



# Structured sequences emerge from random pool when replicated by templated ligation

Patrick W. Kudella<sup>a</sup>, Alexei V. Tkachenko<sup>b</sup>, Annalena Salditt<sup>a</sup>, Sergei Maslov<sup>c,d</sup>, and Dieter Braun<sup>a,1</sup>

<sup>a</sup>Systems Biophysics and Center for NanoScience, Ludwigs-Maximilians-Universität München, 80799 Munich, Germany; <sup>b</sup>Center for Functional Nanomaterials, Brookhaven National Laboratory, Upton, NY 11973; <sup>c</sup>Department of Bioengineering, University of Illinois at Urbana–Champaign, Urbana, IL 61801; and <sup>d</sup>Carl R. Woese Institute for Genomic Biology, University of Illinois at Urbana–Champaign, Urbana, IL 61801

Edited by Eugene V. Koonin, National Institutes of Health, Bethesda, MD, and approved January 20, 2021 (received for review September 7, 2020)

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

origin of life | DNA replication | Darwinian evolution | templated ligation | sequence entropy

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

The problem of nonenzymatic formation of single base nucleotides and short oligomers in settings reminiscent of the primordial soup has been studied before (12–17). However, the continuation of this evolutionary path toward early replication networks would require a preselection mechanism of oligonucleotides (see Fig. 1A), lowering the information entropy of the resulting sequence pool (18–22). In principle, such selection modes include optimization for information storage, local oligomer enrichment (e.g., in hydrogels or in catalytically functional sites).

An important aspect of a selection mechanism is its non-equilibrium driving force. Today's highly evolved cells function through multistep and multicomponent metabolic pathways like glycolysis in the Warburg effect (23) or by specialized enzymes

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

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

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

Formation of progressively longer oligomers from shorter ones requires ligation reactions, a method commonly employed in hairpin-mediated RNA and DNA replication (31, 32). At the origin of life, this might have been achieved by activated oligomers (33, 34) or activation agents (35–37), whereas later on the

## Significance

The structure of life emerged from randomness. This is attributed to selection by molecular Darwinian evolution. This study found that random templated ligation led to the simultaneous elongation and sequence selection of oligomers. Product strands showed highly structured sequence motifs which inhibited self-folding and built self-templating reaction networks. By the reduction of the sequence space, the kinetics of duplex formation increased and led to a faster replication through the ligation process. These findings imply that elementary binding properties of nucleotides can lead to an early selection of sequences even before the onset of Darwinian evolution. This suggests that such a simplification of sequence space could result in faster downstream selection for sequence-based function for the origin of life.

Author contributions: P.W.K., A.V.T., A.S., S.M., and D.B. designed research; P.W.K., A.V.T., A.S., S.M., and D.B. performed research; P.W.K., A.V.T., A.S., S.M., and D.B. contributed new reagents/analytic tools; P.W.K., A.V.T., A.S., S.M., and D.B. analyzed data; and P.W.K., A.V.T., S.M., and D.B. wrote the paper.

The authors declare no competing interest.

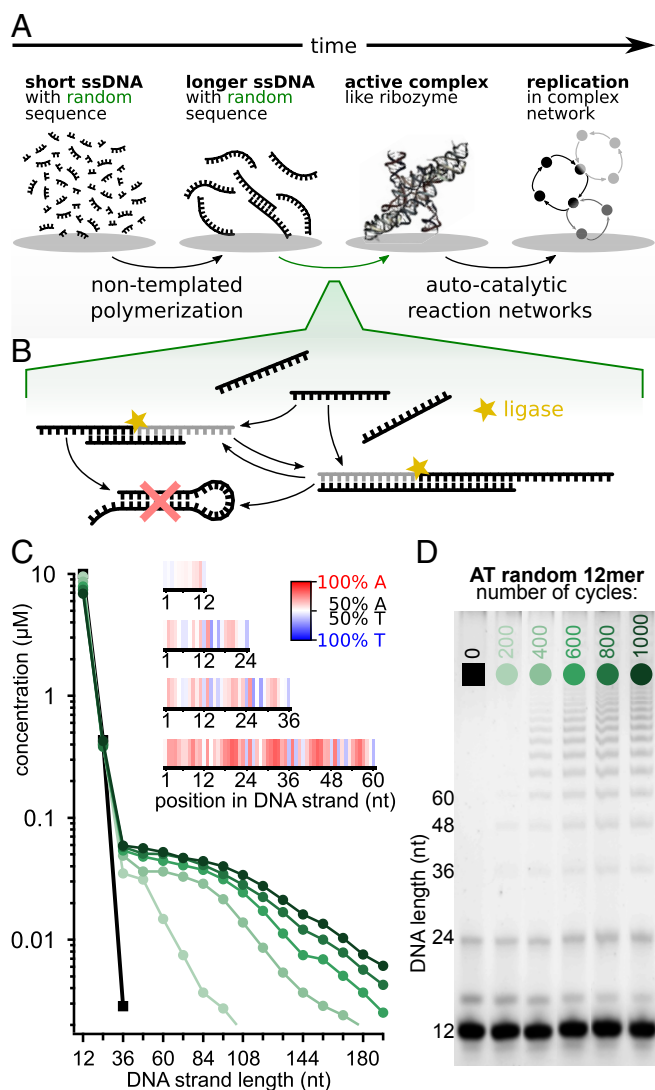
This article is a PNAS Direct Submission.

This open access article is distributed under Creative Commons Attribution-NonCommercial-NoDerivatives License 4.0 (CC BY-NC-ND).

<sup>1</sup>To whom correspondence may be addressed. Email: dieter.braun@lmu.de.

This article contains supporting information online at <https://www.pnas.org/lookup/suppl/doi:10.1073/pnas.2018830118/-DCSupplemental>.

Published February 15, 2021.



**Fig. 1.** Templated ligation of random sequence DNA 12-mers. (A) Before cells evolved, the first ribozymes were thought to perform basic cell functions. In the exponentially vast sequence space, spontaneous emergence of a functional ribozyme is highly unlikely, therefore preselection mechanisms were likely necessary. (B) In our experiment, DNA strands hybridize at low temperatures to form three-dimensional complexes that can be ligated and preserved in the high temperature dissociation steps. The system self-selects for sequences with specific ligation site motifs as well as for strands that continue acting as templates. Hairpin sequences are therefore suppressed. (C) Concentration analysis shows progressively longer strands emerging after multiple temperature cycles. The inset (A-red, T-blue) shows that, although 12-mers (88,009 strands) have essentially random sequences (white), various sequence patterns emerge in longer strands (60-mers, 235,913 strands analyzed). (D) Samples subjected to different number (0 to 1,000) of temperature cycles between 75 °C and 33 °C. Concentration quantification is done on PAGE with SYBR poststained DNA.

formation of simple ribozyme ligases seemed possible (38). Our study is focused on inherent properties of self-assembly by base pairing in random pools of oligomers and not on chemical mechanisms of ligation. Hence, we decided to use TAQ DNA ligase—an evolved enzyme for templated ligation of DNA (21) that is known for its ligation site sequence specificity (39, 40) and lack of sequence-dependent ligation rate (compare *SI Appendix, section 21*). This allowed for fast turnovers of ligation and enabled the observation of sequence dynamics.

## Results

To test templated elongation of polymers in pools of random sequence oligomers, we prepared a 10 μM solution of 12-mer DNA strands composed with nucleobases A and T (sequence space: 4,096) and subjected it to temperature cycling, similar to ref. 21 with 20 s at denaturation temperature of 75 °C and 120 s at ligation temperature of 33 °C. Temperatures were selected according to the melting dynamics of the DNA pool; the time steps were prolonged relative to Toyabe and Braun (21) (*SI Appendix, section 5.3*) because of a greater sequence space. The larger sequence space of full random 12-mers with all four bases did not show any ligation under the same experimental conditions (*SI Appendix, section 5.2*). The sample was split into multiple tubes and exposed to 200, 400, 600, 800, and 1,000 temperature cycles, and one tube was kept at 4 °C for reference, all without influx or outflux of strands. Fig. 1D suggests the maximum depletion of the original 12-mer pool after 1,000 temperature cycles was only about 31%.

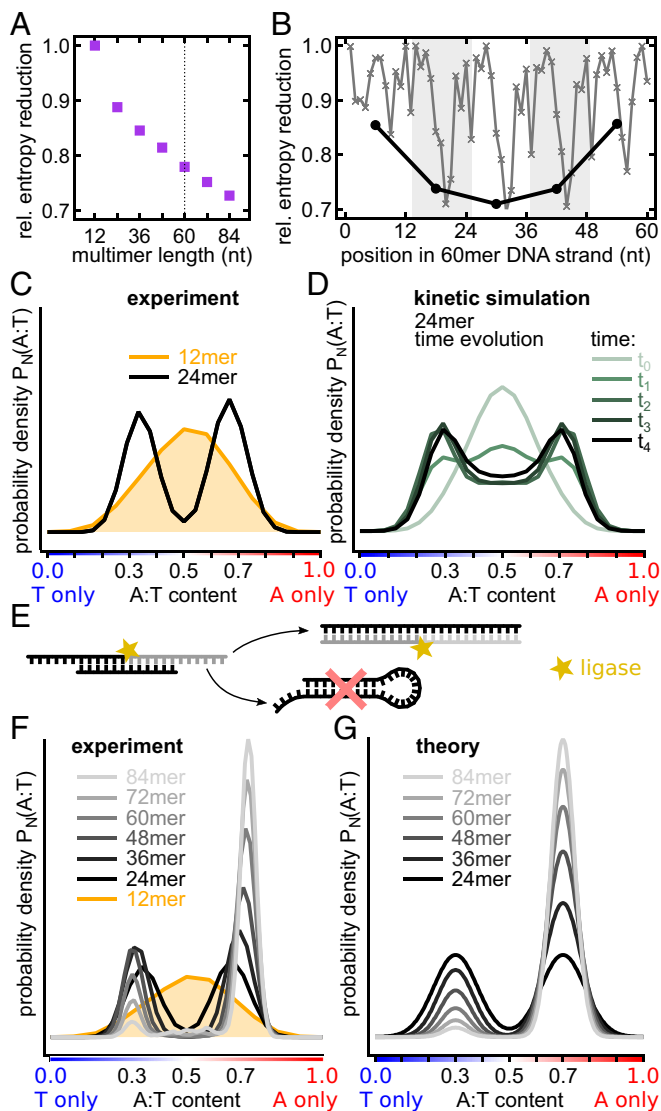
To study the length distributions in our samples, we used polyacrylamide gel electrophoresis (PAGE, Fig. 1D). The first lane is the reference sequence not exposed to temperature cycling, where small amounts of impurities are visible at short lengths (*SI Appendix, section 3.1*). The latter lanes show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

An important property of the initial monomer-pool is its sequence content. Although, for pools with lower sequence complexity, it is possible to show different strand compositions using PAGE (41, 42); a large size of our “monomer” ( $2^{12} = 4,096$ ) and 24-mer product pools (sequence space:  $2^{24} \sim 16.8 \times 10^6$ ) excludes this approach. Thus, we analyzed our final products by next-generation sequencing (NGS) to get insights into product strand compositions.

Plotting the probability of finding a base at a certain position (Fig. 1C, *Inset*) revealed no distinct pattern in 12-mers other than a slight bias toward As. However, longer chains starting with 24-mers developed a strikingly inhomogeneous sequence pattern: bases around ligation sites show a distinct AT-alternating pattern, whereas regions in the middle of individual 12-mers are preferentially enriched with As.

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

In the initial pool of random 12-mers, the A-to-T ratio distribution is shaped binomially, as expected for a random distribution.



**Fig. 2.** Hairpin formation amplifies selection into A-rich and T-rich sequences. (A) Relative entropy reduction as a function of multimer product length: 1 is a random pool, and 0 is a unique sequence. (B) Relative entropy reduction of 60-mer products. Black: Entropy reduction of 12 nt subsequences compared with a random sequence strand of the same length. Gray: Entropy reduction at each nucleotide position showing positional dependence. (C) A gradual development of the bimodal distribution of A:T ratio in chains of different lengths. Whereas the A:T ratio in 12-mers has a single-peaked nearly binomial distribution, 24-mers already have a clearly bimodal distribution peaked at 65:35% (A-type strands) and 35:65% (T-type strands) A:T ratios. (D) Emergence of a bimodal distribution in a kinetic model of templated ligation. (E) Sequences with nearly balanced A:T ratios are prone to the formation of hairpins. In the model in (D) and the experiment, these hairpins prevent strands from acting as templates and substrates for ligation reactions, thereby suppressing the central part of the distribution. (F) A:T ratio distributions in strands of different lengths. As length increases, A-type strands become progressively more abundant in comparison to T-type strands. (G) A:T ratio distributions in a phenomenological model taking into account a slight AT-bias in the initial 12-mer pool resemble experimentally measured ones (E).

However, it dramatically shifted for 24-mer products of ligation: A bimodal distribution of about 65:35% A:T (A-type) as well as the inverse, 35:65% A:T (T-type), was observed with 24-mer products (Fig. 2C). DNA strands composed of only two complementary bases are more prone to formation of single-strand secondary

structures like hairpins than DNAs composed of all four bases. In our templated ligation reaction, we expected that hairpin sequences are not elongated and also not used as template strands because they form catalytically passive Watson–Crick base-paired configuration. A bimodal AT-ratio distribution (Fig. 2D) also emerged in a kinetic computational model in which a pool of random 12-mers was seeded with a small initial amount of random sequence 24-mers. 24-mers that formed hairpins could not act as templates and were therefore less likely to be reproduced (see *SI Appendix*, section 18.2 for details of this model).

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

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

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

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

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

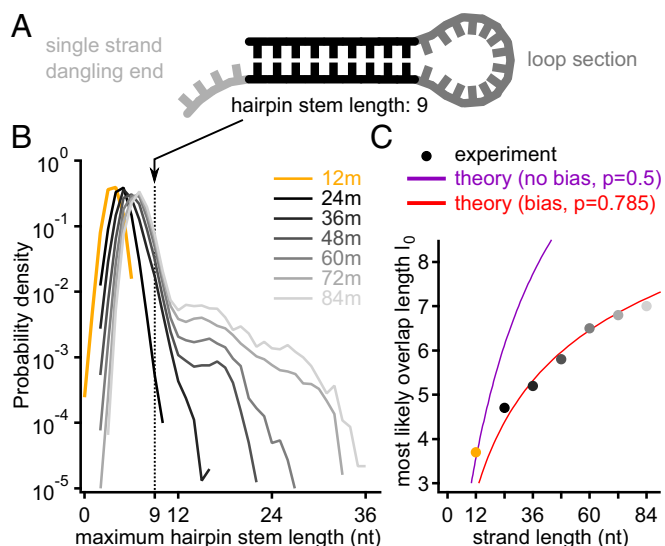
A simple phenomenological model can successfully capture the major features of the observed A:T PDFs for multiple chain lengths. Specifically, we assume both A-type and T-type subpopulations to maximize the sequence entropy, subject to the constraint that the average A:T content is shifted from the midpoint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

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

$$N = 2l_0 + \sqrt{2(2p(1-p))} \frac{l_0}{2}. \quad [3]$$

As one can see in Fig. 3C, this result is in an excellent agreement with experimental data for all the long chains, assuming  $p = 0.785$ . This A:T ratio is indeed comparable to the one observed in the A-type subpopulation. On the other hand, the maximum





**Fig. 3.** Large-scale entropy reduction and sequence correlation per strand. (A) Sketch of a single-strand DNA secondary structure folding on itself, called hairpin. The double stranded part is very similar to a standard duplex DNA. (B) Comparing the PDFs of the maximum hairpin stem length for all strands reveals a group of peaks at around 4 to 7 nt, increasing with the DNA length. Starting with 48-mers, there is a tail visible. These self-similar strands are more abundant the longer the product grows (compare A:T fraction close to  $P = 0.5$  in Fig. 2C). (C) The peak positions as function of the product length follow Eq. 3. The unbiased 12-mers are on the curve with coefficient  $P = 0.5$ , whereas the products starting from 36mers lay on the curve with  $P = 0.785$ . The bias parameter  $p$  is derived from the PDFs in Fig. 2D and describes the A:T-ratio in the strand.

probability length of the longest hairpin for 12-mers is consistent with an unbiased composition,  $p = 0.5$ .

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

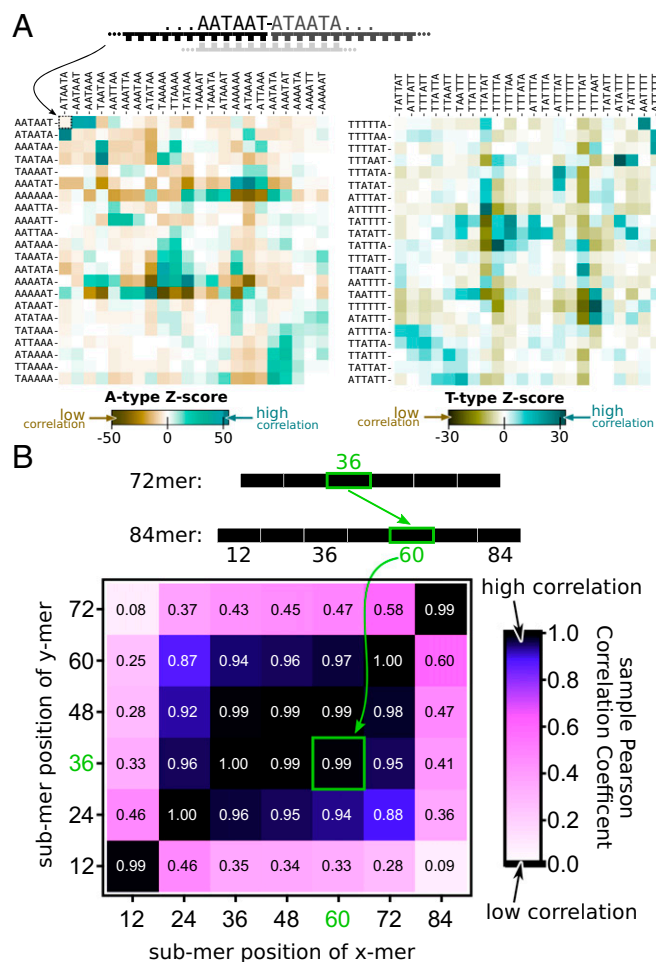
The mechanism of formation of these self-binding sequences may involve recombination of shorter A-type and T-type chains or self-elongation of shorter hairpins. In either case, the hairpin sequence cannot efficiently reproduce by means of template ligation. However, the remainder of the pool would keep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

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

Sequence analysis of the junctions in between the original 12-mer revealed additional information about that landscape, already hinted by patterns seen in Fig. 1B. We characterize pairs of junction-forming sequences with their Z-scores (i.e., probability of their occurrence scaled with its expected value and divided by

the SD calculated in the random binding model) (SI Appendix, section 14).

Fig. 4A shows Z-score heatmaps for junctions within A-type (left panel) and T-type (right panel) subpopulations. More specifically, we show sequences left (row) and right (column) of the junction between the fourth to the fifth 12-mers in the respective 72-mer. These heatmaps reveal a complex landscape of over- and underrepresented junction motifs shown, respectively, in dark teal and dark ochre colors. Emergence of such complex landscape has been theoretically predicted in ref 19. Landscape peaks include repeating A-T motif of alternating bases crossing the ligation site (dark teal peak near the center of each of both



**Fig. 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 with finding it by pure chance and normalized by the SD. Z-scores were calculated for the junction between fourth and the fifth 12-mers in 72-mers of A-type (Left) and T-type (Right), respectively. Other internal junctions in all long chains form very similar landscapes comprised of over- (teal) and underrepresented (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 (note a similarity of landscapes in two panels after this transformation). (B) The matrix of sample Pearson correlation coefficients between abundances of 12-mers in different positions (1–6) inside 72-mers (rows) and 84-mers (columns). Light regions mark low correlations, and 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 12-mer and the last 12-mer of a chain, yet their sequences are similar in chains of different lengths.

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

How similar are selective pressures operating on sequences of different 12-mers within longer chains? Fig. 4B quantifies this similarity in terms of a sample Pearson correlation coefficient (sPCC) between abundances of 12-mer sequences in different positions of long chains of different lengths.

We compare the abundances of  $2^{12} = 4,096$  possible 12-mer sequences in positions 1 through 6 within all 72-mers and compare them with each other and abundances of 12-mers in positions 1 through 7 in all 84-mers. Similar results were obtained for other chains longer than 36 nt. A rectangle of very high correlations ( $>0.9$ ) at the center of the table in Fig. 4B means that very similar sequences get selected at all internal positions of all chains (note that only chains longer than 36 nt have such internally positioned 12-mers). However, the light border of the table means that a rather different subset of 12-mers gets selected in the first and the last position of a multimer. Whatever the nature of selection pressure acting on these 12-mers, it is consistent across oligomers of different lengths as manifested by the high correlation in the lower left and the upper right corner of the table in Fig. 4B.

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

Comparison of abundances of internal 12-mers in A-type and T-type subpopulations predictably yielded no positive correlation and in fact resulted in a slight negative correlation (*SI Appendix, section 11*). However, abundances of reverse complements of sequences from the T-type subpopulation are strongly correlated with those of the A-type, resulting in an sPCC matrix similar to that shown in Fig. 4B (*SI Appendix, Fig. S13*). Therefore, chains in two groups (A-type and T-type) show a considerable degree of reverse complementarity to each other. This fits the elongation and replication mechanism by templated ligation.

To further explore selection capabilities of templated ligation as a function of 12-mer sequences in the initial pool, we conducted three additional experiments referred to as “Replicator,” “Random,” and “Network.” The “Random” experiment started with eight randomly chosen 12-nt sequences served as a control. In the “Replicator” experiment, the pool consisted of eight 12-nt sequences artificially designed for efficient elongation (see below). In the “Network”

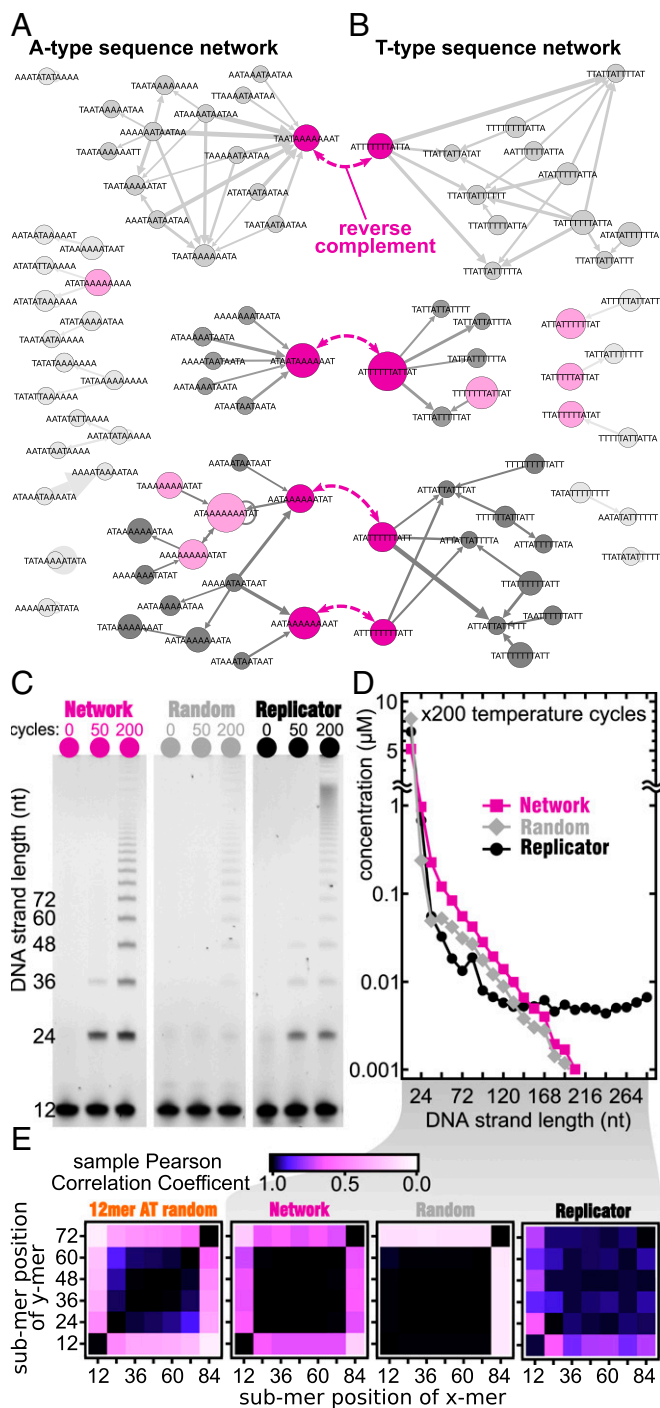
experiment, we populated the pool with eight naturally selected 12-nt sequences commonly found as subsequences of long strands in our original ligation experiment with 4,096 12-mers. To identify these 12-mers, we built a network of the most common 12-mers found in A-type oligomers with a length of more than 48 nt. This network does not include the first and the last 12-mers in a multimer because those are known to be statistically different from the internal ones (Fig. 4B). The circles in Fig. 5A represent unique 12-nt subsequences, and their size describes their Z-scores quantifying their abundance in long chains. The width of the connecting line describes the probability that two subsequences are found one after another in a multimer. The same is done for T-type sequences (Fig. 5B). This representation of a polymer is known as a de Bruijn graph (45) and has been commonly used in DNA fragment analysis and genome assembly (46) and more recently in the context of templated ligation (19).

De Bruijn networks in Fig. 5A break up into several clusters connecting 12-mers with similar subsequences at junctions (TAA-TAA in the top cluster marked by a dark magenta node, ATA-ATA in the middle one, and AAT-AAT in the bottom one). Note that these three common junction subsequences are all related via template shifts. The most common subgraphs found in the A-type network and mirrored among their reverse complements in the T-type network. This pattern is consistent with selection driven by templated ligation (*SI Appendix, section 19*). Among the eight most common subsequences in the A-type network (light and dark magenta nodes in Fig. 5A), four (dark magenta nodes) had a reverse complement among the eight most common subsequences of the T-type network (light and dark magenta nodes in Fig. 5B). These sequences were chosen as the pool of eight 12-mers in the “Network” sample. The “Random” sample consisted of eight 12-mers that were randomly chosen from the 4,096 possible AT-only 12-mers. The “Replicator” sample consisted of eight strands that were built to form three-strand complexes that resemble the assumed first ligation reaction in the pool (*SI Appendix, section 17.1*).

The length distribution of oligomers (Fig. 5D) with concentrations quantified from the PAGE gel image (Fig. 5C) shows that the “Network” sample produced the most product, because the remaining 12-mer sequence concentration was reduced below two other samples down to almost  $5 \mu\text{M}$ . The length distribution in both “Random” and “Network” samples is well described by a piecewise-linear distribution predicted in ref. 18. For short product lengths ranging between 48-mers and 136-mers, the “Random” sample produced more oligomers than the “Replicator” sample. However, for even longer strands, the “Replicator” sample generated the largest number of really long strands because its length distribution reached a plateau around 120-mers. This is probably because of the nature of the eight-sequence pools used here with the “Replicator” one made to form well-aligned double-stranded DNA (dsDNA) that can be properly ligated. According to NUPACK (47), 12-mers in the “Random” sample should not form any complexes that could be subsequently ligated by the TAQ ligase. However, our results shown in Fig. 5C prove the existence of extensive ligation even in the “Random” sample. Presumably, it was initially triggered by a small concentration of complexes formed with low probability, which were subsequently amplified because of the exponential growth of longer strands in our experiment, just like in the “Network” sample.

## Discussion

We experimentally studied templated ligation in a pool of 12-mers made of A and T bases with all possible sequences ( $2^{12} = 4,096$ ), subjected to multiple temperature cycles. To accelerate hypothetical spontaneous ligation reactions operating in the prebiotic world, we employed TAQ DNA ligase in our



**Fig. 5.** Testing self-selection with custom sequence pools. (A) The de Bruijn graph of overrepresented sequence motifs between consecutive 12-mers 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 12-mers are connected by purple dashed arrows. In each network, three families with distinctly similar patterns are observed that each include at least one of the complementary strands. Node sizes reflect relative abundance of 12-mers, and edge thickness denotes the Z-score of the junction between nodes it connects. Light and dark magenta-colored nodes are the eight most abundant 12-mers in each of the two networks. (C) PAGE images of templated ligation of three different samples of 12-mers after different number of temperature cycles (columns): “Replicator”: four substrate 12-mers and four template 12-mers artificially designed for templated ligation, as explained in

experiments. This process produced a complex and heterogeneous ensemble of oligomer products. By performing the NGS of these oligomers, we found that long strands in this ensemble have a significantly lower information entropy compared with a random set of oligomers of the same length. This effect became increasingly more pronounced for longer oligomers (Fig. 2E). The overall reduction in entropy was in line with the theoretical prediction obtained within a simplified model of template-based ligation (19). In that model, the reduction of entropy was due to “mass extinction” in sequence space, with only a very limited (though still exponentially large) set of survivor sequences emerging. In the present experiment, related variation in abundances of different sequences did develop but did not proceed all the way to extinction.

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

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

We selected a subset of eight pairs of mutually complementary 12-mers that appeared anomalously often within longer chains

**SI Appendix; “Random”:** eight random sequence 12-mers randomly selected from the 4,096 possible AT-only 12-mers; **“Network”:** the four most common 12-mers from A-type and another four of T-type shown in dark magenta in A. (D) After 200 temperature cycles, the “Replicator” shows a consistently higher product concentration for all lengths followed by the “Network” sample and then by the “Random” subsamples. In the “Network” and “Random” samples, the length distribution above 48 nt is well described by an exponential distribution as predicted in ref. 18. (E) Pearson correlation matrices between 12-mer abundances within 72-mers and 84-mers in each sample (same as in Fig. 4B). Although the pattern of correlations in the “Network” sample (second from left) resembles that shown in Fig. 4B (reproduced in the left-most subpanel), the “Random” sample (second from right) singles out the last 12-mer but not the first one. The “Replicator” sample (the rightmost subpanel) has its own distinct self-similar pattern of correlations.



and were well connected within the de Bruijn graph. Using this “Network” subset as a new starting pool, we repeated the temperature-cycling experiment and compared it with two other reference systems. One of them were eight randomly selected 12-mers, the other was artificially designed to promote self-elongation. The resulting multimer population in two out of three of these pools followed a near perfect exponential length profile (Fig. 5D). The random pool resulted in a similar behavior to the network one but with significantly lower overall concentration of long chains. Both results are in an excellent agreement with theoretical predictions of ref. 18. A higher concentration of long chains generated by network 12-mers indicates better overall fitness of this set compared with random 12-mers. The “Replicator” set did produce a large number of very long products, presumably by a different mechanism, but a significantly smaller number of products with short and medium lengths. This indicates a lower ability for cross-ligation in both “Replicator” and “Random” sequence pools when compared with the “Network” pool. In *SI Appendix, section 20*, the de Bruijn sequence networks for oligomer products show this difference in elongation fitness clearly: Whereas the “Network” sample forms A-type and T-type groups and is well interconnected, the “Replicator” favors only two sequences.

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

Despite its minimalism, the studied system contains all elements necessary for Darwinian evolution: out of equilibrium conditions, transmission of sequence information from template to substrate strains, reliable reproduction of a subset of oligomer products and the possibility to select from the long fast-growing sequences in the process. At the dawn of life, such pre-Darwinian dynamics would have pushed prebiotic systems toward lower entropy states. A subsequent selection for catalytic function from the replicated structured sequences could then have paved the way toward the eventual emergence of life.

## Materials and Methods

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

**Ligation by DNA Ligase.** For enzymatic ligation of single-stranded DNA (ssDNA), a TAQ DNA ligase from New England Biolabs was used. Chemical reaction conditions were as stated by the manufacturer: 10  $\mu$ M total DNA

concentration in 1 $\times$  ligase buffer. The ligase has a temperature dependent activity and is not active at low (4 to 10  $^{\circ}$ C) and very high temperatures (85 to 95  $^{\circ}$ C). In our experimental system, DNA hybridization characteristics are strongly temperature dependent, as shown in the *SI Appendix*. We expect this to have a stronger influence on the overall length distribution and product concentrations than ligase activity, because the timescale of hybridization is significantly longer than the timescale of ligation (compared in *SI Appendix*). The manufacturer provides activity of the ligase in units/mL, specifically, “One unit is defined as the amount of enzyme required to give 50% ligation of the 12-bp cohesive ends of 1  $\mu$ g of BstEII-digested  $\lambda$  DNA in a total reaction volume of 50  $\mu$ L in 15 min at 45  $^{\circ}$ C.”

**Design of the Random Sequence Pool.** The use of a DNA ligase enables very fast ligation with low error rate. However, not every DNA system is suitable for templated ligation. As stated by the manufacturer, the TAQ ligase does not ligate overhangs which are 4 nt or shorter. Therefore, the shortest possible length of strands is 10-mer, opening up  $4^{10} > 10^6$  different monomer sequences. The resulting pool cannot be sequenced to a reasonable extend. We artificially reduced the sequence space by limiting sequences to only include bases A and T. 10-mer strands with random AT sequence have too low a melting temperature, in a range in which the ligase is not active (compare *SI Appendix*). We found 12-mers with random AT sequences to successfully ligate and to produce longer product strands because of their elevated melting temperature. The monomer sequence space is  $2^{12} = 4,096$ , which is not too large, so we were able to completely sequence it multiple times.

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

**Temperature Cycling.** Temperature cyclers Bio-Rad T100, Bio-Rad CFX96, Analytik Jena qTOWER<sup>3</sup>, and Thermo Fisher Scientific ProFlex PCR System were used to apply alternating dissociation and ligation temperatures to our samples. The dissociation temperature of 75  $^{\circ}$ C was chosen to melt short initially emerging ssDNA of up to 36-mer. In the *SI Appendix*, we also show how a variation of the dissociation temperature changes multimer product distribution in a random sequence templated ligation experiment. Lower dissociation temperatures enable us to run several thousand temperature cycles as the stability of the TAQ DNA ligase is reduced substantially for longer times at 95  $^{\circ}$ C. Time resolution experiments with PAGE analysis demonstrated ligase activity even after 2,000 temperature cycles for a dissociation temperature of 75  $^{\circ}$ C. In experiments screening the ligation temperature (*SI Appendix*), we found that, for ligation temperatures of 25  $^{\circ}$ C, the product length distribution was exponentially falling. For higher ligation temperatures such as 33  $^{\circ}$ C, we found more long sequences but almost no 24-mer and 36-mer sequences. For sequenced samples, we chose a ligation temperature of 25  $^{\circ}$ C because the library preparation kit is better suited for shorter DNA strands. In sequencing data for samples with 33  $^{\circ}$ C, the yield was very low, but the results are similar to the sequencing data of samples with 25  $^{\circ}$ C ligation temperature, but with comparably worse statistics. For dsDNA dissociation in each temperature cycle, the corresponding temperature is held for 20 s with a subsequent 120 s at the ligation temperature.

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

**Sequence Analysis.** Demultiplexing was done by a standard demultiplexing algorithm on servers of the Gen Center Munich running an instance of Galaxy (52) connected to the sequencing machine. *Illumina* sequencing creates three FASTA files, listing the front and the back barcodes and the read sequence for each lane of the flow cell. The demultiplexing algorithm matches

the barcodes of the prepared library DNA to the read sequence and produces a single FASTA file that includes the read quality scores.

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

### RegEx for Searching AT Random Sequences. $(\wedge[\text{ATCG}]\{12\}[\text{ATCG}]\{24\}[\text{ATCG}]\{36\}[\text{ATCG}]\{48\}[\text{ATCG}]\{60\}[\text{ATCG}]\{72\}[\text{ATCG}]\{84\})?(?=[\text{CT}]\{4,\}\text{AGAT})$

RegEx for selecting a maximum of X false reads of G or C in random sequence AT samples:  $\wedge(?![?].*[G|C])\{X,\}\wedge([\text{ATCG}]\{12,\})$ . The sequencing library may have primer-primer dimers and oligomers as well as partial primers that were falsely built in the library preparation step. Because the SWIFT kit is made for longer sequences by design, shorter sequences such as 12-mer in our study may have lower yields and larger error rates for the library kit chemistry. Therefore, the inclusion of sequences with single or multiple false reads can improve the statistics, as long as submotifs with obviously faulty reads are ignored in the analysis.

1. F. H. C. Crick, The origin of the genetic code. *J. Mol. Biol.* **38**, 367–379 (1968).
2. L. E. Orgel, Evolution of the genetic apparatus: A review. *Cold Spring Harb. Symp. Quant. Biol.* **52**, 9–16 (1987).
3. G. Walter, The RNA world. *Nature* **319**, 618 (1986).
4. J. Attwater, A. Wochner, V. B. Pinheiro, A. Coulson, P. Holliger, Ice as a protocellular medium for RNA replication. *Nat. Commun.* **1**, 76 (2010).
5. G. F. Joyce, Toward an alternative biology. *Science* **336**, 307–308 (2012).
6. D. P. Horning, G. F. Joyce, Amplification of RNA by an RNA polymerase ribozyme. *Proc. Natl. Acad. Sci. U.S.A.* **113**, 9786–9791 (2016).
7. K. J. Hertel *et al.*, Numbering system for the hammerhead. *Nucleic Acids Res.* **20**, 3252 (1992).
8. H. W. Pley, K. M. Flaherty, D. B. McKay, Three-dimensional structure of a hammerhead ribozyme. *Nature* **372**, 68–74 (1994).
9. K. R. Birikh, P. A. Heaton, F. Eckstein, The structure, function and application of the hammerhead ribozyme. *Eur. J. Biochem.* **245**, 1–16 (1997).
10. W. G. Scott, J. B. Murray, J. R. P. Arnold, B. L. Stoddard, A. Klug, Capturing the structure of a catalytic RNA intermediate: The hammerhead ribozyme. *Science* **274**, 2065–2069 (1996).
11. J. W. Szostak, D. P. Bartel, Structurally complex highly active RNA ligases derived from random RNA sequences. *Science* **269**, 364–370 (1995).
12. I. A. Kozlov, L. E. Orgel, [Nonenzymatic matrix synthesis of RNA from monomers] [in Russian]. *Mol. Biol. (Mosk.)* **34**, 921–930 (2000).
13. J. Oro, Mechanism of synthesis of adenine from hydrogen cyanide under possible primitive earth conditions. *Nature* **191**, 1193–1194 (1961).
14. R. Lohrmann, Formation of nucleoside 5'-polyphosphates from nucleotides and trimetaphosphate. *J. Mol. Evol.* **6**, 237–252 (1975).
15. G. J. Handschuh, R. Lohrmann, L. E. Orgel, The effect of Mg<sup>2+</sup> and Ca<sup>2+</sup> on urea-catalyzed phosphorylation reactions. *J. Mol. Evol.* **2**, 251–262 (1973).
16. R. Österberg, L. E. Orgel, R. Lohrmann, Further studies of urea-catalyzed phosphorylation reactions. *J. Mol. Evol.* **2**, 231–234 (1973).
17. Z. Liu *et al.*, Harnessing chemical energy for the activation and joining of prebiotic building blocks. *Nat. Chem.* **12**, 1023–1028 (2020).
18. A. V. Tkachenko, S. Maslov, Spontaneous emergence of autocatalytic information-coding polymers. *J. Chem. Phys.* **143**, 045102 (2015).
19. A. V. Tkachenko, S. Maslov, Onset of natural selection in populations of autocatalytic heteropolymers. *J. Chem. Phys.* **149**, 134901 (2018).
20. H. Fellermann, S. Tanaka, S. Rasmussen, Sequence selection by dynamical symmetry breaking in an autocatalytic binary polymer model. *Phys. Rev. E* **96**, 062407 (2017).
21. S. Toyabe, D. Braun, Cooperative ligation breaks sequence symmetry and stabilizes early molecular replication. *Phys. Rev. X* **9**, 011056 (2019).
22. J. M. Horowitz, J. L. England, Spontaneous fine-tuning to environment in many-species chemical reaction networks. *Proc. Natl. Acad. Sci. U.S.A.* **114**, 7565–7570 (2017).
23. P. P. Hsu, D. M. Sabatini, Cancer cell metabolism: Warburg and beyond. *Cell* **134**, 703–707 (2008).
24. P. D. Boyer, The ATP synthase—A splendid molecular machine. *Annu. Rev. Biochem.* **66**, 717–749 (1997).
25. J. A. Baross, S. E. Hoffman, Submarine hydrothermal vents and associated gradient environments as sites for the origin and evolution of life. *Orig. Life Evol. Biosph.* **15**, 327–345 (1985).
26. R. Pascal *et al.*, Towards an evolutionary theory of the origin of life based on kinetics and thermodynamics. *Open Biol.* **3**, 130156 (2013).
27. H. Mutschler, A. Wochner, P. Holliger, Freeze-thaw cycles as drivers of complex ribozyme assembly. *Nat. Chem.* **7**, 502–508 (2015).
28. C. B. Mast, D. Braun, Thermal trap for DNA replication. *Phys. Rev. Lett.* **104**, 188102 (2010).

**Data Availability.** All used data was cited or is reproducible from the study. NGS data files are available upon request.

**ACKNOWLEDGMENTS.** We would like to acknowledge funding by the Simons Foundation (327125 to D.B.), the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation, Project-ID 201269156 – SFB 1032), the Advanced Grant (EvoTrap #787356) PE3, ERC-2017-ADG from the European Research Council, CRC 235 Emergence of Life (Project-ID 364653263), the DFG under Germany's Excellence Strategy EXC-2094-390783311 and the Center for NanoScience. We thank Ulrich Gerland, Joachim Rosenberger, Tobias Göppel, and Bernhard Altaner for their helpful remarks and discussions about hybridization energies, baseline corrections, and the interpretation of multimer product distributions. Our collaboration with this group on the in-depth analysis of the elongation dynamics is currently submitted for review. P.W.K. and D.B. thank Irene Chen and Daniel Duzdevich for help with optimizations of the library preparation protocol and analysis, as well as Stefan Krebs and Marlis Fischalek at the Gene Center Munich for their help with the library preparation and the sequencing of the samples. Additionally, we would like to thank Filiz Civril for her extensive comments on the manuscript. This research was partially done at and used resources of the Center for Functional Nanomaterials, which is a US Department of Energy (DOE) Office of Science Facility, at Brookhaven National Laboratory under Contract No. ~DE-SC0012704.

29. J. D. Crick, F. H. Watson, The complementary structure of deoxyribonucleic acid. *Proc. R. Soc. Lond. Ser. A. Math. Phys. Sci.* **223**, 80–96 (1954).
30. J. SantaLucia, Jr, A unified view of polymer, dumbbell, and oligonucleotide DNA nearest-neighbor thermodynamics. *Proc. Natl. Acad. Sci. U.S.A.* **95**, 1460–1465 (1998).
31. T. Wu, L. E. Orgel, Nonenzymatic template-directed synthesis on oligodeoxycytidylate sequences in hairpin oligonucleotides. *J. Am. Chem. Soc.* **114**, 317–322 (1992).
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).
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).
34. L. Li *et al.*, Enhanced nonenzymatic RNA copying with 2-aminoimidazole activated nucleotides. *J. Am. Chem. Soc.* **139**, 1810–1813 (2017).
35. R. Appel, B. Niemann, W. Schuhn, Synthesis of the first triphosphabutadiene. *Angew. Chem. Int. Ed. Engl.* **119**, 932–935 (1986).
36. D. Sievers, G. Von Kiedrowski, Self-replication of hexadeoxynucleotide analogues: Autocatalysis versus cross-catalysis. *Chemistry* **4**, 629–641 (1998).
37. E. Edeleva *et al.*, Continuous nonenzymatic cross-replication of DNA strands with *in situ* activated DNA oligonucleotides. *Chem. Sci. (Camb.)* **10**, 5807–5814 (2019).
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.* **142**, 15961–15965 (2020).
39. J. Kim, M. Mrksich, Profiling the selectivity of DNA ligases in an array format with mass spectrometry. *Nucleic Acids Res.* **38**, e2 (2010).
40. G. J. S. Lohman *et al.*, A high-throughput assay for the comprehensive profiling of DNA ligase fidelity. *Nucleic Acids Res.* **44**, e14 (2016).
41. S. G. Fischer, L. S. Lerman, DNA fragments differing by single base-pair substitutions are separated in denaturing gradient gels: Correspondence with melting theory. *Proc. Nat. Acad. Sci. U.S.A.* **80**, 1579–1583 (1983).
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. *Nucleic Acids Res.* **13**, 3131–3145 (1985).
43. J. Derr *et al.*, Prebiotically plausible mechanisms increase compositional diversity of nucleic acid sequences. *Nucleic Acids Res.* **40**, 4711–4722 (2012).
44. D. Bartel, J. Szostak, Isolation of new ribozymes from a large pool of random sequences [see comment]. *Science* **261**, 1411–1418 (1993).
45. N. G. de Bruijn, A combinatorial problem. *Proc. Sect. Sci. K. Ned. Akad. van Wet. te Amsterdam* **49**, 758–764 (1946).
46. P. A. Pevzner, H. Tang, M. S. Waterman, An Eulerian path approach to DNA fragment assembly. *Proc. Natl. Acad. Sci. U.S.A.* **98**, 9748–9753 (2001).
47. J. N. Zadeh *et al.*, NUPACK: Analysis and design of nucleic acid systems. *J. Comput. Chem.* **32**, 170–173 (2011).
48. M. Eigen, Selforganization of matter and the evolution of biological macromolecules. *Naturwissenschaften* **58**, 465–523 (1971).
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).
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).
51. M. Morasch *et al.*, Heated gas bubbles enrich, crystallize, dry, phosphorylate and encapsulate prebiotic molecules. *Nat. Chem.* **11**, 779–788 (2019).
52. E. Afgan *et al.*, The Galaxy platform for accessible, reproducible and collaborative biomedical analyses: 2018 update. *Nucleic Acids Res.* **46**, W537–W544 (2018).