outflux of strands.**

115 To study the length distributions in our samples we used polyacrylamide gel electrophoresis
116 (PAGE, Figure 1d). The first lane is the reference sequence not exposed to temperature cycling,
117 where small amounts of impurities are visible at short lengths (SI Section 3.1). The latter lanes
118 show the temperature-cycled samples. As the number of cycles increases, progressively longer
119 strands in multiples of 12 emerge, as the original pool only consisted of 12mers. Figure 1c shows
120 the concentration quantification of each lane (compare SI section 3). For higher cycle counts the
121 total amount of products increases and the concentration as a function of length decreases slower.
122 The behavior of this system is dependent on the time and temperature for both steps in the
123 temperature cycle, the monomer-pool concentration and the sequence space of the pool
124 (SI section 5).

125 An important property of the initial monomer-pool is its sequence content. Although for pools with
126 lower sequence complexity it is possible to show different strand compositions using PAGE(41,
127 42), a large size of our “monomer” ($2^{12}=4096$) and 24mer product pools (sequence space:

128 $2^{24} \approx 16.8 \times 10^6$) excludes this approach. Thus, we analyzed our final products by Next Generation
129 Sequencing (NGS) to get insights into product strand compositions.

130 Plotting the probability of finding a base at a certain position (Figure 1c inset) revealed no distinct
131 pattern in 12mers other than a slight bias towards As. However, longer chains starting with 24mers
132 developed a strikingly inhomogeneous sequence pattern: bases around ligation sites show a
133 distinct AT-alternating pattern, while regions in the middle of individual 12mers are preferentially
134 enriched with As.

135
136 The information entropy of longer chains is expected to be smaller than the entropy of a random
137 sequence strand of the same length, if some sort of selection mechanism is involved(19). We
138 analyzed the entropy reduction for different lengths of products (Figure 2a) as well as the positional
139 dependence of the single base entropy for 60mer products (Figure 2b). The relative entropy
140 reduction is similar to one used in Derr *et. al*(43) where 1 describes a completely random ensemble
141 and 0 an ensemble of only one sequence. Entropy reduction was observed in all analyzed product
142 lengths with a greater reduction observed for longer oligomer lengths. The entropy of each 12mer
143 subsequence was also found to be significantly lower than that of random 12mers (Figure 2b, black
144 line). The central subsequence had the lowest entropy while 12mers located at both ends of chains
145 had relatively higher entropies. This behavior was also observed as a function of nucleotide position
146 within a 12mer suggesting a multi-scale pattern of entropy reduction.

147 In the initial pool of random 12mers the A-to-T ratio distribution is shaped binomially, as expected
148 for a random distribution. However, it dramatically shifted for 24mer products of ligation: a bimodal
149 distribution of about 65:35 % A:T (A-type) as well as the inverse, 35:65 % A:T (T-type) was
150 observed with 24mer products (Figure 2c). DNA strands composed of only two complementary
151 bases are more prone to formation of single-strand secondary structures like hairpins than DNAs
152 composed of all four bases. In our templated ligation reaction, we expected that hairpin-sequences
153 are not elongated and also not used as template-strands because they form catalytically passive
154 Watson-Crick-base-paired configuration. A bimodal AT-ratio distribution (Figure 2d) also emerged
155 in a kinetic computational model in which a pool of random 12mers was seeded with a small initial
156 amount of random sequence 24mers. 24mers that formed hairpins could not act as templates and
157 were therefore less likely to be reproduced (see SI for details of this model, section 18.2).

158 For longer products the bimodal distribution got sharper and centered at approximately
159 70:30 % A:T and 30:70 % A:T (Figure 2e). To compare the distributions of different lengths we
160 computed probability density functions (PDF) of A:T fractions. Each distribution is the sum (integral)
161 over all probabilities P_N to find a certain A:T-fraction $d_{A:T}$ in chains of length N :

$$162 \quad \int P_N(A:T) d_{A:T} = 1. \quad (1)$$

163 The main difference of longer oligomers was a rapid increase of the ratio between the number of
164 A-type and T-type sequences. As oligomers get longer the effect becomes more pronounced. This
165 might be a result of a small bias in the initial pool which has slightly more monomers of A-type than
166 T-type (SI section 9.1).

167 As predicted theoretically(18), the eventual length distribution is approximately exponential. A
168 small A-T bias leads to the respective average chain lengths, \bar{N}_A and \bar{N}_T , to be somewhat different
169 for the two subpopulations. As a result, the bias gets strongly amplified with increasing chain length:

$$170 \quad P_N(A:T) \sim \exp\left(-N\left(\frac{1}{\bar{N}_A} - \frac{1}{\bar{N}_T}\right)\right) = \beta^{-N/12}. \quad (2)$$

171 A simple phenomenological model can successfully capture the major features of the observed A:T
172 PDFs for multiple chain lengths. Specifically, we assume both A-type and T-type sub-populations -
173 to maximize the sequence entropy, subject to the constraint that the average A:T content is shifted

174 from the midpoint (50:50 % composition), by values $\pm x_0$, respectively. This model presented in SI
175 section 18.1, results in a distribution that strongly resembles experimental data, as shown in
176 Figure 2e-f A:T profiles for all chain length are fully parameterized by only two fitting parameters:
177 $\beta = 0.785$, and $x_0 = 0.2$.

178 The proposed mechanism of selection of A-type and T-type subpopulations due to hairpin
179 suppression is further supported by direct sequence analysis. Figure 3b shows PDFs of the longest
180 sequence motifs that would allow hairpin-formation, across the entire pool of sequences of given
181 lengths. While the overall chain length increased by a factor of seven (12 to 84 nt), the most likely
182 hairpin length only grew by a factor of 1.89 (3.7 to 7 nt) (Figure 3b). The observed relationship
183 between the strand length N and the most likely hairpin stem length l_0 can be successfully
184 described by a simple relationship obtained within the above described maximum-entropy model.
185 Specifically, for a random sequence with bias parameter $p = 0.5 + x_0$, one expects N to be related
186 to l_0 as follows (as in Figure 2f):

$$187 \quad N = 2l_0 + \sqrt{2}(2p(1 - p))^{-l_0/2}. \quad (3)$$

188 As one can see in Figure 3c, this result is in an excellent agreement with experimental data for all
189 the long chains, assuming $p=0.785$. This A:T ratio is indeed comparable to the one observed in the
190 A-type subpopulation. On the other hand, the maximum probability length of the longest hairpin for
191 12mers is consistent with an unbiased composition, $p=0.5$.

192 While hairpin formation inhibits the self-reproduction based on template-based ligation, Figure 3b
193 reveals another dramatic feature: a small fraction of chains does feature very long hairpin-forming
194 motifs (seen as shoulders in the distribution function). This effect also reveals itself as small peaks
195 on the 84mer curve in Figure 2e. Those peaks around A:T ratio 0.4, 0.5 and 0.6. stem from
196 subpopulations that have multiple A-types as well as multiple T-type subsequences (see SI
197 section 12) and are prone to hairpin formation.

198 The mechanism of formation of these self-binding sequences may involve recombination of shorter
199 A-type and T-type chains, or self-elongation of shorter hairpins. In either case, the hairpin sequence
200 cannot efficiently reproduce by means of template ligation. However, the remainder of the pool would
201 keep producing them as byproduct. Ironically, for the templated ligation reaction this is a possible
202 failure mode, but those long hairpins may play a key role in the context of origin of life, as precursors
203 of functional motifs. For instance, work by Bartel and Szostak(11, 44) identifies RNA self-binding
204 as crucial for the direct search of ribozymes – those molecules need to fold into non-trivial
205 secondary structures to gain their catalytic function.

206 The separation into A-type and T-type subpopulations only accounts for a small part of the
207 sequence entropy reduction. The emerging ligation landscape in the sequence space is far richer.

208 Sequence analysis of the junctions in-between original 12mer revealed additional information about
209 that landscape, already hinted by patterns seen in Figure 1b. We characterize pairs of junction-
210 forming sequences with their Z-scores, i.e. probability of their occurrence scaled with its expected
211 value and divided by the standard deviation calculated in the random binding model (see SI
212 section 14).

213 Figure 4a shows Z-score heatmaps for junctions within A-type (left panel) and T-type (right panel)
214 subpopulations. More specifically, we show sequences left (row) and right (column) of the junction
215 between the 4th to the 5th 12mers in the respective 72mer. These heatmaps reveal a complex
216 landscape of over- and under-represented junction motifs shown respectively in dark-teal and dark-
217 ocher colors. Emergence of such complex landscape has been theoretically predicted in Ref. (19)
218 landscape peaks include repeating A-T motif of alternating bases crossing the ligation site (dark-
219 teal peak near the center of each of both heatmaps). Relatively rare motifs (valleys) correspond to
220 poly-A and poly-T sequences extending across the junction (dark-ocher areas). One exception to

221 this rule is a relatively abundant poly-A motif at the bottom right of the A-type heatmap (light-teal).
222 Interestingly, these junction sequences had AT-patterns in the beginning of the “left side” and the
223 end of the “right side”. This might provide a clue to the origin of these “abnormal” junction motifs.
224 Indeed, they may have been templated by abundant poly-T sequences in the middle of T-type
225 12mers flanked by alternating A-T motifs. In other words, junctions at templates of poly-A junction
226 motifs may have been shifted by 6 nt relative to substrates. Actually, substrates have no restriction
227 on where they may hybridize on a long template and might happen to have their ligation site in the
228 region of poly-T of the template strand. We call this “ligation site shift”, as explained in SI section 16.
229 Other preferred junction subsequences include repetitions of the AAT motif across the junction (the
230 dark-teal peak in the upper left corner of the left panel). **The origin of the dominant A-T sequence
231 pattern is analyzed with a 12mer pool, sub-motif based Monte-Carlo style templated ligation
232 reaction in SI-section 21. Based in this simulation, small deviations from randomness in the original
233 12mer pool may lead to abundant sequence patterns, especially in the case of a self-similar motif
234 like “AT”, irrespective of a possible sequence bias of the ligation yield of the used ligase.”**

235 How similar are selective pressures operating on sequences of different 12mers within longer
236 chains? Figure 4b quantifies this similarity in terms of sample Pearson-Correlation-Coefficient
237 (sPCC) between abundances of 12mer sequences in different positions of long chains of different
238 lengths.

239 We compare the abundances of $2^{12}=4096$ possible 12mer sequences in positions 1 to 6 within all
240 72mers and compare them to each other and abundances of 12mers in positions 1 to 7 in all
241 84mers. Similar results were obtained for other chains longer than 36 nt. A rectangle of very high
242 correlations (>0.9) at the center of the table in Figure 4b means that very similar sequences get
243 selected at all internal positions of all chains (note that only chains longer than 36nt have such
244 internally positioned 12mers). However, the light border of the table means that a rather different
245 subset of 12mers gets selected in the first and the last position of a multimer. Whatever the nature
246 of selection pressure acting on these 12mers, it is consistent across oligomers of different lengths
247 as manifested by the high correlation in the lower left and the upper right corner of the table in
248 Figure 4b.

249 A simple hypothesis comes to mind: a strand is prolonged and grows in this random sequence
250 templated ligation system as long as the sequences attached to it share similar sequence motifs
251 resulting in high values of sPCC for all internal 12mers. But when a 12mer sequence that is similar
252 to the start- or end-subsequence is attached, the growth in that direction stops.

253 Comparison of abundances of internal 12mers in A-type and T-type subpopulations predictably
254 yielded no positive correlation and in fact resulted in a slight negative correlation (see SI
255 section 11). However, abundances of reverse complements of sequences from the T-type
256 subpopulation are strongly correlated with those of the A-type resulting in a sPCC matrix similar to
257 that shown in Figure 4b (see SI-Figure 13). Therefore, chains in two groups (A-type and T-type)
258 show a considerable degree of reverse complementarity to each other. This fits the elongation and
259 replication mechanism by templated ligation.

260 To further explore selection capabilities of templated ligation as a function of 12mer sequences in
261 the initial pool we conducted three additional experiments referred to as “Replicator”, “Random”
262 and “Network”. The “Random” experiment started with eight randomly chosen 12 nt sequences
263 served as a control. In the “Replicator” experiment the pool consisted of eight 12 nt sequences
264 artificially designed for efficient elongation (see below). In the “Network” experiment we populated
265 the pool with eight naturally selected 12 nt sequences commonly found as subsequences of long
266 strands in our original ligation experiment with 4096 12mers. To identify these 12mers, we built a
267 network of the most common 12mers found in A-type oligomers with length of more than 48 nt.
268 This network does not include the first and the last 12mers, in a multimer as those are known to be
269 statistically different from the internal ones (see Figure 4b). The circles in Figure 5a represent
270 unique 12 nt subsequences while their size describes their Z-scores quantifying their abundance
271 in long chains. The width of the connecting line describes the probability that two subsequences

272 are found one after another in a multimer. The same is done for T-type sequences (Figure 5b). This
273 representation of a polymer is known as de Bruijn graph(45) and has been commonly used in DNA
274 fragment analysis and genome assembly(46) and more recently in the context of templated
275 ligation(19).

276 De Bruijn networks in Figure 5a break up into several clusters connecting 12mers with similar
277 subsequences at junctions (TAA-TAA in the top cluster marked by a dark-magenta node, ATA-ATA
278 in the middle one, and AAT-AAT in the bottom one). Note that these three common junction
279 subsequences are all related via template shifts. The most common subgraphs found in the A-type
280 network and mirrored among their reverse complements in the T-type network. This pattern is
281 consistent with selection driven by templated ligation (see SI section 19). Among the eight most
282 common subsequences in the A-type network (light and dark magenta nodes in Figure 5a), four
283 (dark magenta nodes) had a reverse complement among the eight most common subsequences
284 of the T-type network (light and dark magenta nodes in Figure 5b). These sequences were chosen
285 as the pool of eight 12mers in the “Network” sample. The “Random” sample consisted of eight
286 12mers which were randomly chosen from the 4096 possible AT-only 12mers. The “Replicator”
287 sample consisted of eight strands that were built to form three-strand complexes that resemble the
288 assumed first ligation reaction in the pool (SI section 17.1).

289 The length distribution of oligomers (Figure 5d) with concentrations quantified from the PAGE gel
290 image (Figure 5c) shows that the “Network” sample produced the most product, as the remaining
291 12mer sequence concentration was reduced below two other samples down to almost 5 μ M. The
292 length distribution in both “Random” and “Network” samples is well described by a piecewise-linear
293 distribution predicted in Ref. (18). For short product lengths ranging between 48mers up to 136mers
294 the “Random” sample produced more oligomers than the “Replicator” sample. However, for even
295 longer strands, the “Replicator” sample generated the largest number of really long strands since
296 its length distribution reached a plateau around 120mers. This is probably due to the nature of the
297 eight-sequences pools used here with the “Replicator” one made to form well aligned dsDNA that
298 can be properly ligated. According to NUPACK(47), 12mers in the “Random” sample should not
299 form any complexes that could be subsequently ligated by the TAQ ligase. However, our results
300 shown in Figure 5c prove the existence of extensive ligation even in the “Random” sample.
301 Presumably, it was initially triggered by small concentration of complexes formed with low
302 probability, which were subsequently amplified due to the exponential growth of longer strands in
303 our experiment, just like in the “Network” sample.

304 Discussion

305 We experimentally studied templated ligation in a pool of 12mers made of A and T bases with all
306 possible sequences ($2^{12}=4096$), subjected to multiple temperature cycles. To accelerate
307 hypothetical spontaneous ligation reactions operating in the prebiotic world, we employed TAQ
308 DNA ligase in our experiments. This process produced a complex and heterogeneous ensemble
309 of oligomer products. By performing the “next generation sequencing” of these oligomers, we found
310 that long strands in this ensemble have a significantly lower information entropy compared to a
311 random set of oligomers of the same length. This effect became increasingly more pronounced for
312 longer oligomers (Figure 2e). The overall reduction in entropy was in line with the theoretical
313 prediction obtained within a simplified model of template-based ligation(19). In that model, the
314 reduction of entropy was due to “mass extinction” in sequence space, with only a very limited
315 (though still exponentially large) set of survivor sequences emerging. In the present experiment
316 related variation in abundances of different sequences did develop but didn’t proceed all the way
317 to extinction.

318 Several patterns can be easily spotted in the pool of surviving sequences. In particular, multimer
319 strands predominantly fell in one of two groups: A-type or T-type each characterized by about 70 %
320 of either base A or T (Figure 2c, d). The initially single-peaked approximately binomial A:T-ratio
321 distribution in random monomers changed into a bimodal one in longer chains. We attribute this
322 separation into two subpopulations to the fact that such composition bias suppresses the formation

323 of internal hairpins and other secondary structures. The self-hybridization reduces the activity of
324 both template and substrate chains leading to a lower rate of ligation. The adaptation by separation
325 into two subpopulations was reproduced by a kinetic model in which activities of reacting strands
326 were corrected for hairpin formation, with realistic account for its thermodynamic cost. This model
327 produced a bimodal distribution of A-content in 24mers, in qualitative agreement with the
328 experimental data. Furthermore, the eventual distribution of longer oligomer lengths could be
329 successfully captured by the maximum entropy distribution, subject to the constraint of fixed
330 average composition of A- and T-type subpopulations. Another remarkable observation is that
331 although formation of hairpins was suppressed through the mechanism above, a small but
332 noticeable fraction of oligomers have extremely long stretches of internal hairpins. The likely
333 mechanisms of their formation are either ligation of a pair of nearly complementary chains from A-
334 type and T-type subpopulations, or self-elongation of such oligomers.

335 Another common pattern was a distinct AT-alternating pattern around the ligation site, as can be
336 seen in Figure 1b. Those AT-alternating motifs first appeared in 24mers, and remained very
337 common in longer chains. These features accounted for some of the reduction in sequence entropy,
338 but did not account for all of the selection at ligation sites, where, as demonstrated by the Z-score
339 analysis, a rich ligation landscape has developed (Figure 4a, b). Not only some 12mers within
340 longer chains were far more abundant than average, but there were also pairs of those that
341 preferentially follow each other, as demonstrated by de Bruijn graphs in Figure 5a, b.

342 We selected a subset of eight pairs of mutually complementary 12mers that appeared anomalously
343 often within longer chains and were well connected within the de Bruijn graph. Using this “Network”
344 subset as a new starting pool, we repeated the temperature-cycling experiment, and compared it
345 to two other reference systems. One of them were eight randomly selected 12mers, the other was
346 artificially designed to promote self-elongation. The resulting multimer population in two out of three
347 of these pools followed a near perfect exponential length profile (Figure 5d). The random pool
348 resulted in a similar behavior to the network one but with significantly lower overall concentration
349 of long chains. Both results are in an excellent agreement with theoretical predictions of
350 reference(18). A higher concentration of long chains generated by network 12mers indicates better
351 overall fitness of this set compared to random 12mers. The “Replicator” set did produce a large
352 number of very long products, presumably by a different mechanism, but a significantly smaller
353 number of products with short and medium lengths. This indicates lower autocatalytic ability in both
354 “Replicator” and “Random” sequence pools when compared to the “Network” pool. In SI section 20
355 the de Bruijn sequence networks for oligomer products show this difference in elongation fitness
356 clearly: while the “Network” sample forms A-type and T-type groups and is well interconnected, the
357 “Replicator” favors only two sequences.

358 For emergence of life on early earth, oligomers needed to spontaneously show an evolution-like
359 behavior and create structure from randomness. *We think this might be difficult for base-by-base*
360 *replication reactions because of the Eigen error catastrophe(48). Emerging strands are either*
361 *accurate copies of the template strand or they become more and more random due to the*
362 *incorporated errors every time a strand is replicated. Thus, the system loses information and*
363 *function over time. But even if the replication fidelity would be below the error threshold and*
364 *replicated strands were perfect copies of the original strand template, the emergence of a fittest*
365 *sequence from a random initial pool would require Darwinian selection of function over a potentially*
366 *very large sequence space.* In contrast, we here followed templated ligation from a pool of random
367 12mer strands made from two bases under temperature oscillations. Both the cooperation of
368 sequences and the usage of ligation instead of base-by-base replication distinguishes this work
369 from(48) and lead to ligated sequences that were highly structured. *Those sequences could*
370 *physically be* selected by length using temperature differences(28, 49–51). This combination of
371 mechanisms would have a dynamics very similar to Darwinian evolution

372 Despite its minimalism, the studied system contains all elements necessary for Darwinian evolution:
373 out of equilibrium conditions, transmission of sequence information from template to substrate

374 strains, reliable reproduction of a subset of oligomer products and the possibility to select from the
375 long fast growing sequences in the process. At the dawn of life, such pre-Darwinian dynamics
376 would have pushed prebiotic systems towards lower entropy states. A subsequent selection for
377 catalytic function from the replicated structured sequences could then have paved the way towards
378 the eventual emergence of life.

379 **Materials and Methods**

380 **Nomenclature**

381 **Oligomer:** a product from the templated ligation reaction with a length of a multiple of 12 nt.
382 **Subsequence:** 12mer long sequence in between two ligation sites or in the beginning or end of a
383 multimer. **Submotif:** a sequence of a certain length x . In contrast to a subsequence, a submotif
384 can start at any position in a mono- or oligomers, not only at ligation sites, or the sequence start.
385 **Ligation site:** in particular, the bond between two monomer or multimer strands. In context of
386 sequence motifs, it refers to the region around this bond (± 1 to 6 bases).

387

388 **Ligation by DNA ligase**

389 For enzymatic ligation of ssDNA a TAQ DNA ligase from *New England Biolabs* was used. Chemical
390 reaction conditions were as stated by the manufacturer: 10 μ M total DNA concentration in 1x ligase
391 buffer. The ligase has a temperature dependent activity and is not active at low (4-10 °C) and very
392 high temperatures (85-95 °C). In our experimental system DNA hybridization characteristics are
393 strongly temperature dependent, as shown in the SI. We expect this to have stronger influence on
394 the overall length distribution and product concentrations than ligase activity, as the timescale of
395 hybridization is significantly longer than the timescale of ligation (compared in SI). The
396 manufacturer provides activity of the ligase in units/ml, specifically: “one unit is defined as the
397 amount of enzyme required to give 50 % ligation of the 12-base pair cohesive ends of 1 μ g of
398 BstEII-digested λ DNA in a total reaction volume of 50 μ l in 15 minutes at 45 °C”.

399

400 **Design of the random sequence pool**

401 The use of a DNA ligase enables very fast ligation with low error rate. But not every DNA system
402 is suitable for templated ligation. As stated by the manufacturer, the TAQ ligase does not ligate
403 overhangs which are 4 nt or shorter. Therefore, the shortest possible length of strands is 10mer,
404 opening up $4^{10} > 10^6$ different monomer sequences. The resulting pool cannot be sequenced to a
405 reasonable extend. We artificially reduced the sequence space by limiting sequences to only
406 include bases adenosine (A) and thymine (T). 10mer strands with random AT sequence have too
407 low melting temperature, in a range where the ligase is not active (compare SI). We found 12mers
408 with random AT sequences to successfully ligate and to produce longer product strands due to
409 their elevated melting temperature. The monomer sequence space is $2^{12}=4096$ is not too large, so
410 that we were able to completely sequence it multiple times.

411 The DNA was produced as 5'-WWWWWWWWWWWW-3' with a 5' POH modification by
412 *biomers.net*. “W” denotes base A or T with the same probability. We analyze the “randomness” of
413 this pool in the SI.

414

415 **Temperature Cycling**

416 Temperature cyclers *Bio-Rad T100*, *Bio-Rad CFX96*, *Analytik Jena qTOWER³* and *Thermo Fisher*
417 *Scientific ProFlex PCR System* were used to apply alternating dissociation and ligation
418 temperatures to our samples. The dissociation temperature of 75 °C was chosen, to melt short
419 initially emerging ssDNA of up to 36mer. In the SI we also show how a variation of the dissociation
420 temperature changes multimer product distribution in a random sequence templated ligation
421 experiment. Lower dissociation temperatures enable us to run several thousand temperature
422 cycles, as the stability of the TAQ DNA ligase is reduced substantially for longer times at 95 °C.
423 Time resolution experiments with PAGE-analysis demonstrated ligase activity even after
424 2000 temperature cycles for a dissociation temperature of 75 °C. In experiments screening the
425 ligation temperature (see SI), we found that for ligation temperatures of 25 °C the product length

426 distribution is exponentially falling. For higher ligation temperatures such as 33 °C we find more
427 long sequences, but almost no 24mer and 36mer sequences. For sequenced samples we chose a
428 ligation temperature of 25 °C because the library preparation kit is better suited for shorter DNA
429 strands. In sequencing data for samples with 33 °C the yield was very low, but the results are
430 similar to the sequencing data of samples with 25 °C ligation temperature, but with comparably
431 worse statistics. For dsDNA dissociation in each temperature cycle the corresponding temperature
432 is held for 20 s with subsequent 120 s at the ligation temperature.
433

434 Sequencing by Next Generation Sequencing (NGS)

435 For sequencing we used the Accel-NGS 1S Plus DNA Library Kit from *Swift Biosciences*. The
436 sequencing was done using a HiSeq 2500 DNA sequencer from *Illumina*. The kit was used as
437 stated in the manufacturer's manual. All volumes were divided by four to achieve more output from
438 a limited supply of chemicals. Library preparation was done in four steps: first a random sequence
439 CT-tail was added to the 3' end of the DNA by (probably, the manufacturer does not give information
440 about this step) a terminal transferase. In a single 15 min ligation step the back primer sequence
441 (starting with AGAT...) was ligated to the 3' end of the random CT-stretch. In the second step a
442 single cycle PCR was used to produce the reverse complement and to leave double stranded DNA
443 with a single A overhang. Step three ligated the start primer to the 5' end of the DNA. Step four
444 added barcode indices to both ends of the DNA by a PCR reaction. This step was done several
445 times to result in the desired amount of DNA for sequencing.
446

447 Sequence Analysis

448 Demultiplexing was done by a standard demultiplexing algorithm on servers of the Gen Center
449 Munich running an instance of Galaxy(52) connected to the sequencing machine. *Illumina*-
450 sequencing creates three FASTA-files, listing the front and the back barcodes and the read
451 sequence, for each lane of the flow cell. The demultiplexing-algorithm matches the barcodes of the
452 prepared library DNA to the read sequence and produces a single FASTA file including the read
453 quality scores.

454 The sequence-data was analyzed with a custom written *LabVIEW* software. The main challenge
455 was to separate the read sequences from the attached primers. The start primer is automatically
456 cut in the demultiplexing step. The end primer is cut with an algorithm based on regular expression
457 (RegEx) pattern matching. With RegEx we first search for multiples of the monomer length. If these
458 structures were followed by at least four bases of C or T followed by the sequence AGAT we
459 concluded that we found a relevant sequence. The 3'-primer was cut and the resulting sequence
460 saved for analysis.

461 RegEx for searching AT random sequences:

```
462 (^[ATCG]{12}|[ATCG]{24}|[ATCG]{36}|[ATCG]{48}|[ATCG]{60}|[ATCG]{72}|[ATCG]{84})(?=[CT]{4,}AGAT))
```

463 RegEx for selecting a maximum of X false reads of G or C in random sequence AT samples:

```
464 ^(?!(?:.*?(G|C)){X,})^[ATCG]{12,}
```

465 The sequenced library may have primer-primer dimers and
466 oligomers as well as partial primers that were falsely built in the library preparation step. As the
467 SWIFT kit is made for longer sequences by design, shorter sequences such as 12mer in our study
468 may have lower yields and larger error rates for the library kit chemistry. Therefore, the inclusion of
469 sequences with a single or multiple false reads can improve the statistics, as long as submotifs with
obviously faulty reads are ignored in the analysis.

470 Acknowledgments

471 The authors would like to acknowledge funding by the Simons Foundation (327125 to D.B.), the
472 Deutsche Forschungsgemeinschaft (DFG, German Research Foundation)– Project-ID 201269156
473 – SFB 1032, the Advanced Grant (EvoTrap #787356) PE3, ERC-2017-ADG from the European
474 Research Council, CRC 235 Emergence of Life (Project-ID 364653263) and the Center for
475 NanoScience (CeNS). We would like to thank Ulrich Gerland, Joachim Rosenberger, Tobias
476 Göppel, and Bernhard Altaner for their helpful remarks and discussions about hybridization
477 energies, baseline corrections and the interpretation of multimer product distributions. Our
478 collaboration with this group on the in-depth analysis of the elongation dynamics is currently

479 submitted for review. P.W.K and D.B. thank Irene Chen and Daniel Duzdevich for help with
480 optimizations of the library preparation protocol and analysis, as well as Stefan Krebs and Marlis
481 Fischalek at the Gene Center Munich for their help with the library preparation and the sequencing
482 of the samples. Additionally, we like to thank Filiz Civril for her extensive comments on the
483 manuscript. This research was partially done at, and used resources of the Center for Functional
484 Nanomaterials, which is a U.S. DOE Office of Science Facility, at Brookhaven National Laboratory
485 under Contract No.~DE-SC0012704.
486
487

488

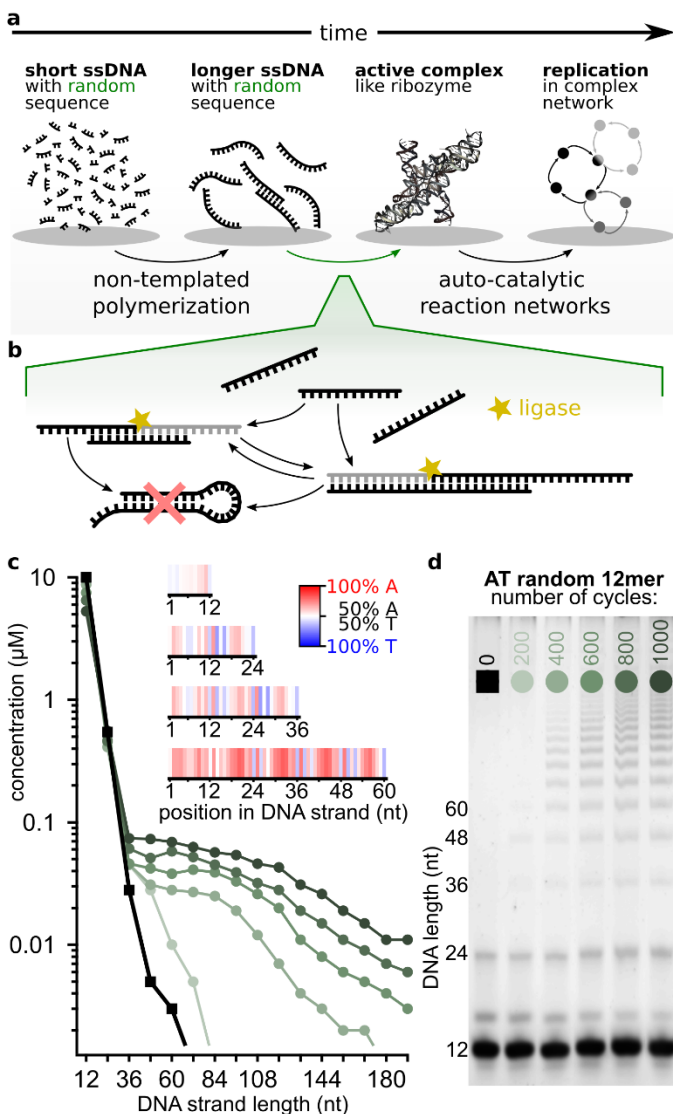
References

- 489 1. F. H. C. Crick, The origin of the genetic code. *J. Mol. Biol.* **38**, 367–379 (1968).
- 490 2. L. E. Orgel, Evolution of the genetic apparatus: A review. *Cold Spring Harb. Symp. Quant.*
491 *Biol.* **52**, 9–16 (1987).
- 492 3. G. Walter, The RNA World. *Nature* **319**, 618 (1986).
- 493 4. J. Attwater, A. Wochner, V. B. Pinheiro, A. Coulson, P. Holliger, Ice as a protocellular
494 medium for RNA replication. *Nat. Commun.* **1**, 1–8 (2010).
- 495 5. G. F. Joyce, Toward an alternative biology. *Science (80-.)*. **336**, 307–308 (2012).
- 496 6. D. P. Horning, G. F. Joyce, Amplification of RNA by an RNA polymerase ribozyme. *Proc.*
497 *Natl. Acad. Sci.* **113**, 9786–9791 (2016).
- 498 7. K. J. Hertel, *et al.*, Numbering system for the hammerhead. *Nucleic Acids Res.* **20**, 3252
499 (1992).
- 500 8. H. W. Pley, K. M. Flaherty, D. B. McKay, Three-dimensional structure of a hammerhead
501 ribozyme. *Nature* **372**, 68–74 (1994).
- 502 9. K. R. Birikh, P. A. Heaton, F. Eckstein, The structure, function and application of the
503 hammerhead ribozyme. *Eur. J. Biochem.* **245**, 1–16 (1997).
- 504 10. W. G. Scott, J. B. Murray, J. R. P. Arnold, B. L. Stoddard, A. Klug, Capturing the structure
505 of a catalytic RNA intermediate: The hammerhead ribozyme. *Science (80-.)*. **274**, 2065–
506 2069 (1996).
- 507 11. J. W. Szostak, D. P. Bartel, Structurally Complex Highly Active RNA Ligases Derived from
508 Random RNA Sequences (1995).
- 509 12. J. A. Kozlov, L. E. Orgel, Nonenzymatic template-directed synthesis of RNA from
510 monomers. *Mol. Biol.* **34**, 921–930 (2000).
- 511 13. J. Oró, Mechanism of synthesis of adenine from hydrogen cyanide under possible
512 primitive earth conditions. *Nature* (1961).
- 513 14. R. Lohrmann, Formation of nucleoside 5'-polyphosphates from nucleotides and
514 trimetaphosphate. *J. Mol. Evol.* (1975).
- 515 15. G. J. Handschuh, R. Lohrmann, L. E. Orgel, The effect of Mg²⁺ and Ca²⁺ on urea-
516 catalyzed phosphorylation reactions. *J. Mol. Evol.* (1973).
- 517 16. R. Österberg, L. E. Orgel, R. Lohrmann, Further studies of urea-catalyzed phosphorylation
518 reactions. *J. Mol. Evol.* (1973).
- 519 17. Z. Liu, *et al.*, Harnessing chemical energy for the activation and joining of prebiotic
520 building blocks. *Nat. Chem.* (2020) <https://doi.org/10.1038/s41557-020-00564-3>.
- 521 18. A. V. Tkachenko, S. Maslov, Spontaneous emergence of autocatalytic information-coding
522 polymers. *J. Chem. Phys.* **143**, 045102 (2015).
- 523 19. A. V. Tkachenko, S. Maslov, Onset of natural selection in populations of autocatalytic
524 heteropolymers. *J. Chem. Phys.* **149** (2018).
- 525 20. H. Fellermann, S. Tanaka, S. Rasmussen, Sequence selection by dynamical symmetry
526 breaking in an autocatalytic binary polymer model. *Phys. Rev. E* **96**, 1–14 (2017).
- 527 21. S. Toyabe, D. Braun, Cooperative Ligation Breaks Sequence Symmetry and Stabilizes
528 Early Molecular Replication. *Phys. Rev. X* **9**, 011056 (2019).
- 529 22. J. M. Horowitz, J. L. England, Spontaneous fine-tuning to environment in many-species
530 chemical reaction networks. *Proc. Natl. Acad. Sci. U. S. A.* (2017).
- 531 23. P. P. Hsu, D. M. Sabatini, Cancer cell metabolism: Warburg and beyond. *Cell* **134**, 703–
532 707 (2008).
- 533 24. P. D. Boyer, The ATP synthase - a splendid molecular machine. *Annu. Rev. Biochem.* **66**,
534 717–749 (1997).
- 535 25. J. A. Baross, S. E. Hoffman, Submarine hydrothermal vents and associated gradient
536 environments as sites for the origin and evolution of life. *Orig. Life Evol. Biosph.* (1985).
- 537 26. R. Pascal, *et al.*, Towards an evolutionary theory of the origin of life based on kinetics and
538 thermodynamics. 1–9 (2013).
- 539 27. H. Mutschler, A. Wochner, P. Holliger, Freeze-thaw cycles as drivers of complex ribozyme
540 assembly. *Nat. Chem.* **7**, 502–508 (2015).
- 541 28. C. B. Mast, D. Braun, Thermal trap for DNA replication. *Phys. Rev. Lett.* **104**, 1–4 (2010).

- 542 29. J. D. Crick, F. H., Watson, The complementary structure of deoxyribonucleic acid. *Proc. R. Soc. London. Ser. A. Math. Phys. Sci.* **223**, 80–96 (1954).
- 543
- 544 30. J. SantaLucia, A unified view of polymer, dumbbell, and oligonucleotide DNA nearest-neighbor thermodynamics. *Proc. Natl. Acad. Sci. U. S. A.* **95**, 1460–1465 (1998).
- 545
- 546 31. T. Wu, L. E. Orgel, Nonenzymic template-directed synthesis on oligodeoxycytidylate sequences in hairpin oligonucleotides. **LII**, 317–322 (1992).
- 547
- 548 32. R. Rohatgi, D. P. Bartel, J. W. Szostak, Kinetic and mechanistic analysis of nonenzymatic, template-directed oligoribonucleotide ligation. *J. Am. Chem. Soc.* **118**, 3332–3339 (1996).
- 549
- 550 33. A. C. Fahrenbach, *et al.*, Common and Potentially Prebiotic Origin for Precursors of Nucleotide Synthesis and Activation. *J. Am. Chem. Soc.* **139**, 8780–8783 (2017).
- 551
- 552 34. L. Li, *et al.*, Enhanced nonenzymatic RNA copying with 2-aminoimidazole activated nucleotides. *J. Am. Chem. Soc.* **139**, 1810–1813 (2017).
- 553
- 554 35. R. Appel, B. Niemann, W. Schuhn, Synthesis of the First Triphosphabutadiene. *Angew. Chem. Inf. Ed. Engl.* **119**, 932–935 (1986).
- 555
- 556 36. D. Sievers, G. Von Kiedrowski, Self-replication of hexadeoxynucleotide analogues: Autocatalysis versus cross-catalysis. *Chem. - A Eur. J.* **4**, 629–641 (1998).
- 557
- 558 37. E. Edeleva, *et al.*, Continuous nonenzymatic cross-replication of DNA strands with in situ activated DNA oligonucleotides. *Chem. Sci.* **10**, 5807–5814 (2019).
- 559
- 560 38. L. Zhou, D. K. O’Flaherty, J. W. Szostak, Assembly of a Ribozyme Ligase from Short Oligomers by Nonenzymatic Ligation. *J. Am. Chem. Soc.* (2020)
- 561
- 562 <https://doi.org/10.1021/jacs.0c06722>.
- 563 39. J. Kim, M. Mrksich, Profiling the selectivity of DNA ligases in an array format with mass spectrometry. **38**, 1–10 (2010).
- 564
- 565 40. G. J. S. Lohman, *et al.*, A high-throughput assay for the comprehensive profiling of DNA ligase fidelity. **44** (2016).
- 566
- 567 41. S. G. Fischer, L. S. Lerman, DNA fragments differing by single base-pair substitutions. *Proc. Nat. Acad. Science, Biochem.* **80**, 1579–1583 (1983).
- 568
- 569 42. R. M. Myers, S. G. Fischer, L. S. Lerman, T. Maniatis, Nearly all single base substitutions in DNA fragments joined to a GC-clamp can be detected by denaturing gradient gel electrophoresis. **13**, 3131–3145 (1985).
- 570
- 571
- 572 43. J. Derr, *et al.*, Prebiotically plausible mechanisms increase compositional diversity of nucleic acid sequences. *Nucleic Acids Res.* **40**, 4711–4722 (2012).
- 573
- 574 44. D. Bartel, J. Szostak, Isolation of new ribozymes from a large pool of random sequences [see comment]. *Science (80-)*. **261**, 1411–1418 (1993).
- 575
- 576 45. N. G. de Bruijn, A combinatorial problem. *Proc. Sect. Sci. K. Ned. Akad. van Wet. te Amsterdam* **49**, 758–764 (1946).
- 577
- 578 46. P. A. Pevzner, H. Tang, M. S. Waterman, An Eulerian path approach to DNA fragment assembly. *Proc. Natl. Acad. Sci. U. S. A.* (2001).
- 579
- 580 47. J. N. Zadeh, *et al.*, NUPACK: Analysis and design of nucleic acid systems. *J. Comput. Chem.* **32**, 170–173 (2011).
- 581
- 582 48. M. Eigen, Selforganization of Matter and the Evolution of Biological Macromolecules. *Naturwissenschaften* **10**, 465–523 (1971).
- 583
- 584 49. M. Kreysing, L. Keil, S. Lanzmich, D. Braun, Heat flux across an open pore enables the continuous replication and selection of oligonucleotides towards increasing length. *Nat. Chem.* **7**, 203–208 (2015).
- 585
- 586
- 587 50. L. Keil, M. Hartmann, S. Lanzmich, D. Braun, Probing of molecular replication and accumulation in shallow heat gradients through numerical simulations. *Phys. Chem. Chem. Phys.* **18**, 20153–20159 (2016).
- 588
- 589
- 590 51. M. Morasch, *et al.*, Heated gas bubbles enrich, crystallize, dry, phosphorylate and encapsulate prebiotic molecules. *Nat. Chem.* **11**, 779–788 (2019).
- 591
- 592 52. E. Afgan, *et al.*, The Galaxy platform for accessible, reproducible and collaborative biomedical analyses: 2018 update. *Nucleic Acids Res.* (2018).
- 593
- 594

595
596

Figures and Tables



597
598

599 **Figure 1.** Autocatalytic templated ligation of DNA 12mers.

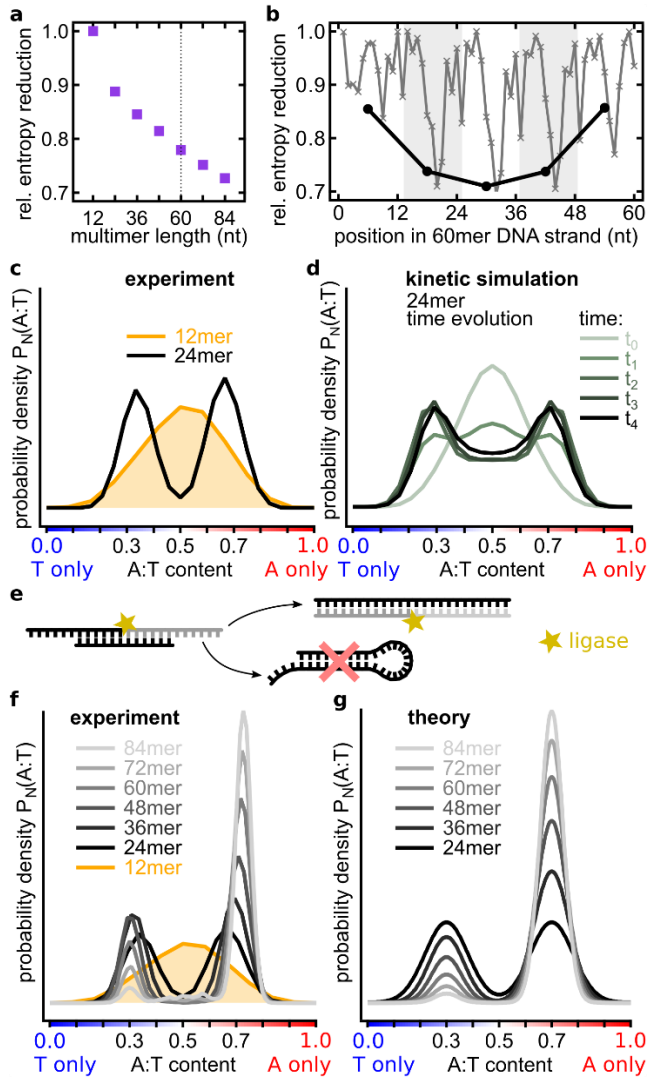
600 **a** Before cells evolved, the first ribozymes were thought to perform basic cell functions. In the
601 exponentially vast sequence space, spontaneous emergence of a functional ribozyme is highly
602 unlikely, therefore pre-selection mechanisms were likely necessary.

603 **b** In our experiment, DNA strands hybridize at low temperatures to form 3D complexes which can
604 be ligated and preserved in the high temperature dissociation steps. The system self-selects for
605 sequences with specific ligation site motifs as well as for strands that continue acting as templates.
606 Hairpin sequences are therefore suppressed.

607 **c** Concentration analysis shows progressively longer strands emerging after multiple temperature
608 cycles. The inset (A-red, T-blue) shows that while 12mers (88009 strands) have essentially random
609 sequences (white), various sequence patterns emerge in longer strands (60mers, 235913 strands
610 analyzed).

611 **d** Samples subjected to different number (0-1000) of temperature cycles between 75 °C and 33 °C.
612 Concentration quantification is done on PAGE with SYBR post-stained DNA

613
614



615
616

617 **Figure 2.** Hairpin formation amplifies selection into A-rich and T-rich sequences.
618 **a** Relative entropy reduction as a function of multimer product length: 1 – a random pool and 0 – a
619 unique sequence.

620 **b** Relative entropy reduction of 60mer products. Black: Entropy reduction of 12 nt subsequences
621 compared to a random sequence strand of the same length. Grey: Entropy reduction at each
622 nucleotide position showing positional dependence.

623 **c** A gradual development of the bimodal distribution of A:T ratio in chains of different lengths. While
624 the A:T ratio in 12mers has a single-peaked nearly binomial distribution, 24mers already have a
625 clearly bimodal distribution peaked at 65:35 % (A-type strands) and 35:65 % (T-type strands) A:T
626 ratios.

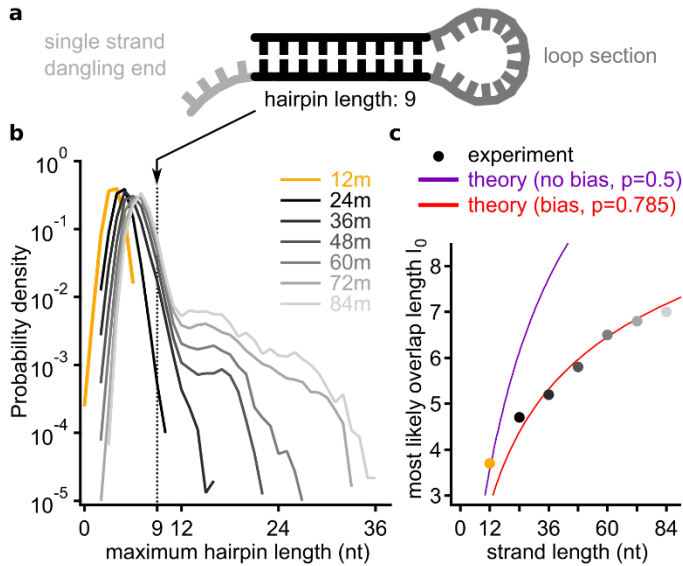
627 **d** Emergence of a bimodal distribution in a kinetic model of templated ligation.

628 **e** Sequences with nearly balanced A:T ratios are prone to formation of hairpins. In the model in d
629 and the experiment, these hairpins prevent strands from acting as templates and substrates for
630 ligation reactions thereby suppressing the central part of the distribution.

631 **f** A:T ratio distributions in strands of different length. As length increases A-type strands become
632 progressively more abundant in comparison to T-type strands.

633 **g** A:T ratio distributions in a phenomenological model taking into account a slight AT-bias in the
634 initial 12mer pool resemble experimentally measured ones (panel e).

635



636

637

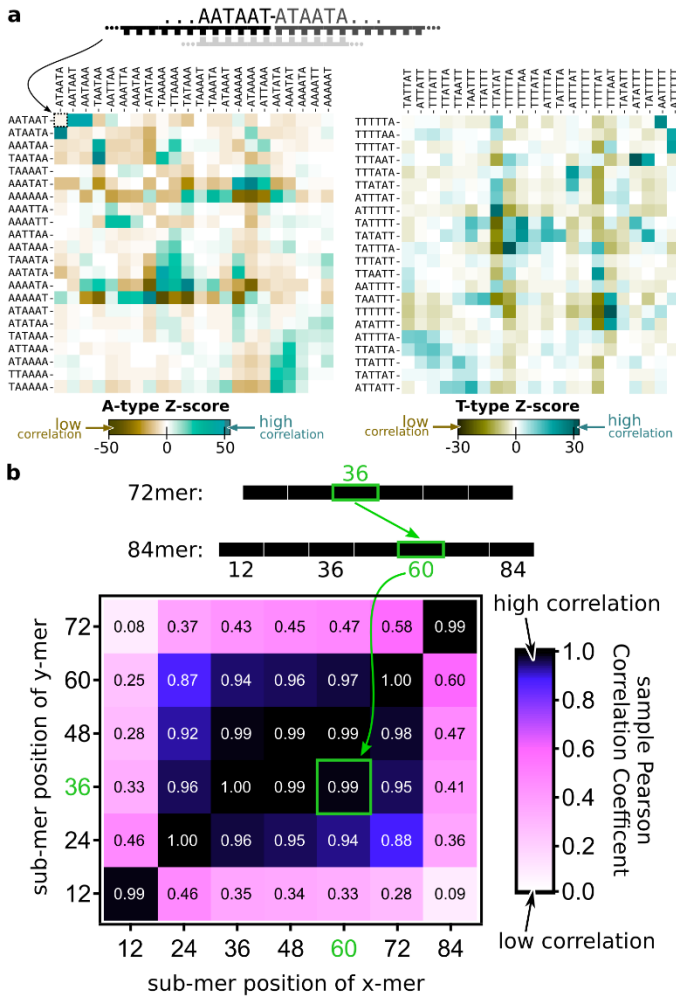
638 **Figure 3.** Large scale entropy reduction and sequence correlation per strand.

639 **a** Sketch of a single strand DNA secondary structure folding on itself, called hairpin. The double
640 stranded part is very similar to a standard duplex DNA.

641 **b** Comparing the PDFs of the maximum hairpin stem length for all strands reveals a group of peaks
642 at around 4 to 7 nt, increasing with the DNA length. Starting for 48mers, there is a tail visible: these
643 self-similar strands are more abundant, the longer the product grows (compare A:T fraction close
644 to $p=0.5$ in Figure 2c).

645 **c** The peak-positions as function of the product length follow equation (3). The unbiased 12mers
646 are on the curve with coefficient $p=0.5$, whereas the products starting from 36mers lay on the curve
647 with $p=0.785$. The bias parameter p is derived from the PDFs in Figure 2d and describes the A:T-
648 ratio in the strand.

649

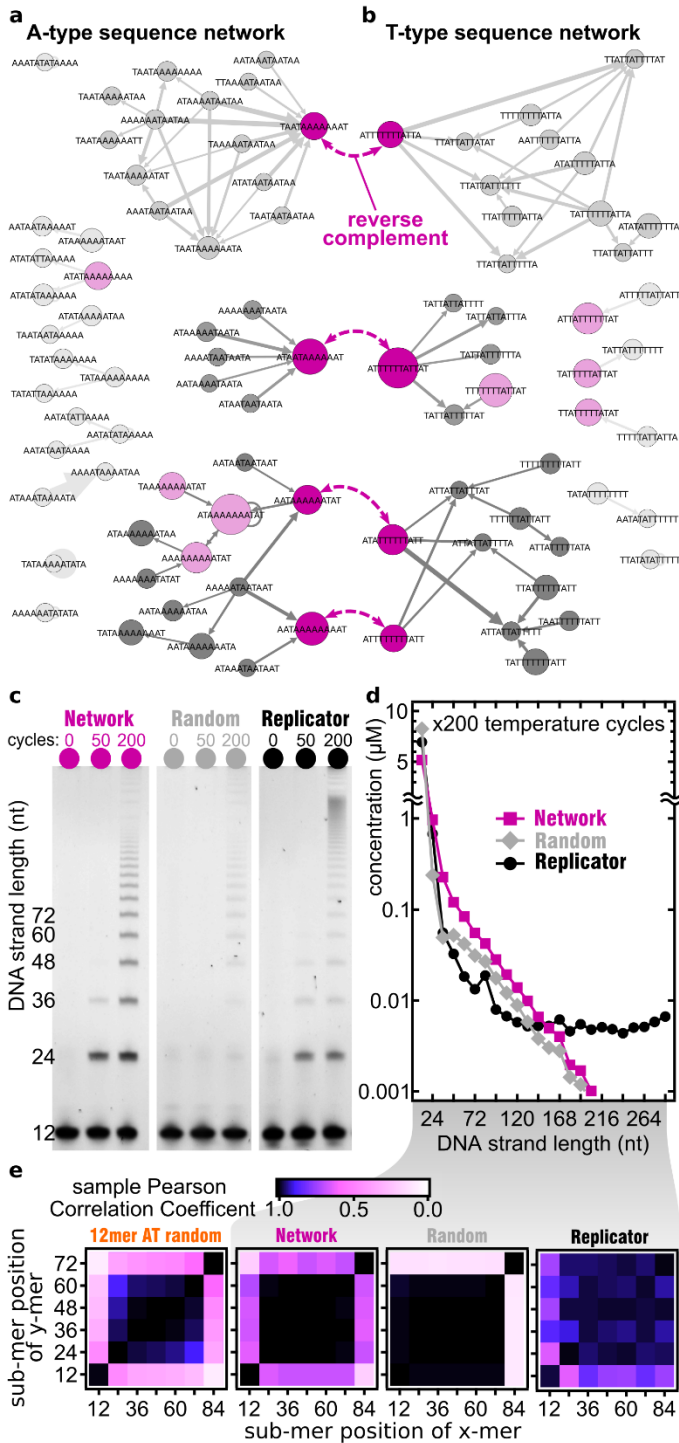


650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667

Figure 4. Emergent landscape of junction sequences.

a The heatmap of Z-scores quantifying the probability to find a junction between a 6 nt sequence listed in rows followed by the 6 nt sequence listed in columns compared to finding it by pure chance and normalized by the standard deviation. Z-scores were calculated for the junction between 4th to the 5th 12mers in 72mers of A-type (left) and T-type (right) respectively. Other internal junctions in all long chains form very similar landscapes composed of over- (teal) and under-represented (ocher) sequences and described in detail in the text. T-type sequences complementary to A-type sequences correspond to the 90° clockwise rotation of the left panel (note a similarity of landscapes in two panels after this transformation).

b The matrix of sample Pearson Correlation coefficients between abundances of 12mers in different positions (1 to 6) inside 72mers (rows) and 84mers (columns). Light regions mark low correlations, dark regions mark high correlations. Very high correlations (>0.9) at the center of the table mean that very similar sequences get selected at all internal positions of chains of different lengths. Different selection pressures operate on the first 12mer and the last 12mer of a chain, yet their sequences are similar in chains of different lengths.



668
669
670
671
672
673
674
675

Figure 5. Testing self-selection with custom sequence pools.

a The de Bruijn graph of overrepresented sequence motifs between consecutive 12mers found in long oligomers. All internal junctions of A-type sequences >48 nt are shown, except the first and the last. All analyzed strands have a Z-score >30 and are sequenced at least 20 times.

b The same de Bruijn graph but for T-type sequences with Z-score >15 and sequenced at least 10 times. Four pairs of most common reverse complementary 12mers are connected by purple dashed

676 arrows. In each network three families with distinctly similar patterns are observed, that each
677 include at least one of the complementary strands. Node sizes reflect relative abundance of
678 12mers, edge thickness denotes the Z-score of the junction between nodes it connects. Light and
679 dark magenta-colored nodes are eight most abundant 12mers in each of two networks.

680 **c** PAGE images of templated ligation of three different samples of 12mers after different number of
681 temperature cycles (columns): “Replicator”: four substrate 12mers and four template 12mers
682 artificially designed for templated ligation, as explained in SI, “Random”: eight random sequence
683 12mers randomly selected from the 4096 possible AT-only 12mers, “Network”: four most common
684 12mers from A-type and another four of T-type shown in dark magenta color in panel a.

685 **d** After 200 temperature cycles, the “Replicator” shows a consistently higher product concentration
686 for all lengths followed by the “Network” sample and then by the “Random” subsamples. In the
687 “Network” and “Random” samples the length distribution above 48nt is well described by an
688 exponential distribution as predicted in Ref. (18).

689 **e** Pearson correlation matrices between 12mer abundances within 72mers and 84mers in each
690 sample (same as in Figure 4b). While the pattern of correlations in the “Network” sample (second
691 from left) resembles that shown in Figure 4b (reproduced in the leftmost subpanel), the “Random”
692 sample (second from right) singles out the last 12mer but not the first one. The “Replicator” sample
693 (the rightmost subpanel) has its own distinct self-similar pattern of correlations.