

# 1 Ligation of random oligomers leads to emergence of autocatalytic 2 sequence network

3 Patrick W. Kudella<sup>1</sup>, Alexei V. Tkachenko<sup>2</sup>, Sergei Maslov<sup>3,4</sup>, Dieter Braun\*<sup>1</sup>

4 <sup>1</sup>Systems Biophysics and Center for NanoScience, Ludwigs-Maximilian-Universität München, 80799 Munich, Germany

5 <sup>2</sup>Center for Functional Nanomaterials, Brookhaven National Laboratory, Upton, New York 11973, USA

6 <sup>3</sup>Department of Bioengineering, University of Illinois at Urbana-Champaign, Urbana, Illinois 61801, USA

7 <sup>4</sup>Carl R. Woese Institute for Genomic Biology, University of Illinois at Urbana-Champaign, 1206 West Gregory Drive, Urbana  
8 Illinois 61801, USA

9

## 10 ABSTRACT

11 The emergence of longer information-carrying and functional nucleotide polymers from random short  
12 strands was a major stepping stone at the dawn of life. But the formation of those polymers under  
13 temperature oscillation required some form of selection. A plausible mechanism is template-based  
14 ligation where theoretical work already suggested a reduction in information entropy.

15 Here, we show how nontrivial sequence patterns emerge in a system of random 12mer DNA sequences  
16 subject to enzyme-based templated ligation reaction and temperature cycling. The strands acted both  
17 as a template and substrates of the reaction and thereby formed longer oligomers. The selection for  
18 templating sequences leads to the development of a multiscale ligation landscape. A position-  
19 dependent sequence pattern emerged with a segregation into mutually complementary pools of A-  
20 rich and T-rich sequences. Even without selection for function, the base pairing of DNA with ligation  
21 showed a dynamics resembling Darwinian evolution.

## 22 BACKGROUND

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

31 The problem of non-enzymatic formation of single base nucleotides and short oligomers in settings  
32 reminiscent of the primordial soup has been studied before<sup>12-16</sup>. However, the continuation of this  
33 evolutionary path towards early replication networks would require a pre-selection mechanism of  
34 oligonucleotides (as shown in Fig. 1a), lowering the information entropy of the resulting sequence  
35 pool<sup>17-20</sup>. In principle, such selection modes include optimization for information storage, local  
36 oligomer enrichment e.g. in hydrogels or in catalytically functional sites.

37 An important aspect of a selection mechanism is its non-equilibrium driving force. Today's highly  
38 evolved cells function through multistep and multicomponent metabolic pathways like glycolysis in the  
39 Warburg effect<sup>21</sup> or by specialized enzymes like ATP synthase which provide energy-rich adenosine  
40 triphosphate (ATP)<sup>22</sup>. In contrast, it is widely assumed<sup>3,4,23-26</sup> that selection mechanisms for molecular  
41 evolution at the dawn of life must have been much simpler, e.g. mediated by random binding between  
42 biomolecules subject to non-equilibrium driving forces such as fluid flow and cyclic changes in  
43 temperature.

44 Here, we explored the possibility of a significant reduction of sequence entropy driven by templated  
45 ligation<sup>17</sup> and mediated by Watson-Crick base pairing<sup>27</sup>. Starting from a random pool of  
46 oligonucleotides we observed a gradual formation of longer chains showing reproducible sequence  
47 landscape inhibiting self-folding and promoting templated ligation. Here we argue, that base pairing  
48 combined with ligation chemistry, can trigger processes that have many features of the Darwinian  
49 evolution.

50 As a model oligomer we decided to use DNA instead of RNA since the focus of our study is on base  
51 pairing which is very similar for both<sup>28</sup>. We start our experiments with a random pool of 12mers formed  
52 of bases A (adenosine) and T (thymine). This binary code facilitates binding between molecules and  
53 allows us to sample the whole sequence space in microliter volumes ( $2^{12} \ll 10 \mu\text{M} * 20 \mu\text{l} = 10^{14}$ ).

54 Formation of progressively longer oligomers from shorter ones requires ligation reactions, a method  
55 commonly employed in hairpin-mediated RNA and DNA replication<sup>29,30</sup>. At the origin of life, this might  
56 have been achieved by activated oligomers<sup>31,32</sup> or activation agents<sup>33-35</sup>. Our study is focused on  
57 inherent properties of self-assembly by base pairing in random pools of oligomers and not on chemical  
58 mechanisms of ligation. Hence, we decided to use TAQ DNA ligase - an evolved enzyme for templated  
59 ligation of DNA<sup>19</sup>. This allowed for fast turnovers of ligation and enabled the observation of sequence  
60 dynamics.

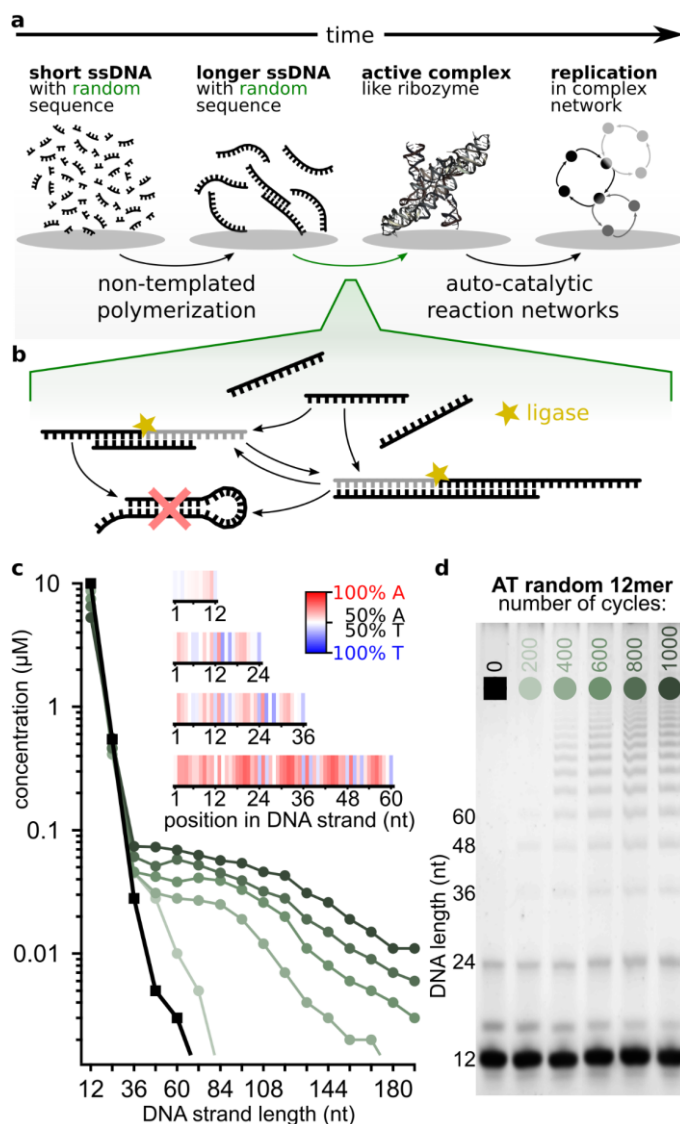
## 61 RESULTS

62 To test templated elongation of polymers in pools of random sequence oligomers, we prepared a  
63 10  $\mu\text{M}$  solution of 12mer DNA strands composed of nucleotides A and T (sequence space: 4096) and  
64 subjected it to temperature cycling, similar to reference<sup>19</sup> with 20 s at denaturation temperature of  
65 75 °C and 120 s at ligation temperature of 33 °C. Temperatures were selected according to the melting  
66 dynamics of the DNA pool; the time steps were prolonged relative to Toyabe and Braun (SI section 5.3)  
67 because of a greater sequence space. The sample was split into multiple tubes and exposed to 200,  
68 400, 600, 800, 1000 temperature cycles, with one tube kept at 4 °C for reference.

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

78 An important property of the initial monomer-pool is its sequence content. Although for pools with  
79 lower sequence complexity it is possible to show different strand compositions using PAGE<sup>36,37</sup>, a large  
80 size of our “monomer” ( $2^{12}=4096$ ) and 24mer product pools (sequence space:  $2^{24}\approx 16.8\times 10^6$ ) excludes  
81 this approach. Thus, we analyzed our final products by Next Generation Sequencing (NGS) to get  
82 insights into product strand compositions.

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



88

89 **Fig. 1, Autocatalytic templated ligation of DNA 12mers.**

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

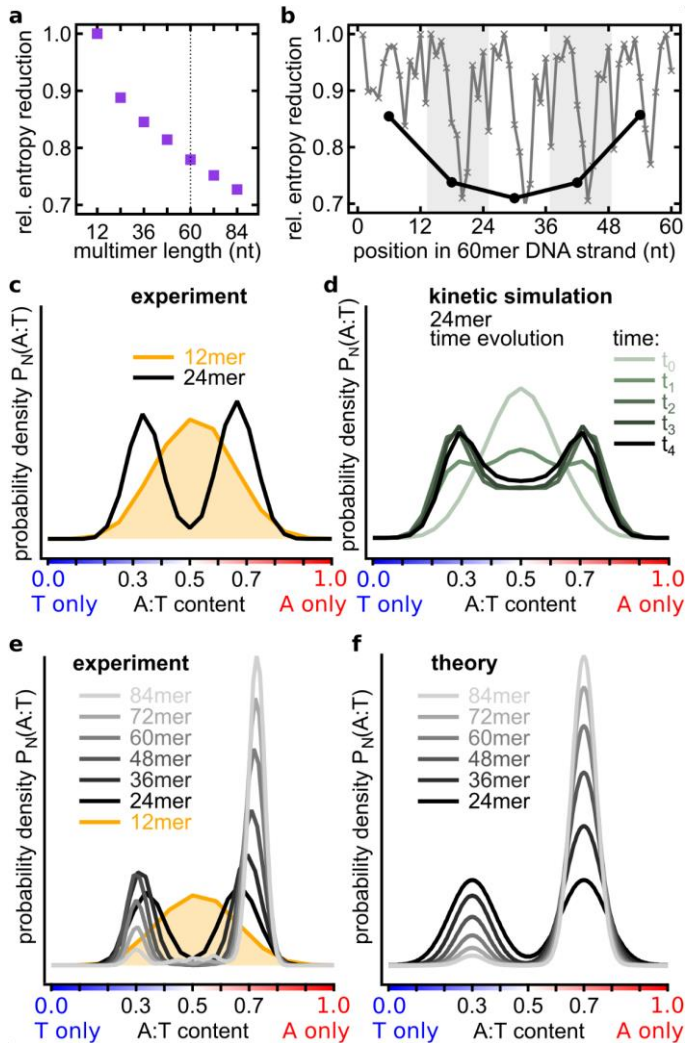
93 *b* In our experiment, DNA strands hybridize at low temperatures to form 3D complexes which can be ligated and preserved  
 94 in the high temperature dissociation steps. The system self-selects for sequences with specific ligation site motifs as well as  
 95 for strands that continue acting as templates.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

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

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

101 The information entropy of longer chains is expected to be smaller than the entropy of a random  
 102 sequence strand of the same length, if some sort of selection mechanism is involved<sup>17</sup>. We analyzed  
 103 the entropy reduction for different lengths of products (Fig. 2a) as well as the positional dependence  
 104 of the single base entropy for 60mer products (Fig. 2b). The relative entropy reduction is similar to one  
 105 used in Derr *et al.*<sup>18</sup> where 1 describes a completely random ensemble and 0 an ensemble of only one  
 106 sequence. Entropy reduction was observed in all analyzed product lengths with a greater reduction  
 107 observed for longer oligomer lengths. The entropy of each 12mer subsequence was also found to be  
 108 significantly lower than that of random 12mers (Fig. 2b, black line). The central subsequence had the  
 109 lowest entropy while 12mers located at both ends of chains had relatively higher entropies. This

110 behavior was also observed as a function of nucleotide position within a 12mer suggesting a multi-  
 111 scale pattern of entropy reduction.



112

113 **Fig. 2, Hairpin formation amplifies selection into A-rich and T-rich sequences.**

114 **a** Relative entropy reduction as a function of multimer product length: 1 – a random pool and 0 – a unique sequence.

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

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

120 **d** Emergence of a bimodal distribution in a kinetic model of templated ligation. Sequences with nearly balanced A:T ratios  
 121 are prone to formation of hairpins. In the model these hairpins prevent strands from acting as templates and substrates for  
 122 ligation reactions thereby suppressing the central part of the distribution.

123 **e** A:T ratio distributions in strands of different length. As length increases A-type strands become progressively more  
 124 abundant in comparison to T-type strands.

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

127 In the initial pool of random 12mers the A-to-T ratio distribution is shaped binomially, as expected for  
 128 a random distribution. However, it dramatically shifted for 24mer products of ligation: a bimodal  
 129 distribution of about 65:35 % A:T (A-type) as well as the inverse, 35:65 % A:T (T-type) was observed  
 130 with 24mer products (Fig. 2c). DNA strands composed of only two complementary bases are more  
 131 prone to formation of single-strand secondary structures like hairpins than DNAs composed of all four  
 132 bases. In our templated ligation reaction, we expected that hairpin-sequences are not elongated and  
 133 also not used as template-strands because they form catalytically passive Watson-Crick-base-paired

134 configuration. A bimodal AT-ratio distribution (Fig. 2d) also emerged in a kinetic computational model  
135 in which a pool of random 12mers was seeded with a small initial amount of random sequence 24mers.  
136 24mers that formed hairpins could not act as templates and were therefore less likely to be reproduced  
137 (see SI for details of this model, section 18.2).

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

$$142 \quad \int P_N(A:T) d_{A:T} = 1. \quad (1)$$

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

147 As predicted theoretically<sup>39</sup>, the eventual length distribution is approximately exponential. A small A-  
148 T bias leads to the respective average chain lengths,  $\bar{N}_A$  and  $\bar{N}_T$ , to be somewhat different for the two  
149 subpopulations. As a result, the bias gets strongly amplified with increasing chain length:

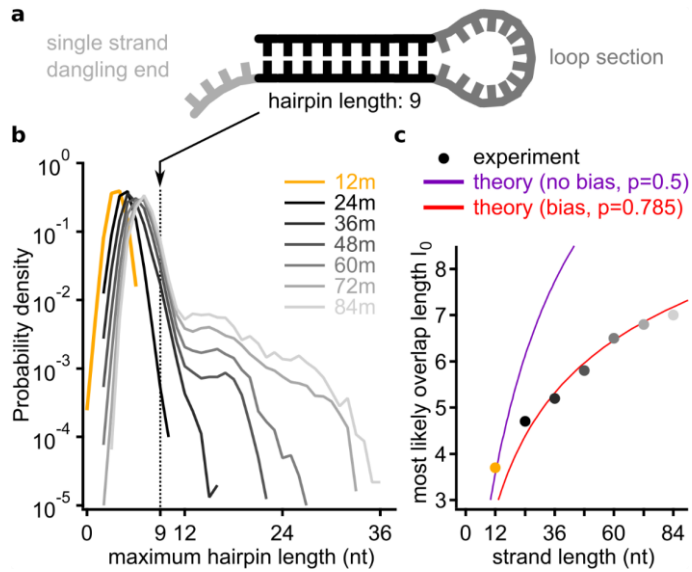
$$150 \quad P_N(A:T) \sim \exp\left(-N\left(\frac{1}{\bar{N}_A} - \frac{1}{\bar{N}_T}\right)\right) = \beta^{-N/12}. \quad (2)$$

151 A simple phenomenological model can successfully capture the major features of the observed A:T  
152 PDFs for multiple chain lengths. Specifically, we assume both A-type and T-type sub-populations -to  
153 maximize the sequence entropy, subject to the constraint that the average A:T content is shifted from  
154 the midpoint (50:50 % composition), by values  $\pm x_0$ , respectively. This model presented in SI  
155 section 18.1, results in a distribution that strongly resembles experimental data, as shown in Fig. 2e-f  
156 A:T profiles for all chain length are fully parameterized by only two fitting parameters:  $\beta = 0.785$ , and  
157  $x_0 = 0.2$ .

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

$$166 \quad N = 2l_0 + \sqrt{2}(2p(1-p))^{-l_0/2}. \quad (3)$$

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



171

172 **Fig. 3, Large scale entropy reduction and sequence correlation per strand.**

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

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

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

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

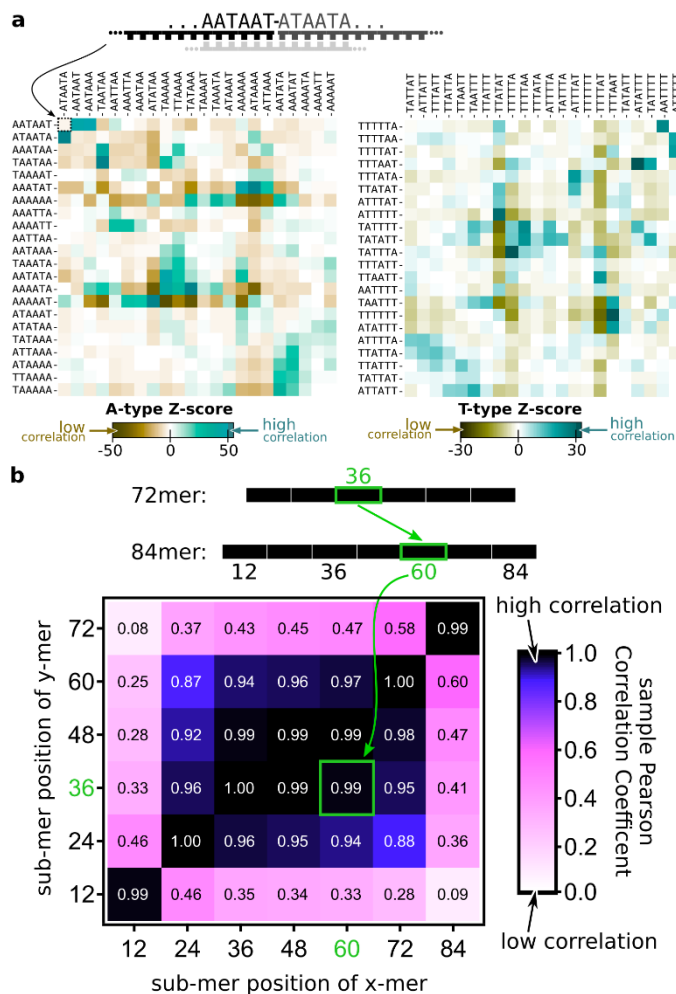
187 The mechanism of formation of these self-binding sequences may involve recombination of shorter A-  
188 type and T-type chains, or self-elongation of shorter hairpins. In either case, the hairpin sequence  
189 cannot efficiently reproduce by means of template ligation. However, the remainder of the pool would  
190 keep producing them as byproduct. Ironically, for the templated ligation reaction this is a possible  
191 failure mode, but those long hairpins may play a key role in the context of origin of life, as precursors  
192 of functional motifs. For instance, work by Bartel and Szostak<sup>11,40</sup> identifies RNA self-binding as crucial  
193 for the direct search of ribozymes – those molecules need to fold into non-trivial secondary structures  
194 to gain their catalytic function.

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

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

201 Fig. 4a shows Z-score heatmaps for junctions within A-type (left panel) and T-type (right panel)  
202 subpopulations. More specifically, we show sequences left (row) and right (column) of the junction  
203 between the 4th to the 5th 12mers in the respective 72mer. These heatmaps reveal a complex  
204 landscape of over- and under-represented junction motifs shown respectively in dark-teal and dark-

205 other colors. Emergence of such complex landscape has been theoretically predicted in Ref. <sup>17</sup>  
 206 Landscape peaks include repeating A-T motif of alternating bases crossing the ligation site (dark-teal  
 207 peak near the center of each of both heatmaps). Relatively rare motifs (valleys) correspond to poly-A  
 208 and poly-T sequences extending across the junction (dark-ocher areas). One exception to this rule is a  
 209 relatively abundant poly-A motif at the bottom right of the A-type heatmap (light-teal). Interestingly,  
 210 these junction sequences had AT-patterns in the beginning of the “left side” and the end of the “right  
 211 side”. This might provide a clue to the origin of these “abnormal” junction motifs. Indeed, they may  
 212 have been templated by abundant poly-T sequences in the middle of T-type 12mers flanked by  
 213 alternating A-T motifs. In other words, junctions at templates of poly-A junction motifs may have been  
 214 shifted by 6 nt relative to substrates. Actually, substrates have no restriction on where they may  
 215 hybridize on a long template and might happen to have their ligation site in the region of poly-T of the  
 216 template strand. We call this “ligation site shift”, as explained in SI section 16. Other preferred junction  
 217 subsequences include repetitions of the AAT motif across the junction (the dark-teal peak in the upper  
 218 left corner of the left panel).



219

220 **Fig. 4, Emergent landscape of junction sequences.**

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

228 *b* The matrix of sample Pearson Correlation coefficients between abundances of 12mers in different positions (1 to 6) inside  
229 72mers (rows) and 84mers (columns). Light regions mark low correlations, dark regions mark high correlations. Very high  
230 correlations (>0.9) at the center of the table mean that very similar sequences get selected at all internal positions of chains  
231 of different lengths. Different selection pressures operate on the first 12mer and the last 12mer of a chain, yet their  
232 sequences are similar in chains of different lengths.

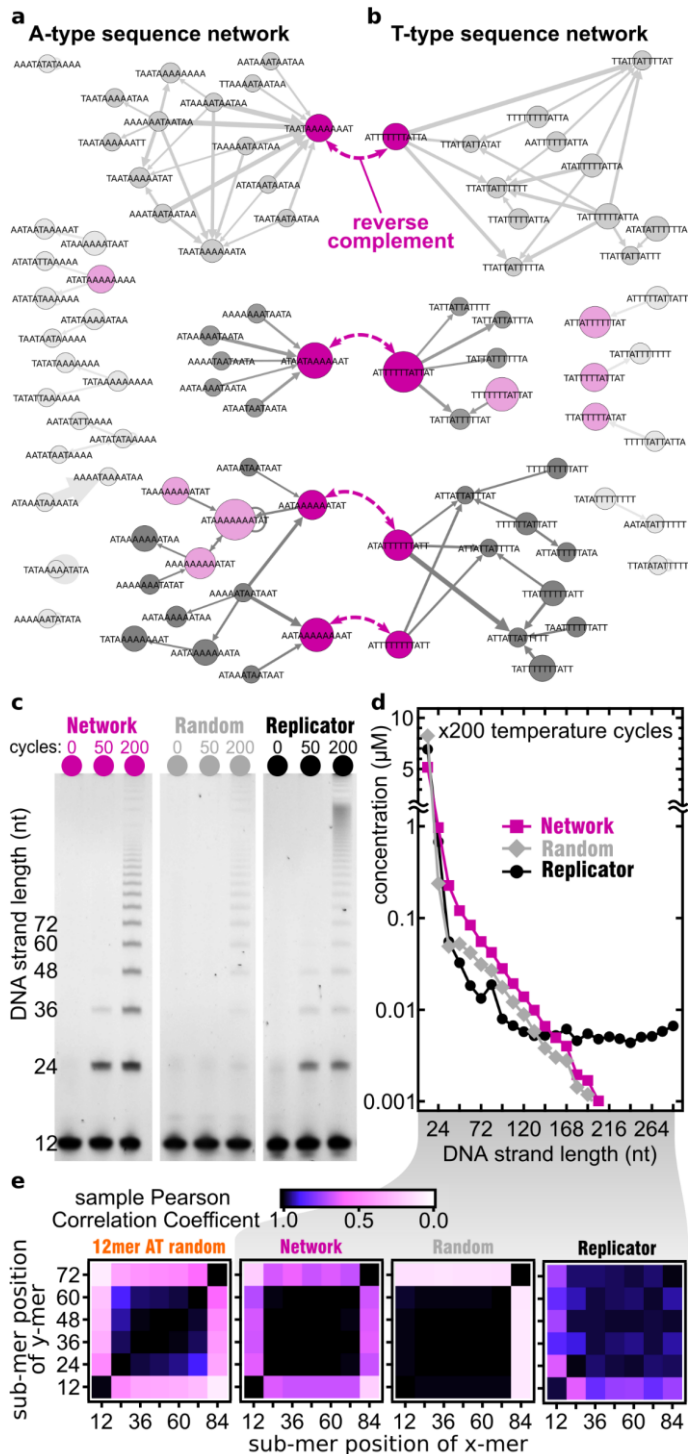
233 How similar are selective pressures operating on sequences of different 12mers within longer chains?  
234 Fig. 4b quantifies this similarity in terms of sample Pearson-Correlation-Coefficient (sPCC) between  
235 abundances of 12mer sequences in different positions of long chains of different lengths.  
236 We compare the abundances of  $2^{12}=4096$  possible 12mer sequences in positions 1 to 6 within all  
237 72mers and compare them to each other and abundances of 12mers in positions 1 to 7 in all 84mers.  
238 Similar results were obtained for other chains longer than 36 nt. A rectangle of very high correlations  
239 (>0.9) at the center of the table in Fig. 4b means that very similar sequences get selected at all internal  
240 positions of all chains (note that only chains longer than 36nt have such internally positioned 12mers).  
241 However, the light border of the table means that a rather different subset of 12mers gets selected in  
242 the first and the last position of a multimer. Whatever the nature of selection pressure acting on these  
243 12mers, it is consistent across oligomers of different lengths as manifested by the high correlation in  
244 the lower left and the upper right corner of the table in Fig. 4b.

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

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

256 To further explore selection capabilities of templated ligation as a function of 12mer sequences in the  
257 initial pool we conducted three additional experiments referred to as “Replicator”, “Random” and  
258 “Network”. The “Random” experiment started with eight randomly chosen 12 nt sequences served as  
259 a control. In the “Replicator” experiment the pool consisted of eight 12 nt sequences artificially  
260 designed for efficient elongation (see below). In the “Network” experiment we populated the pool  
261 with eight naturally selected 12 nt sequences commonly found as subsequences of long strands in our  
262 original ligation experiment with 4096 12mers. To identify these 12mers, we built a network of the  
263 most common 12mers found in A-type oligomers with length of more than 48 nt. This network does  
264 not include the first and the last 12mers, in a multimer as those are known to be statistically different  
265 from the internal ones (see Fig. 4b). The circles in Fig. 5a represent unique 12 nt subsequences while  
266 their size describes their Z-scores quantifying their abundance in long chains. The width of the  
267 connecting line describes the probability that two subsequences are found one after another in a  
268 multimer. The same is done for T-type sequences (Fig. 5b). This representation of a polymer is known  
269 as de Bruijn graph<sup>41</sup> and has been commonly used in DNA fragment analysis and genome assembly<sup>42</sup>  
270 and more recently in the context of templated ligation<sup>17</sup>.





271

272 **Fig. 5, testing self-selection with custom sequence pools.**

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

276 **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  
277 common reverse complementary 12mers are connected by purple dashed arrows. In each network three families with  
278 distinctly similar patterns are observed, that each include at least one of the complementary strands. Node sizes reflect  
279 relative abundance of 12mers, edge thickness denotes the Z-score of the junction between nodes it connects. Light and dark  
280 magenta-colored nodes are eight most abundant 12mers in each of two networks.

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

285 **d** After 200 temperature cycles, the “Replicator” shows a consistently higher product concentration for all lengths followed

286 *by the “Network” sample and then by the “Random” subsamples. In the “Network” and “Random” samples the length*  
287 *distribution above 48nt is well described by an exponential distribution as predicted in Ref<sup>39</sup>.*  
288 *e Pearson correlation matrices between 12mer abundances within 72mers and 84mers in each sample (same as in Fig. 4b).*  
289 *While the pattern of correlations in the “Network” sample (second from left) resembles that shown in Fig. 4b (reproduced in*  
290 *the leftmost subpanel), the “Random” sample (second from right) singles out the last 12mer but not the first one. The*  
291 *“Replicator” sample (the rightmost subpanel) has its own distinct self-similar pattern of correlations.*

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

305 The length distribution of oligomers (Fig. 5d) with concentrations quantified from the PAGE gel image  
306 (Fig. 5c) shows that the “Network” sample produced the most product, as the remaining 12mer  
307 sequence concentration was reduced below two other samples down to almost 5  $\mu\text{M}$ . The length  
308 distribution in both “Random” and “Network” samples is well described by a piecewise-linear  
309 distribution predicted in Ref<sup>39</sup>. For short product lengths ranging between 48mers up to 136mers the  
310 “Random” sample produced more oligomers than the “Replicator” sample. However, for even longer  
311 strands, the “Replicator” sample generated the largest number of really long strands since its length  
312 distribution reached a plateau around 120mers. This is probably due to the nature of the eight-  
313 sequences pools used here with the “Replicator” one made to form well aligned dsDNA that can be  
314 properly ligated. According to NUPACK<sup>43</sup>, 12mers in the “Random” sample should not form any  
315 complexes that could be subsequently ligated by the TAQ ligase. However, our results shown in Fig. 5c  
316 prove the existence of extensive ligation even in the “Random” sample. Presumably, it was initially  
317 triggered by small concentration of complexes formed with low probability, which were subsequently  
318 amplified due to the exponential growth of longer strands in our experiment, just like in the “Network”  
319 sample.

## 320 DISCUSSION

321 We experimentally studied templated ligation in a pool of 12mers made of A and T bases with all  
322 possible sequences ( $2^{12}=4096$ ), subjected to multiple temperature cycles. To accelerate hypothetical  
323 spontaneous ligation reactions operating in the prebiotic world, we employed TAQ DNA ligase in our  
324 experiments. This process produced a complex and heterogeneous ensemble of oligomer products. By  
325 performing the “next generation sequencing” of these oligomers, we found that long strands in this  
326 ensemble have a significantly lower information entropy compared to a random set of oligomers of  
327 the same length. This effect became increasingly more pronounced for longer oligomers (Fig. 2e). The  
328 overall reduction in entropy was in line with the theoretical prediction obtained within a simplified  
329 model of template-based ligation<sup>17</sup>. In that model, the reduction of entropy was due to “mass  
330 extinction” in sequence space, with only a very limited (though still exponentially large) set of survivor  
331 sequences emerging. In the present experiment related variation in abundances of different sequences  
332 did develop but didn’t proceed all the way to extinction.

333 Several patterns can be easily spotted in the pool of surviving sequences. In particular, multimer  
334 strands predominantly fell in one of two groups: A-type or T-type each characterized by about 70 % of  
335 either base A or T (Fig. 2c, d). The initially single-peaked approximately binomial A:T-ratio distribution  
336 in random monomers changed into a bimodal one in longer chains. We attribute this separation into  
337 two subpopulations to the fact that such composition bias suppresses the formation of internal  
338 hairpins and other secondary structures. The self-hybridization reduces the activity of both template  
339 and substrate chains leading to a lower rate of ligation. The adaptation by separation into two  
340 subpopulations was reproduced by a kinetic model in which activities of reacting strands were  
341 corrected for hairpin formation, with realistic account for its thermodynamic cost. This model  
342 produced a bimodal distribution of A-content in 24mers, in qualitative agreement with the  
343 experimental data. Furthermore, the eventual distribution of longer oligomer lengths could be  
344 successfully captured by the maximum entropy distribution, subject to the constraint of fixed average  
345 composition of A- and T-type subpopulations. Another remarkable observation is that although  
346 formation of hairpins was suppressed through the mechanism above, a small but noticeable fraction  
347 of oligomers have extremely long stretches of internal hairpins. The likely mechanisms of their  
348 formation are either ligation of a pair of nearly complementary chains from A-type and T-type  
349 subpopulations, or self-elongation of such oligomers.

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

357 We selected a subset of eight pairs of mutually complementary 12mers that appeared anomalously  
358 often within longer chains and were well connected within the de Bruijn graph. Using this “Network”  
359 subset as a new starting pool, we repeated the temperature-cycling experiment, and compared it to  
360 two other reference systems. One of them were eight randomly selected 12mers, the other was  
361 artificially designed to promote self-elongation. The resulting multimer population in two out of three  
362 of these pools followed a near perfect exponential length profile (Fig. 5d). The random pool resulted  
363 in a similar behavior to the network one but with significantly lower overall concentration of long  
364 chains. Both results are in an excellent agreement with theoretical predictions of reference<sup>39</sup>. A higher  
365 concentration of long chains generated by network 12mers indicates better overall fitness of this set  
366 compared to random 12mers. The “Replicator” set did produce a large number of very long products,  
367 presumably by a different mechanism, but a significantly smaller number of products with short and  
368 medium lengths. This indicates lower autocatalytic ability in both “Replicator” and “Random”  
369 sequence pools when compared to the “Network” pool.

370 For emergence of life on early earth, random oligomers needed to act in an evolution-like behavior.  
371 Here, we followed templated ligation of random 12mer strands made from two bases under  
372 temperature oscillations. Despite its minimalism, the system contains all elements necessary for  
373 Darwinian evolution: out of equilibrium conditions, transmission of sequence information from  
374 template to substrate strands, reliable reproduction of a subset of oligomer products and selection of  
375 fast growing sequences in the process. At the dawn of life, pre-Darwinian dynamics would have been  
376 important to push prebiotic systems towards lower entropy states. Such pre-selection for catalytic  
377 function could have paved the way towards eventual emergence of life.

## 378 METHODS

### 379 Nomenclature

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

386

### 387 Ligation by DNA ligase

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

398

### 399 Design of the random sequence pool

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

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

413

### 414 Temperature Cycling

415 Temperature cyclers *Bio-Rad* T100, *Bio-Rad* CFX96, *Analytik Jena* qTOWER<sup>3</sup> and *Thermo Fisher Scientific*  
416 ProFlex PCR System were used to apply alternating dissociation and ligation temperatures to our  
417 samples. The dissociation temperature of 75  $^{\circ}$ C was chosen, to melt short initially emerging ssDNA of

418 up to 36mer. In the SI we also show how a variation of the dissociation temperature changes multimer  
419 product distribution in a random sequence templated ligation experiment. Lower dissociation  
420 temperatures enable us to run several thousand temperature cycles, as the stability of the TAQ DNA  
421 ligase is reduced substantially for longer times at 95 °C. Time resolution experiments with PAGE-  
422 analysis demonstrated ligase activity even after 2000 temperature cycles for a dissociation  
423 temperature of 75 °C. In experiments screening the ligation temperature (see SI), we found that for  
424 ligation temperatures of 25 °C the product length distribution is exponentially falling. For higher  
425 ligation temperatures such as 33 °C we find more long sequences, but almost no 24mer and 36mer  
426 sequences. For sequenced samples we chose a ligation temperature of 25 °C because the library  
427 preparation kit is better suited for shorter DNA strands. In sequencing data for samples with 33 °C the  
428 yield was very low, but the results are similar to the sequencing data of samples with 25 °C ligation  
429 temperature, but with comparably worse statistics. For dsDNA dissociation in each temperature cycle  
430 the corresponding temperature is held for 20 s with subsequent 120 s at the ligation temperature.  
431

### 432 Sequencing by Next Generation Sequencing (NGS)

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

### 445 Sequence Analysis

446 Demultiplexing was done by a standard demultiplexing algorithm on servers of the Gen Center Munich  
447 running an instance of Galaxy<sup>44</sup> connected to the sequencing machine. *Illumina*-sequencing creates  
448 three FASTA-files, listing the front and the back barcodes and the read sequence, for each lane of the  
449 flow cell. The demultiplexing-algorithm matches the barcodes of the prepared library DNA to the read  
450 sequence and produces a single FASTA file including the read quality scores.

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

457 RegEx for searching AT random sequences:

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

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

460 `^(?!(?:.*?(G|C)){X,})^[ATCG]{12,}`. The sequenced library may have primer-primer dimers and

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

## 467 BACKMATTER

### 468 **Competing Interests**

469 The authors declare that they have no competing interests.

470

### 471 **Author's contribution**

472 P.W.K. performed the experiments, prepared the library for sequencing, performed the  
473 demultiplexing, the analysis, programmed the analysis software, analyzed the data, drafted and wrote  
474 the manuscript. A.V.T and S.M. performed the theoretical analysis and analyzed the data in context of  
475 their already published theoretical work, drafted graphs, drafted and wrote the manuscript. D.B.  
476 contrived the experiment, guided the experimental progress, analyzed data and drafted the  
477 manuscript.

478

### 479 **Acknowledgements**

480 The authors would like to acknowledge funding by the Deutsche Forschungsgemeinschaft (DFG,  
481 German Research Foundation)– Project-ID 201269156 – SFB 1032, the Advanced Grant (EvoTrap  
482 #787356) PE3, ERC-2017-ADG from the European Research Council, CRC 235 Emergence of Life  
483 (Project-ID 364653263) and the Center for NanoScience (CeNS). We would like to thank Ulrich  
484 Gerland, Tobias Göppel, Joachim Rosenberger and Bernhard Altaner for their helpful remarks and  
485 discussions about hybridization energies, baseline corrections and interpretation of multimer  
486 product distributions. P.W.K and D.B. thank Stefan Krebs and Marlis Fischalek at the Gene Center  
487 Munich for help with the library preparation and the sequencing the samples and Annalena Salditt  
488 and Filiz Civril for comments on the manuscript. This research was partially done at, and used  
489 resources of the Center for Functional Nanomaterials, which is a U.S.

490 DOE Office of Science Facility, at Brookhaven National Laboratory under Contract No.~DE-

491 SC0012704.

492

## 493 LITERATURE

- 494 1. Crick, F. H. C. The origin of the genetic code. *J. Mol. Biol.* **38**, 367–379 (1968).
- 495 2. Orgel, L. E. Evolution of the genetic apparatus: A review. *Cold Spring Harb. Symp. Quant. Biol.*  
496 **52**, 9–16 (1987).
- 497 3. Walter, G. The RNA World. *Nature* **319**, 618 (1986).
- 498 4. Attwater, J., Wochner, A., Pinheiro, V. B., Coulson, A. & Holliger, P. Ice as a protocellular medium  
499 for RNA replication. *Nat. Commun.* **1**, 1–8 (2010).
- 500 5. Joyce, G. F. Toward an alternative biology. *Science (80-. )*. **336**, 307–308 (2012).
- 501 6. Horning, D. P. & Joyce, G. F. Amplification of RNA by an RNA polymerase ribozyme. *Proc. Natl.*  
502 *Acad. Sci.* **113**, 9786–9791 (2016).
- 503 7. Hertel, K. J. *et al.* Numbering system for the hammerhead. *Nucleic Acids Res.* **20**, 3252 (1992).
- 504 8. Pley, H. W., Flaherty, K. M. & McKay, D. B. Three-dimensional structure of a hammerhead  
505 ribozyme. *Nature* **372**, 68–74 (1994).
- 506 9. Birikh, K. R., Heaton, P. A. & Eckstein, F. The structure, function and application of the  
507 hammerhead ribozyme. *Eur. J. Biochem.* **245**, 1–16 (1997).
- 508 10. Scott, W. G., Murray, J. B., Arnold, J. R. P., Stoddard, B. L. & Klug, A. Capturing the structure of  
509 a catalytic RNA intermediate: The hammerhead ribozyme. *Science (80-. )*. **274**, 2065–2069  
510 (1996).
- 511 11. Szostak, J. W. & Bartel, D. P. Structurally Complex Highly Active RNA Ligases Derived from  
512 Random RNA Sequences. (1995).
- 513 12. Kozlov, J. A. & Orgel, L. E. Nonenzymatic template-directed synthesis of RNA from monomers.  
514 *Mol. Biol.* **34**, 921–930 (2000).
- 515 13. Oró, J. Mechanism of synthesis of adenine from hydrogen cyanide under possible primitive  
516 earth conditions. *Nature* (1961). doi:10.1038/1911193a0

- 517 14. Lohrmann, R. Formation of nucleoside 5'-polyphosphates from nucleotides and  
518 trimetaphosphate. *J. Mol. Evol.* (1975). doi:10.1007/BF01794633
- 519 15. Handschuh, G. J., Lohrmann, R. & Orgel, L. E. The effect of Mg<sup>2+</sup> and Ca<sup>2+</sup> on urea-catalyzed  
520 phosphorylation reactions. *J. Mol. Evol.* (1973). doi:10.1007/BF01654094
- 521 16. Österberg, R., Orgel, L. E. & Lohrmann, R. Further studies of urea-catalyzed phosphorylation  
522 reactions. *Journal of Molecular Evolution* (1973). doi:10.1007/BF01654004
- 523 17. Tkachenko, A. V. & Maslov, S. Onset of natural selection in populations of autocatalytic  
524 heteropolymers. *J. Chem. Phys.* **149**, (2018).
- 525 18. Fellermann, H., Tanaka, S. & Rasmussen, S. Sequence selection by dynamical symmetry  
526 breaking in an autocatalytic binary polymer model. *Phys. Rev. E* **96**, 1–14 (2017).
- 527 19. Toyabe, S. & Braun, D. Cooperative Ligation Breaks Sequence Symmetry and Stabilizes Early  
528 Molecular Replication. *Phys. Rev. X* **9**, 011056 (2019).
- 529 20. Horowitz, J. M. & England, J. L. Spontaneous fine-tuning to environment in many-species  
530 chemical reaction networks. *Proc. Natl. Acad. Sci. U. S. A.* (2017). doi:10.1073/pnas.1700617114
- 531 21. Hsu, P. P. & Sabatini, D. M. Cancer cell metabolism: Warburg and beyond. *Cell* **134**, 703–707  
532 (2008).
- 533 22. Boyer, P. D. The ATP synthase - a splendid molecular machine. *Annu. Rev. Biochem.* **66**, 717–  
534 749 (1997).
- 535 23. Baross, J. A. & Hoffman, S. E. Submarine hydrothermal vents and associated gradient  
536 environments as sites for the origin and evolution of life. *Orig. Life Evol. Biosph.* (1985).  
537 doi:10.1007/BF01808177
- 538 24. Pascal, R. *et al.* Towards an evolutionary theory of the origin of life based on kinetics and  
539 thermodynamics. 1–9 (2013).
- 540 25. Mutschler, H., Wochner, A. & Holliger, P. Freeze-thaw cycles as drivers of complex ribozyme  
541 assembly. *Nat. Chem.* **7**, 502–508 (2015).
- 542 26. Mast, C. B. & Braun, D. Thermal trap for DNA replication. *Phys. Rev. Lett.* **104**, 1–4 (2010).
- 543 27. Crick, F. H., Watson, J. D. The complementary structure of deoxyribonucleic acid. *Proc. R. Soc.*  
544 *London. Ser. A. Math. Phys. Sci.* **223**, 80–96 (1954).
- 545 28. SantaLucia, J. A unified view of polymer, dumbbell, and oligonucleotide DNA nearest-neighbor  
546 thermodynamics. *Proc. Natl. Acad. Sci. U. S. A.* **95**, 1460–1465 (1998).
- 547 29. Wu, T. & Orgel, L. E. Nonenzymic template-directed synthesis on oligodeoxycytidylate  
548 sequences in hairpin oligonucleotides. **LII**, 317–322 (1992).
- 549 30. Rohatgi, R., Bartel, D. P. & Szostak, J. W. Kinetic and mechanistic analysis of nonenzymatic,  
550 template-directed oligoribonucleotide ligation. *J. Am. Chem. Soc.* **118**, 3332–3339 (1996).
- 551 31. Fahrenbach, A. C. *et al.* Common and Potentially Prebiotic Origin for Precursors of Nucleotide  
552 Synthesis and Activation. *J. Am. Chem. Soc.* **139**, 8780–8783 (2017).
- 553 32. Li, L. *et al.* Enhanced nonenzymatic RNA copying with 2-aminoimidazole activated nucleotides.  
554 *J. Am. Chem. Soc.* **139**, 1810–1813 (2017).
- 555 33. Appel, R., Niemann, B. & Schuhn, W. Synthesis of the First Triphosphabutadiene. *Angew. Chem.*  
556 *Inf. Ed. Engl.* **119**, 932–935 (1986).
- 557 34. Sievers, D. & Von Kiedrowski, G. Self-replication of hexadeoxynucleotide analogues:  
558 Autocatalysis versus cross-catalysis. *Chem. - A Eur. J.* **4**, 629–641 (1998).
- 559 35. Edeleva, E. *et al.* Continuous nonenzymatic cross-replication of DNA strands with in situ  
560 activated DNA oligonucleotides. *Chem. Sci.* **10**, 5807–5814 (2019).
- 561 36. Fischer, S. G. & Lerman, L. S. DNA fragments differing by single base-pair substitutions. *Proc.*  
562 *Nat. Acad. Science, Biochem.* **80**, 1579–1583 (1983).
- 563 37. Myers, R. M., Fischer, S. G., Lerman, L. S. & Maniatis, T. Nearly all single base substitutions in  
564 DNA fragments joined to a GC-clamp can be detected by denaturing gradient gel  
565 electrophoresis. **13**, 3131–3145 (1985).
- 566 38. Derr, J. *et al.* Prebiotically plausible mechanisms increase compositional diversity of nucleic acid  
567 sequences. *Nucleic Acids Res.* **40**, 4711–4722 (2012).
- 568 39. Tkachenko, A. V & Maslov, S. Spontaneous emergence of autocatalytic information-coding

- 569 polymers. *J. Chem. Phys.* **143**, 045102 (2015).  
570 40. Bartel, D. & Szostak, J. Isolation of new ribozymes from a large pool of random sequences [see  
571 comment]. *Science (80-. )*. **261**, 1411–1418 (1993).  
572 41. de Bruijn, N. G. A combinatorial problem. *Proc. Sect. Sci. K. Ned. Akad. van Wet. te Amsterdam*  
573 **49**, 758–764 (1946).  
574 42. Pevzner, P. A., Tang, H. & Waterman, M. S. An Eulerian path approach to DNA fragment  
575 assembly. *Proc. Natl. Acad. Sci. U. S. A.* (2001). doi:10.1073/pnas.171285098  
576 43. Zadeh, J. N. *et al.* NUPACK: Analysis and design of nucleic acid systems. *J. Comput. Chem.* **32**,  
577 170–173 (2011).  
578 44. Afgan, E. *et al.* The Galaxy platform for accessible, reproducible and collaborative biomedical  
579 analyses: 2018 update. *Nucleic Acids Res.* (2018). doi:10.1093/nar/gky379  
580