

Onset of natural selection in populations of autocatalytic heteropolymers

Alexei V. Tkachenko, and Sergei Maslov

Citation: *The Journal of Chemical Physics* **149**, 134901 (2018); doi: 10.1063/1.5048488

View online: <https://doi.org/10.1063/1.5048488>

View Table of Contents: <http://aip.scitation.org/toc/jcp/149/13>

Published by the [American Institute of Physics](#)

Articles you may be interested in

[Encoding and selecting coarse-grain mapping operators with hierarchical graphs](#)

The Journal of Chemical Physics **149**, 134106 (2018); 10.1063/1.5040114

[\$\pi\$ - \$\pi\$ stacking vs. C-H/ \$\pi\$ interaction: Excimer formation and charge resonance stabilization in van der Waals clusters of 9,9'-dimethylfluorene](#)

The Journal of Chemical Physics **149**, 134314 (2018); 10.1063/1.5044648

[Polytetrahedral structure and glass-forming ability of simulated Ni-Zr alloys](#)

The Journal of Chemical Physics **149**, 134501 (2018); 10.1063/1.5041325

[Exponential parameterization of wave functions for quantum dynamics: Time-dependent Hartree in second quantization](#)

The Journal of Chemical Physics **149**, 134110 (2018); 10.1063/1.5049344

[Importance sampling large deviations in nonequilibrium steady states. I](#)

The Journal of Chemical Physics **148**, 124120 (2018); 10.1063/1.5003151

[Non-statistical intermolecular energy transfer from vibrationally excited benzene in a mixed nitrogen-benzene bath](#)

The Journal of Chemical Physics **149**, 134101 (2018); 10.1063/1.5043139

PHYSICS TODAY

WHITEPAPERS

ADVANCED LIGHT CURE ADHESIVES

Take a closer look at what these environmentally friendly adhesive systems can do

READ NOW

PRESENTED BY
 **MASTERBOND**
ADHESIVES | SEALANTS | COATINGS

Onset of natural selection in populations of autocatalytic heteropolymers

Alexei V. Tkachenko^{1,a)} and Sergei Maslov^{2,3,b)}

¹Center for Functional Nanomaterials, Brookhaven National Laboratory, Upton, New York 11973, USA

²Department of Bioengineering, University of Illinois at Urbana-Champaign, 1270 Digital Computer Laboratory, MC-278, Urbana, Illinois 61801, USA

³Carl R. Woese Institute for Genomic Biology, University of Illinois, Urbana-Champaign, Illinois 61801, USA

(Received 15 July 2018; accepted 10 September 2018; published online 4 October 2018)

Reduction of information entropy along with ever-increasing complexity is among the key signatures of life. Understanding the onset of such behavior in the early prebiotic world is essential for solving the problem of the origin of life. Here we study a general problem of heteropolymers capable of template-assisted ligation based on Watson-Crick-like hybridization. The system is driven off-equilibrium by cyclic changes in the environment. We model the dynamics of 2-mers, i.e., sequential pairs of specific monomers within the heteropolymer population. While the possible number of them is Z^2 (where Z is the number of monomer types), we observe that most of the 2-mers get extinct, leaving no more than $2Z$ survivors. This leads to a dramatic reduction of the information entropy in the sequence space. Our numerical results are supported by a general mathematical analysis of the competition of growing polymers for constituent monomers. This natural-selection-like process ultimately results in a limited subset of polymer sequences. Importantly, the set of surviving sequences depends on initial concentrations of monomers and remains exponentially large (2^L down from Z^L for length L) in each of realizations. Thus, an inhomogeneity in initial conditions allows for a massively parallel search of the sequence space for biologically functional polymers, such as ribozymes. We also propose potential experimental implementations of our model in the contexts of either biopolymers or artificial nano-structures. *Published by AIP Publishing.* <https://doi.org/10.1063/1.5048488>

I. INTRODUCTION

The second law of thermodynamics states that the entropy of a closed system increases with time. Life represents a remarkable example of the opposite trend taking place in an open, non-equilibrium system.¹ Indeed, both information and thermodynamic entropies decrease in the course of Darwinian evolution, reflecting ever-increasing complexity of living organisms and their communities.² Interestingly, both the second law and the concept of entropy were introduced by Clausius in the 1850s, roughly at the same time as Darwin developed and published his seminal work. A century later the connection between life and entropy was highlighted in the classical work of Schrödinger titled “What is Life?”³ According to him, living systems are characterized by their ability to “feed on” and store the negative entropy (which he referred to as “negentropy”).² In the same work, Schrödinger effectively predicted the existence of information-storing molecules such as DNA. Soon after, Brillouin established⁴ the connection between the thermodynamic negentropy and its information cousin defined by Shannon.⁵

The emergence of life from non-living matter is one of the greatest mysteries of fundamental science. In addition, the search for artificial self-replicating nano- and micro-scale

systems is an exciting field with potential engineering applications.^{6–9} The central challenge in both of these fields is to come up with a simple, physically realizable self-replicating system obeying the laws of thermodynamics, yet ultimately capable of Darwinian evolution.

Chemical networks of molecules engaged in mutual catalysis have long been considered a plausible form of the prebiotic world.^{10–13} Furthermore, a set of mutually catalyzing RNA-based enzymes (ribozymes) is one of the best known examples of experimentally realized autonomous self-replication. This is viewed as major evidence supporting the RNA-world hypothesis (see, e.g., Refs. 14–18). The ribozyme activity requires relatively long polymers made of hundreds of nucleotides with carefully designed sequences, whose spontaneous emergence by pure chance is nearly impossible. Thus, to make the first steps toward explanation of the origin of life, one needs to come up with a much simpler system capable of spontaneous reduction of the information entropy, which would ultimately set the stage for Darwinian evolution, e.g., toward functional ribozymes and/or autocatalytic metabolic cycles.

A promising candidate for such a mechanism is provided by template-assisted ligation. In this process, pairs of polymers are brought together via hybridization with a complementary template chain and eventually ligated to form a longer chain [see Figs. 1(a) and 1(b)]. Unlike the non-templated reversible step-growth polymerization used in Ref. 19, this mechanism naturally involves transmission of sequence information from

^{a)}Electronic mail: oleksiyt@bnl.gov

^{b)}Electronic mail: ssmaslov@gmail.com

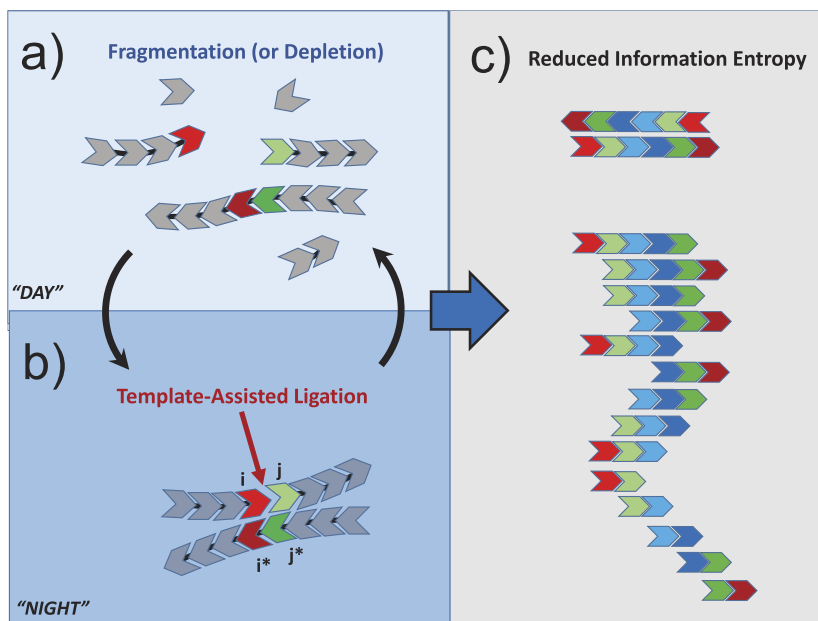


FIG. 1. A conceptual illustration of our model. The population of heteropolymers is cycled between day (a) and night (b) phases. During the night phase, polymer chains undergo template-assisted ligation, joining left and right ends i and j , respectively, to form a new 2-mer ij . This process is assisted by a complementary 2-mer j^*i^* . The process results in a reduced information entropy of chain sequences, which ends up being dominated by a small subset of 2-mers as illustrated in (c).

the template to the newly ligated chain, thus opening an exciting possibility of long-term memory and evolvability. An early conceptual model involving template-assisted polymerization was proposed by Anderson and colleagues.^{20,21} It has also been a subject of several more recent experimental and theoretical studies.^{22,24,25} In particular, the model by Hordijk *et al.*²² makes a connection between the classical Kauffman model of autocatalytic sets¹² and polymer systems capable of template-assisted ligation (see Ref. 23 for further development of that approach). Recently, we theoretically established²⁶ that a cyclically driven system of this type is capable of producing long, mutually catalyzing chains starting from a primordial soup dominated by monomers. A conceptually similar model combining templated and non-templated ligation has been recently used²⁷ to describe the chiral symmetry breaking in a mixture of autocatalytic polymers. In the current study, we focus on the statistics of sequences of these chains and discover that the dynamics of the system naturally results in a dramatic reduction of the information entropy in the sequence space.

II. RESULTS

A. Model

Here we further develop the model introduced in Ref. 26. It describes the emergence of heteropolymers out of the “primordial soup” of monomers by virtue of template-assisted ligation. Our system is driven out of equilibrium by cyclic changes in physical conditions such as temperature, salt concentration, pH, etc. (see Fig. 1).

We consider a general case of information-coding heteropolymers composed of Z types of monomers capable of making $Z/2$ mutually complementary pairs. Polymerization occurs during the “night” phase of each cycle when existing heteropolymers may serve as templates for ligation of pairs of chains to form longer ones. When the end groups

of two substrate chains are positioned next to each other by virtue of hybridization with the template, a new covalent bond connecting these end groups is formed at a certain rate [see Fig. 1(b)]. During the “day” phase of each cycle, all hybridized pairs dissociate and individual chains are fully dispersed [see Fig. 1(a)].

One of the key results of our previous work²⁶ is the existence of the optimal hybridization overlap length k_0 for template-substrate binding. In this work, for the sake of simplicity, we assume that a single pair of complementary monomers is sufficient to bind a substrate to a template. This can be interpreted as if each of Z monomers in the present model is in fact a “word” composed of k_0 smaller elementary letters, e.g., RNA or DNA bases. Within this interpretation, the number Z of such “composite monomers” can be exponentially large: $Z = z^{k_0}$, where z is the number of elementary letters ($z = 4$ in the case of RNA). As in Ref. 26, we ignore the process of spontaneous, non-templated ligation.^{19,25}

In our model, monomer types are labeled in such a way that type i is complementary to type i^* . One of the key concepts in our analysis is that of a “2-mer” ij referring to a monomer i immediately followed by the monomer j and found anywhere within any heteropolymer. Note that, similar to DNA/RNA complementary strands, polymers in our system are assumed to be directional and anti-parallel when hybridized to each other. Therefore, a 2-mer j^*i^* formed from monomers j^* and i^* is complementary to the 2-mer ij . It can serve as a template catalyzing the ligation of two substrate chains with monomers i and j located at their appropriate ends [see Fig. 1(b)].

Let d_{ij} denote the overall concentration of 2-mers of type ij , i.e., the total number of consecutive monomers of types i and j found anywhere within any chain, divided by the volume of the system. We will refer to the $Z \times Z$ matrix formed by all d_{ij} as the 2-mer matrix. Let r_i denote the concentration of all chains ending with a monomer of type i at their right end, while

l_j is the concentration of all chains starting with a monomer of type j at their left end. When two ends i and j of such chains meet due to hybridization with a complementary template j^*i^* , they are ligated at a certain rate to form a new 2-mer ij . We describe this process by a three-body mass-action term $\lambda_{ij} \cdot r_i(t) \cdot l_j(t) \cdot d_{j^*i^*}(t)$. Here λ_{ij} is the ligation rate averaged over the duration of the day-night cycle with the understanding that ligation happens only during the night phase. 2-mers ij in our system are assumed to spontaneously break up at a rate β_{ij} . Thus here we extend our original model by introducing an explicit sequence dependence of ligation (λ_{ij}) and breakage (β_{ij}) rates. Master equations, describing the slow dynamics in our system occurring over multiple day/night cycles, are

$$\dot{d}_{ij}(t) = \lambda_{ij} r_i(t) l_j(t) d_{j^*i^*}(t) - \beta_{ij} d_{ij}(t). \quad (1)$$

This mass-action description implies that our system stays well below the saturation regime during the night phase. In other words, we assume that template-substrate hybridization probability is determined by the association rate and not by the competition of multiple different substrates for the same binding site on a chain. This is realized when the duration of the night phase of the cycle is shorter than the typical association time for hybridization. Importantly, this regime also ensures that there is no template poisoning; i.e., the probability of two complementary 2-mers binding each other (and thus loosing their catalytic activity) remains low.

One could write a similar set of kinetic equations describing the dynamics of concentrations of “left” and “right” ends of chains, $l_i(t)$ and $r_i(t)$. Instead, we use the conservation of overall concentrations of monomers of each type to obtain the explicit algebraic expressions for $l_i(t)$ and $r_i(t)$ in terms of the 2-mer matrix,

$$\begin{aligned} l_i(t) &= c_i - \sum_k d_{ki}(t), \\ r_i(t) &= c_i - \sum_k d_{ik}(t). \end{aligned} \quad (2)$$

Here c_i is the overall concentration of monomers of type i in the pool, both free and bound. At the start, only free monomers are present ($d_{ij} = 0$), and thus, the initial conditions are given by $l_i(0) = r_i(0) = c_i$.

Our model allows for an alternative interpretation that does not involve breakage of intra-polymer bonds. One can show that Eqs. (1) and (2) also describe a system subject to uniform dilution at rate β and the influx of fresh monomers at rates ϕ_i . In this case, the dilution adds terms $-\beta \cdot d_{ij}$ to the rhs of Eq. (1). The dynamics of individual monomer concentrations (both bound and unbound), c_i , is given by equations $\dot{c}_i(t) = \phi_i - \beta c_i(t)$. After a brief transient regime, all monomer concentrations $c_i(t)$ reach a steady state value $\phi_i/\beta = c_i$ so that Eq. (2) become automatically satisfied. In the light of this interpretation, below we focus on the case of all β_{ij} equal to each other. Without loss of generality, they can all be set to unity: $\beta_{ij} = 1$. This defines the fundamental time scale in our system as either the average lifetime of a single bond or the inverse of the dilution rate.

B. Spontaneous entropy reduction

In our previous study,²⁶ we worked within the random sequence approximation (RSA). If all monomers have identical total concentrations $c_i = c$, this approximation corresponds to all 2-mer concentrations $d_{ij}(t)$ being equal to each other. For general initial conditions, these elements would be proportional to $c_i \cdot c_j$. The key hypothesis proposed but not tested in Ref. 26 is that the system dynamics would eventually favor the survival of a subset of the “fittest” sequences at the expense of the others, thus breaking the random sequence approximation. Here we test this hypothesis by simulating the dynamics of the model given by Eqs. (1) and (2) with $Z = 20$. We start with a system characterized by a weak variation in individual ligation rates λ_{ij} and concentrations c_i . We choose them from a log-normal distribution with their logarithms having standard deviation 0.1 and mean values of 0 and $\log 3$, respectively. Our choice of parameters is motivated by the need to understand the limit of infinitesimally weak variation of rates and concentrations. For this combination of parameters, Eq. (1) are initially linearly unstable with respect to formation of all 2-mers. However, no 2-mer would be formed until either it or its complementary partner is present in the system at least in some infinitesimal “seed” concentration. Once such a seed is introduced, the corresponding pair of mutually complementary 2-mers ij and j^*i^* would be exponentially amplified. In our simulations, we used the same small seed concentration of 10^{-4} for each of Z^2 individual 2-mers.

The key parameter we use to quantify the emergent complexity in our system is the information entropy of 2-mers based on their relative concentrations $\tilde{d}_{ij} = d_{ij} / \sum_{kl} d_{kl}$ and defined in the standard Boltzmann-Shannon manner,

$$S(t) = - \sum_{kl} \tilde{d}_{kl}(t) \log \tilde{d}_{kl}(t). \quad (3)$$

Figure 2 shows the time dependence of this entropy in 5 different realizations of λ_{ij} and c_i . The entropy starts at its maximal value $\log(Z^2)$, and after a brief dip followed by a rebound, it steadily *declines* as a function of time. Such behavior is a remarkable manifestation of the non-equilibrium nature of our system, as the entropy changes in the direction opposite to that dictated by the second law of thermodynamics. To reveal the source of this entropy dynamics, in Figs. 2(b) and 2(c), we show the 2-mer matrix at two time points during our simulations. At $t = 2$, all of 2-mers have grown from their seed concentrations to substantial concentrations. Remarkably, the subsequent dynamics leads to a *complete extinction* of the majority of 2-mers ultimately giving rise to the 2-mer matrix at $t = 8000$ shown in Fig. 2(c). The time dependence of the logarithm of the number of surviving 2-mers is shown as red lines in Fig. 2(a). The ultimate number of survivors, 36 ± 4 , is just below $2Z = 40$ (out of $Z^2 = 400$) represented by the lower horizontal dotted line at $\log 2Z$ in Fig. 2(a).

C. Competition between 2-mers and the number of survivors

The observed behavior can be understood from the analysis of Eqs. (1) and (2). For a fixed set of concentrations l_i

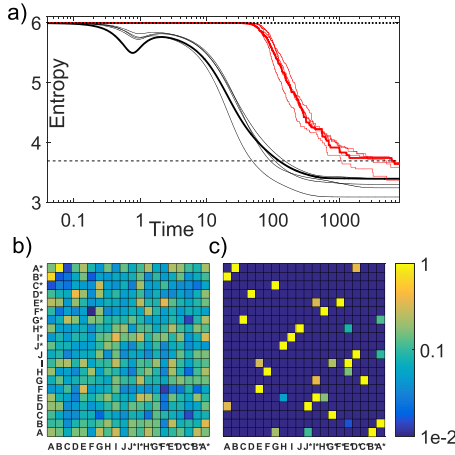


FIG. 2. The entropy of 2-mer concentrations. (a) The information entropy S (black lines) given by Eq. (3) and the natural logarithm of the number of surviving 2-mers N (red lines) plotted vs time in 5 different realizations of our model with logarithms of both λ_{ij} and c_i normally distributed with standard deviation of 0.1 and mean values of $\log 1$ and $\log 3$, respectively. (b) The heatmap visualizing \log_{10} of concentrations of 2-mers at $t = 2$ (the second maximum of the entropy) in one of these realizations highlighted by thick black and red lines in panel (a). (c) The same heatmap in the steady state at $t = 8000$ where the entropy is saturated at its lowest point.

and r_i , Eq. (1) form a set of linear kinetic equations with respect to 2-mer concentrations d_{ij} . Furthermore, this set of Z^2 equations breaks into independent blocks of equations describing the dynamics of mutually complementary 2-mers d_{ij} and $d_{j^*i^*}$. For a small subset of self-complementary 2-mers ii^* , occupying a diagonal of the 2-mer matrix, such a block is represented by a single equation. In all other cases, it involves a pair of equations for d_{ij} and $d_{j^*i^*}$ coupled via a 2×2 matrix,

$$\begin{pmatrix} \dot{d}_{ij} \\ \dot{d}_{j^*i^*} \end{pmatrix} = \begin{pmatrix} -1 & \lambda_{ij}r_i l_j \\ \lambda_{j^*i^*}r_{j^*} l_{i^*} & -1 \end{pmatrix} \begin{pmatrix} d_{ij} \\ d_{j^*i^*} \end{pmatrix}. \quad (4)$$

Because the trace of the matrix is always negative, at least one of the eigenvalues has to have a negative real part, while the real part of the other one could be positive, negative, or zero depending on the value of the matrix determinant Δ_{ij} . A negative value of the determinant $\Delta_{ij} < 0$ corresponds to a positive eigenvalue and hence to the exponential growth of two complementary 2-mer concentrations observed at the initial stage. As growing 2-mers gradually deplete r_i , r_{j^*} , l_{i^*} , and l_j , Δ_{ij} increases and may eventually turn positive. In this case, both eigenvalues become negative. This triggers the exponential decay of concentrations and ultimate extinction of the corresponding pair of 2-mers. A small subset of 2-mers survive and reach the steady state. For these survivors, the determinant *has to become exactly zero*: $\Delta_{ij} = 0$. These conditions for surviving 2-mers can be rewritten as

$$\Delta_{ij} \equiv 1 - \lambda_{ij}r_i \cdot l_j \cdot \lambda_{j^*i^*}l_{i^*}r_{j^*} = 0, \quad (5)$$

while for all extinct 2-mers, $\Delta_{ij} > 0$.

Now we can put the upper bound on the number of surviving 2-mers in the steady state of the system. This is accomplished by comparing the total number of constraints given by Eq. (5) to the number of independent variables. Since the ligation rate matrix λ_{ij} is fixed, the only variable parameters

in Eq. (5) are the left and right end concentrations $l_i(t)$ and $r_i(t)$, respectively. While naively, the number of such variables is $2Z$, Eq. (5) always contain them in combinations $r_i \cdot l_{i^*}$. Therefore, for the purpose of our counting argument, only these Z products should be considered as independent variables. The number of constraints [Eqs. (5)] that are simultaneously satisfied cannot be greater than Z . Each of these equations corresponds to either a pair of mutually complementary 2-mers or a single self-complementary 2-mer. Denoting the total number of surviving 2-mers as N and the number of self-complementary surviving 2-mers as N_{sc} , the number of equations for surviving 2-mers is given by $(N - N_{sc})/2 + N_{sc} = (N + N_{sc})/2$, which has to be smaller than Z —the number of independent variables. Thus the upper bound on the number of surviving 2-mers is given by

$$N \leq 2Z - N_{sc}. \quad (6)$$

Note that for large Z the number of surviving 2-mers is dramatically lower than Z^2 —the total number of possible ones. This explains the entropy reduction observed numerically (see Fig. 1). The parameters of the system were chosen in such a way that initially all Z^2 2-mers grow exponentially. Since the rate of this early exponential growth depends on λ_{ij} and c_i , it differs from one 2-mer to another. This results in a transient behavior where the inhomogeneity of 2-mer concentrations is amplified giving rise to an early decrease in entropy [see the dip around $t = 1$ in Fig. 2(a)]. As concentrations l_i and r_i start to get gradually depleted, the growth saturates, giving time for slower-growing 2-mers to catch up with the faster-growing ones around $t = 2$ [see Fig. 2(b)]. As a consequence, the entropy recovers close to its maximal value. After that, a new process starts in which 2-mers actively compete with each other for the remaining left and right ends. When the determinant Δ_{ij} for a particular pair of 2-mers ij and j^*i^* changes its sign to positive, that pair of 2-mers starts to exponentially decay and eventually goes extinct. This process continues until N , the number of remaining 2-mers with $\Delta_{ij} = 0$, falls below the upper bound given by Eq. (6). These surviving 2-mers are visible as bright spots in the heatmap in Fig. 2(c).

D. Graph-theoretical representations

A useful visualization of the emergent state of the system is the so-called de Bruijn graph shown in Fig. 3(a). It represents each of Z monomer types as a vertex and each of N surviving 2-mers ij as a directed edge connecting vertices i and j . The weight of every edge is proportional to the steady state 2-mer concentration d_{ij} . A de Bruijn graph is a common representation of heteropolymer ensembles such as DNA sequences of all chromosomes in the genome of an organism. It is straightforward to construct it from a known pool of sequences. However, the inverse problem of reconstruction of the statistics of a sequence pool from a de Bruijn graph is highly non trivial. Each of the polymers in the pool can be represented as a walk on this graph. The simplest case is when the consecutive steps of this walk are uncorrelated with each other. This means that the walk is a random Markov process

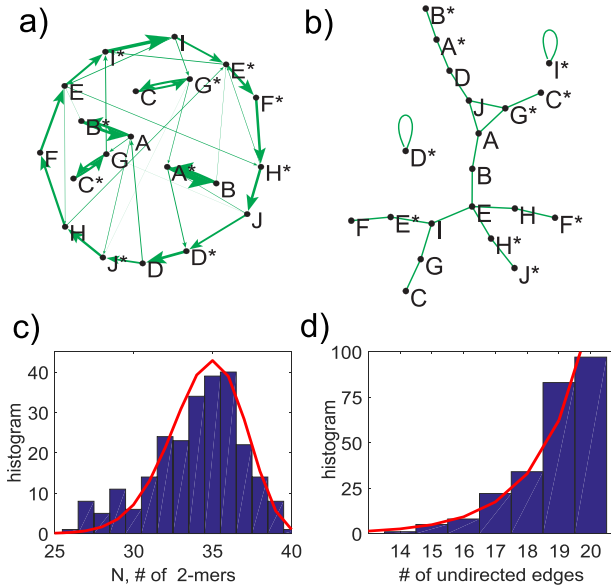


FIG. 3. Network representations. (a) The de Bruijn directed graph with $Z = 20$ nodes corresponding to monomers and edges corresponding to surviving 2-mers. The thickness of each edge scales with 2-mer's concentration. Polymer sequences in our pool are walks on this graph. (b) Undirected graph representation of the system constructed as described in the text. Each edge represents two mutually complementary 2-mers, for instance, $F - E^*$ stands for $F \rightarrow E$ and $E^* \rightarrow F$. (c) The histogram of the number N of surviving 2-mers (directed edges in the de Bruijn graph). (d) The histogram of the number of undirected edges $(N + N_{sc})/2$ in 250 realizations of the model with different λ_{ij} and c_i .

with the probability of a step $i \rightarrow j$ given by d_{ij}/c_i , while the probability of the termination of a polymer at vertex i given by $1 - \sum_j d_{ij}/c_i = r_i/c_i$. This gives rise to an exponential distribution of polymer lengths, the same as in our previous study²⁶ where it was obtained within the random sequence approximation. The average length of chains is the ratio of the total number of all monomers to the total number of right (or, equivalently, left) ends of chains,

$$\langle L \rangle = \sum c_i / \sum r_i = \sum c_i / \sum l_i. \quad (7)$$

Since in the steady state for surviving dimers $\lambda_{ij}r_i l_j \sim 1$, one has $l \simeq r1/\sqrt{\lambda}$. Hence, the average chain length can be estimated as $\langle L \rangle \sim c \cdot \sqrt{\lambda}$, which again is similar to the results of Ref. 26.

Note that the entropy defined above and plotted in Fig. 2(a) is *exactly* the information entropy of a pool of polymer sequences generated by such a Markov process.³⁰ Longer-range correlations between different 2-mers in the polymer sequence are not captured by the present model but could in principle emerge due to effects outlined in Sec. III. Such correlations would lead to further reduction of the information entropy in the system.

The de Bruijn graph can be complemented by another, more compact graphical representation, which is specific to our system. Since mutually complementary 2-mers always appear in pairs ij and j^*i^* (with the exception of self-complementary 2-mers ii^*), each such pair can be depicted as a single undirected edge connecting vertices i to j^* . In this representation, each edge represents two 2-mers, while each vertex i stands for either i or i^* monomer, depending on whether it is the

first or the second letter within the 2-mer. In Appendix A, we show that this undirected graph [see Fig. 2(b)] has a number of remarkable properties. First, it is a so-called “pseudoforest”.²⁸ each of its individual connected components contains no more than one cycle. This allows us to refine and give a topological interpretation to Eq. (6): $N = 2Z - N_{sc} - 2N_{trees}$, where N_{trees} is the number of trees (components without cycles) in the pseudoforest. Second, only the odd-length cycles (1,3,5, etc.) are allowed in this graph.

Figure 3(c) shows the distribution of the number of surviving 2-mers [or equivalently of directed edges in the de Bruijn graph shown in Fig. 3(a)] in 250 realizations of the system with different values of λ_{ij} and c_i . Figure 3(d) shows the distribution of the number $(N + N_{sc})/2$ of edges in undirected graphs such as the one shown in Fig. 3(b) for the same set of realizations. As discussed above, the deviation of this last quantity down from Z is equal to the number of trees in the pseudoforest. As shown in Fig. 3(d), it can be approximated by an exponential distribution with the average around 1.6 [the red line in Fig. 3(c)]. At the same time, the distribution of the number of surviving 2-mers (N) has a peak around $36 < 2Z = 40$. The quantity $2Z - N$ is always positive and approximately follows a Poisson distribution with the average of 5.5 [the red line in Fig. 3(c)].

E. Variability of the set of surviving 2-mers

The set of surviving 2-mers and their concentrations depend on a number of parameters: ligation rates λ_{ij} , total monomer concentrations c_i , and, possibly, seed concentrations of individual 2-mers.

We analyzed the sensitivity of the steady state of our system with respect to all of these parameters one-by-one. First we fixed both λ_{ij} and c_i and analyzed the final 2-mer concentrations for a large number of random realizations of Z^2 small (but positive) seed concentrations. We found the final state to be reproducible as long as all seed concentrations are non-zero. Note that due to the autocatalytic nature of 2-mer dynamics given by Eq. (1), a pair of complementary 2-mers with zero seed concentrations would never emerge on their own. To search for alternative stable states that are not accessible starting with small random seed concentrations of all 2-mers, we performed an additional test that was explicitly biased toward finding other stable solutions (if they exist). The following protocol was implemented: (i) we determined the set of surviving 2-mers for a specific set of λ_{ij} and c_i and rerun our dynamics with the initial seed concentrations of these 2-mers being artificially set to 0. As expected, the resulting set of survivors did not include any 2-mers from the excluded set. (ii) Following that, we added 2-mers from the excluded set at very small seed concentrations. In some cases, that was sufficient for them to completely take over the system, thereby returning it to exactly the same steady state as in our standard protocol. However, in a significant fraction of cases, we observed the appearance of a new steady state in the system. Multistability was also observed in other models of prebiotic evolution such as that in Ref. 29.

Next, we fixed the ligation rates λ_{ij} to their values used to construct the heatmaps shown in Fig. 2 and networks

in Fig. 3. We then simulated 100 realizations of the system with c_i pulled from a log-normal distribution $P(c_i) \sim \exp(-[\log(c_i/c)]^2/\sigma_c)/c_i$ with $c = 3$ and $\sigma_c = 0.1$. Figure 4(a) shows the heatmap of the fraction of realizations of c_i in which each individual 2-mer survives in the steady state. In Fig. 4(b), we present the same results in the form of the histogram (blue bars). The majority of 2-mers [324 out of 400 visible as the leftmost bar in Fig. 4(b)] got extinct in all of 100 realizations. While the number of surviving 2-mers in each realization never exceeded $2Z = 40$ [see Fig. 2(c)], the overall number of 2-mers that survived in at least one realization of c_i was substantially larger: 76. Out of this set, 20 “universal survivors” were present in all 100 realizations. Furthermore, as can be seen by comparing heatmaps in Figs. 2(d) and 4(a), these 20 universal surviving 2-mers typically have high steady state concentrations d_{ij} . To further investigate the correlation between 2-mer concentration and its likelihood to survive, we analyzed a larger set of 250 realizations of the system with the same fixed ligation matrix but variable monomer concentrations c_i . The resulting distribution of 2-mer concentrations shown in Fig. 4(c) is clearly bimodal. The bi-modality is also apparent from an example of the de Bruijn graph shown in Fig. 3(a), where approximately half of all links (thicker lines) correspond to more abundant 2-mers, while the other half (thinner lines) correspond to 2-mers present at low concentrations. The high-concentration peak of the distribution in Fig. 4(c) is dominated by the contribution from 20 universal survivors shown as the red line. Note that despite an increased number of realizations,

the set and number of universal survivors stayed the same as in Figs. 4(a) and 4(b).

We further investigate the variability of the set of surviving 2-mers as a function of the width σ_c of the log-normal distribution of monomer concentrations. As shown in Fig. 4(d), the number of universal survivors (red line) systematically decreases with σ_c , ultimately reaching 0 for $\sigma_c \geq 0.75$. Consistent with this trend, the bimodality of the concentration distribution in Fig. 4(c) also disappears for larger values of σ_c . Note that all numbers of 2-mers shown in Fig. 4(d) were normalized by $2Z$, i.e., the upper bound of the number of surviving 2-mers in each realization. The average number of survivors in a single realization (green line) does not have a notable dependence on σ_c . The blue line in Fig. 4(d) shows the number of 2-mers (normalized by $2Z$) that survived at least once among 100 realizations of c_i . Note that this curve grows significantly with σ_c ultimately reaching the value as high as 5. This corresponds to half of all $Z^2 = 400$ possible 2-mers having a chance to survive in at least one of these realizations.

III. DISCUSSION

The major conclusion following from this study is that our model system of mutually catalyzing heteropolymers has a natural tendency toward spontaneous reduction of the information entropy. This represents an effective *reversal* of the second law of thermodynamics in this class of systems. While “violations” of the second law are indeed expected in externally driven non-equilibrium systems, the observed “reversal” has much greater implications. Both living organisms and other self-organized systems such as human culture, economics, and technology are characterized by an ever increasing complexity, indicating the ongoing reduction in the information entropy.

The thermodynamic entropy of a system of heteropolymers is composed of two distinct parts:³⁰ (i) the translational and configurational entropy of polymer chains and (ii) the information entropy associated with sequence statistics. Our current model hints at a hierarchical scenario of entropy reduction in populations of heteropolymers. First, the translational entropy is reduced due to template-based polymerization as studied in our previous work²⁶ within the Random Sequence Approximation (RSA). Then RSA breaks down at the level of 2-mers due to their competition with each other for a limited resource of monomers. Such symmetry breaking in the sequence space results in a dramatic reduction of the information entropy (ii). At this point, sequences of the entire pool of chains are generated as Markovian random walks on the de Bruijn graph [Fig. 3(a)]. One can imagine a further reduction of the information entropy due to emergence of correlations between consecutive steps of this walk.

There are multiple physical scenarios outside the scope of the current model that would lead to such longer-range correlations in the sequence space. They include, e.g., a dependence of the association rates in Eq. (1) on the lengths of the three chains involved in the process of template-assisted ligation. Another intriguing scenario is a spontaneous emergence of chains with weak catalytic activity above and beyond their role as templates for ligation. For instance, some sequences

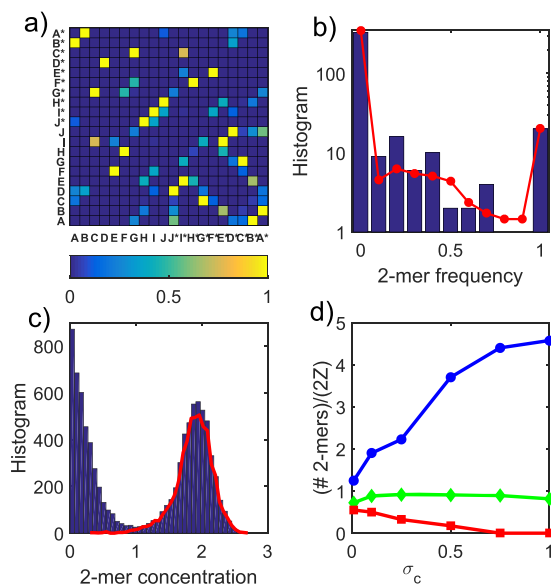


FIG. 4. Variability of surviving 2-mers for different realizations. (a) The heatmap of survival frequency of 2-mer in the steady state of the system with fixed λ_{ij} (the same as in Figs. 2 and 3) and 100 different realizations of c_i . Note the similarity to Fig. 2(c). (b) The histogram of these survival frequencies (blue bars) as well as its average (red line) over 10 different realizations of λ_{ij} . (c) The histogram of 2-mer concentrations in 250 realizations of c_i for fixed λ_{ij} . The red line shows the contribution from 20 “universal survivors”—2-mers present across all realizations. (d) The number of universal survivors (red line), the average number of survivors in a single realization (green line), and the number of 2-mers present in at least one among hundred realizations of c_i (blue) each normalized by $2Z$. The x-axis is the width σ_c of the log-normal distribution of monomer concentrations.

may facilitate breakage or conversely promote ligation reactions, either among some specific sequences or universally. Such sequences would provide a missing link between the prebiotic soup considered here and the emergence of the first ribozymes in the RNA world scenario.

A common pattern in functional RNA-based structures such as ribozymes and ribosomes is the presence of hairpins and loops. Note, however, that the mass-action term in Eq. (1) assumes that the template and the two substrates belong to three different chains thus ignoring these higher order structures. A proper model description taking them into account is a natural next step in the development of our approach. An important question for a future study is whether our system naturally evolves toward or away from loops and hairpins.

An interesting feature of the sequence statistics emerging in our model is that the entropy does not reach its absolute minimum that would correspond to a unique “master sequence.” In the de Bruijn representation, such a master sequence would look, e.g., like a single cycle (or several unconnected cycles) in which for every monomer there is unique right neighbor following it on every chain. In that case, the walk on the de Bruijn graph would be completely deterministic. By contrast, in our model, each monomer typically has two possible right neighbors in the de Bruijn graph. In the limit of relatively small variations $\sigma_c = 0.1$ used to construct Fig. 3(a), stronger links corresponding to “universal survivors,” produce the c_i -independent backbone of the graph akin to the master sequence. Conversely, weaker links allow for infrequent hopping between different parts of the graph resulting in deviations from this master sequence. The opposite limit of large variations in monomer concentrations is especially promising from the point of view of further evolution of our system. In that limit, there is no dominant master sequence. This dramatically expands the explored region in the sequence space: now, for every step of the Markovian walk, there are on average two comparable probabilities for selecting the next monomer. As a result, the number of possible sequences of length L in our model, $\sim 2^L$, remains exponentially large, yet dramatically reduced compared to its random sequence limit, Z^L . Furthermore, in the limit of large σ_c , different realizations of c_i give rise to significantly different sets of surviving 2-mers. Note that, unlike ligation rates λ_{ij} , monomer concentrations c_i (or equivalently their influxes ϕ_i) could vary significantly from one spatial location to another (see Ref. 31 for a study of spatial inhomogeneity in the prebiotic context). This allows for an effective exploration of various regions (of size $\sim 2^L$ each) of the global sequence space, rather than converging to the same subset of sequences over and over.

Our model describes a simple yet general mechanism for spontaneous entropy reduction in a system capable of template-assisted ligation. There are multiple possible experimental realizations of such systems based on either traditional DNA/RNA biochemistry or artificial micro/nanostructures. The most direct implementation of our model would be a system composed of Z words made of a string of nucleotides bound together by strong (e.g., DNA-type) bonds. They have to be designed to form $Z/2$ mutually complementary pairs that are orthogonal to each other; i.e., words from different pairs

have no long overlaps. These words would play the role of composite monomers in our model that could be subsequently connected to each other with weaker, breakable (e.g., RNA-type) bonds.³² As discussed earlier, our model is directly applicable to the scenario in which all bonds are unbreakable, while the whole system is uniformly diluted and fresh (unbound) monomers are supplied at a constant rate. This greatly expands possibilities for its experimental implementation. In particular, one may now use purely DNA-based words not worrying about how to break their rather stable bonds. Similarly, our dynamics can be realized using the DNA origami nanoblocks introduced in Ref. 6. It should be emphasized that in order to achieve the behavior described by our model, the experiments need to be conducted well below the saturation regime; i.e., the night phase should be shorter than a typical association time for hybridization.

Yet another interpretation of our model does not involve any long polymers at all. In this case, 2-mers are represented by physical dimers made of only two monomeric units incapable of forming longer chains. Our model predicts that, even in this simple system, the compositional entropy would drop due to extinction of most of the dimers leaving no more than $2Z$ survivors. On the one hand, this further extends possibilities for experimental implementation. For instance, one could construct the DNA-based system described above but limited to chains no longer than 2 words. This would greatly reduce the complexity of the screening process. On the other hand, this dimer interpretation has an intriguing connection to the Kauffman model of autocatalytic chemical reaction networks.¹² In the Kauffman case, some of the molecules in the pool are capable of catalyzing the synthesis of others from “raw materials” (abundant small metabolites) ultimately resulting in the emergence of metabolic autocatalytic cycles in large systems. A recent model of such chemical reactions³³ shows that the system self-organizes to a state finely tuned to the external driving force. This can be interpreted as the maximization of the rate of negentropy adsorption from the environment. In our case, dimers correspond to mutually catalytic entities, while monomers represent raw materials. While in the current implementation of our model, the catalytic cycles could only involve two mutually complementary dimers, it is straightforward to generalize the model to allow for catalytic activity of any dimer toward any other dimer or dimers. That version of the model allows for autocatalytic cycles of length longer than two. We expect our findings about the reduction of entropy to be fully transferable to that case. Thus, our model has a potential of bridging the gap between two traditionally competitive visions of the Origin of Life: “information first” and “metabolism first.”

ACKNOWLEDGMENTS

We would like to thank Paul Higgs who has made a number of valuable comments on this manuscript and, in particular, persuaded us to look harder for multiple stable states in our model. This research used resources of the Center for Functional Nanomaterials, which is a U.S. DOE Office of Science User Facility, at the Brookhaven National Laboratory under Contract No. DE-SC0012704.

APPENDIX A: TOPOLOGICAL ANALYSIS OF THE UNDIRECTED GRAPH

Here we discuss the undirected graph representation of the pool of heteropolymers. Since our system always contains pairs of mutually complementary 2-mers ij and j^*i^* (with the exception of self-complementary 2-mers ii^*), one can represent this pair with a single undirected edge connecting i to j^* . These edges form an undirected graph shown in Fig. 3(b). Note that due to these rules an edge connecting vertices i and k in this graph represents a pair of 2-mers ik^* and ki^* shown as two edges in Fig. 3(a) or two matrix elements in Fig. 2(c). For simplicity, in Fig. 3(b) we did not assign weights to these symmetric edges. Each edge of this graph corresponds to exactly one equation in the set of Eq. (5). Hence, the number of undirected edges is equal to $(N + N_{sc})/2$ and according to Eq. (6) it cannot exceed Z . On the other hand, based on network topology, this number can be expressed as $Z - N_{comp} + N_{cycles}$. Here N_{comp} is the number of connected components of the graph, while N_{cycles} is the number of independent cycles defined as the minimal number of edges one needs to cut to remove all cycles. The inequality (6) means that the number of independent cycles cannot be larger than the number of components. Furthermore, this inequality must be also satisfied for each of the individual connected components of the graph because the number of equations (edges) cannot exceed the number of independent variables (the number of vertices in this component). In other words, *each of the components may contain no more than one cycle*. Graphs with this property are known as “pseudoforests.”²⁸ For unicyclic components (i.e., those that include exactly one cycle), the numbers of edges and vertices are equal to each other, while for each of the tree (cycle-free) components, the difference between these two numbers is equal to 1. This gives a topological interpretation to the number of surviving 2-mers N in terms of the number of tree-like components of the undirected graph, N_{trees} ,

$$N = 2Z - N_{sc} - 2N_{trees}, \quad (A1)$$

which automatically leads to inequality (6).

One can also demonstrate that only the cycles of odd lengths (1,3,5, etc.) are allowed in our system. Indeed, for a hypothetical even-length cycle $i_1 - i_2^* - i_3 - i_4^* - \dots - i_{n-1} - i_n^* - i_1$, one can construct a combination of Eq. (5) of the following form:

$$\frac{\Lambda_{i_1 i_2} \Lambda_{i_3 i_4} \dots \Lambda_{i_{n-1} i_n}}{\Lambda_{i_2 i_3} \dots \Lambda_{i_{n-2} i_{n-1}} \Lambda_{i_n i_1}} = 1. \quad (A2)$$

Here $\Lambda_{ij} = 1 - \Delta_{ij} = \lambda_{ij} \cdot \lambda_{j^* i^*} \cdot r_i \cdot l_j \cdot l_{i^*} \cdot r_{j^*}$. All the variables l_i and r_i at the left-hand-side of this equation cancel, making it an invariant that depends only on ligation rates. Therefore this equation cannot be satisfied for a generic matrix λ_{ij} , which rules out the existence of even cycles in most of the cases.

APPENDIX B: ADDITION OF NON-TEMPLATED LIGATION TO OUR MODEL

To test the robustness of our results, we explored a variant of our model in which in addition to template-assisted

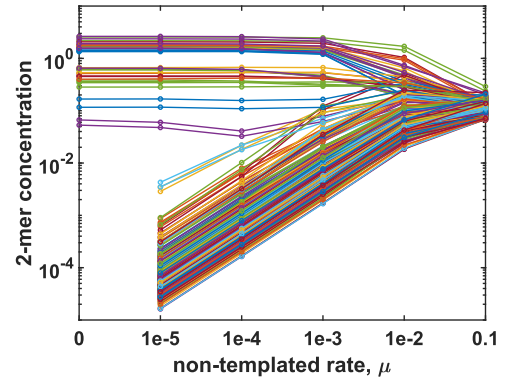


FIG. 5. Steady state concentrations of 2-mers plotted vs the rate μ of non-templated ligation. All other parameters of the model are the same as in Fig. 1. Note that all Z^2 2-mers survive for any $\mu > 0$, while the original set of 38 survivors for $\mu = 0$ continue to dominate in terms of abundances for low values of μ .

ligation, polymer chains can also merge spontaneously without the help of any template. Such a process is described by an additional positive term $\mu \cdot r_i l_j$ in Eqs. (1). The equations thereby become

$$\dot{d}_{ij}(t) = (\lambda_{ij} d_{j^* i^*}(t) + \mu) r_i(t) l_j(t) - \beta_{ij} d_{ij}(t). \quad (B1)$$

We computed the evolution of our system for several values of μ ranging from 10^{-5} to 0.1. The results are shown in Fig. 5. As one can see, our conclusions are robust with respect to introduction of a non-zero (but small) non-templated ligation rate μ . While all Z^2 2-mers now survive in the steady state, their distribution remains clearly bimodal for small values of μ . Figure 5 traces the value of the steady state 2-mer concentrations as a function of μ . When non-templated ligation becomes strong, the distinction between original survivors and the rest of 2-mers gradually disappears.

¹J. Lovelock, *GAIA—A New Look at Life on Earth* (Oxford University Press, New York, 1979).

²Strictly speaking, due to the exchange of energy and material with the environment, the appropriate measure of the reduction of thermodynamic entropy by living systems is the local accumulation of Gibbs free energy. This was pointed out by Schrödinger in Ref. 3.

³E. Schrödinger, *What Is Life—The Physical Aspect of the Living Cell* (Cambridge University Press, Cambridge, 1944).

⁴L. Brillouin, *Science and Information Theory* (Dover, Mineola, New York, 1956).

⁵C. E. Shannon, “A mathematical theory of communication,” *Bell Syst. Tech. J.* **27**, 379 (1948).

⁶T. Wang, R. Sha, R. Dreyfus, M. E. Leunissen, C. Maass, D. J. Pine, P. M. Chaikin, and N. C. Seeman, *Nature* **478**, 225 (2011).

⁷Z. Zeravcic and M. P. Brenner, *Proc. Natl. Acad. Sci. U. S. A.* **111**, 1748 (2014).

⁸Z. Zeravcic and M. P. Brenner, “Spontaneous emergence of catalytic cycles with colloidal spheres,” *Proc. Natl. Acad. Sci. U. S. A.* **114**, 4342 (2017).

⁹J. M. Dempster, R. Zhang, and M. de la Cruz, “Self-replication with magnetic dipolar colloids,” *Phys. Rev. E* **92**, 042305 (2015).

¹⁰M. Eigen and P. Schuster, *Naturwissenschaften* **64**, 541 (1977).

¹¹F. J. Dyson, *J. Mol. Evol.* **18**, 344 (1982).

¹²S. A. Kauffman, “Autocatalytic sets of proteins,” *J. Theor. Biol.* **119**, 1 (1986).

¹³S. Jain and S. Krishna, “Autocatalytic sets and the growth of complexity in an evolutionary model,” *Phys. Rev. Lett.* **81**, 5684 (1998).

¹⁴M. P. Robertson and G. F. Joyce, “The origins of the RNA world,” *Cold Spring Harbor Perspect. Biol.* **4**, a003608 (2012).

- ¹⁵J. A. Doudna and J. W. Szostak, "RNA-catalysed synthesis of complementary-strand RNA," *Nature* **339**, 519 (1989).
- ¹⁶T. A. Lincoln and G. F. Joyce, "Self-sustained replication of an RNA enzyme," *Science* **323**, 1229 (2009).
- ¹⁷W. Gilbert, "Origin of life: The RNA world," *Nature* **319**, 618 (1986).
- ¹⁸L. E. Orgel, "Prebiotic chemistry and the origin of the RNA world," *Crit. Rev. Biochem. Mol. Biol.* **39**, 99 (2004).
- ¹⁹C. B. Mast, S. Schink, U. Gerland, and D. Braun, *Proc. Natl. Acad. Sci. U. S. A.* **110**, 8030 (2013).
- ²⁰P. W. Anderson, "Suggested model for prebiotic evolution: The use of chaos," *Proc. Natl. Acad. Sci. U. S. A.* **80**, 3386 (1983).
- ²¹P. W. Anderson and D. L. Stein, in *Self-Organizing Systems*, edited by F. E. Yates, A. Garfinkel, D. O. Walter, and G. B. Yates (Springer US, 1987), pp. 445–457.
- ²²W. Hordijk, S. A. Kauffman, and M. Steel, "Required levels of catalysis for emergence of autocatalytic sets in models of chemical reaction systems," *Int. J. Mol. Sci.* **12**, 3085 (2011).
- ²³H. Fellermann, S. Tanaka, and S. Rasmussen, "Sequence selection by dynamical symmetry breaking in an autocatalytic binary polymer model," *Phys. Rev. E* **96**, 062407 (2017).
- ²⁴R. Rohatgi, D. P. Bartel, and J. W. Szostak, "Kinetic and mechanistic analysis of nonenzymatic, template-directed oligoribonucleotide ligation," *J. Am. Chem. Soc.* **118**, 3332 (1996).
- ²⁵J. Derr, M. L. Manapat, S. Rajamani, K. Leu, R. Xulvi-Brunet, I. Joseph, M. A. Nowak, and I. A. Chen, "Prebiotically plausible mechanisms increase compositional diversity of nucleic acid sequences," *Nucleic Acids Res.* **40**, 4711 (2012).
- ²⁶A. V. Tkachenko and S. Maslov, *J. Chem. Phys.* **143**, 045102 (2015).
- ²⁷A. S. Tupper, K. Shi, and P. G. Higgs, "The role of templating in the emergence of RNA from the prebiotic chemical mixture," *Life* **7**, 41 (2017).
- ²⁸G. B. Dantzig, *Linear Programming and Extensions* (Princeton University Press, Princeton, NJ, 1963).
- ²⁹P. G. Higgs, "Chemical evolution and the evolutionary definition of life," *J. Mol. Evol.* **84**, 225 (2017).
- ³⁰D. Andrieux and P. Gaspard, "Nonequilibrium generation of information in copolymerization processes," *Proc. Natl. Acad. Sci. U. S. A.* **105**, 9516–9521 (2008).
- ³¹S. Walker, M. A. Grover, and N. V. Hud, "Universal sequence replication, reversible polymerization and early functional biopolymers: A model for the initiation of prebiotic sequence evolution," *PLoS One* **7**, e34166 (2012).
- ³²Y. Li and R. R. Breaker, "Kinetics of RNA degradation by specific base catalysis of transesterification involving the 2' hydroxyl group," *J. Am. Chem. Soc.* **121**, 5364 (1999).
- ³³J. M. Horowitz and J. L. England, "Spontaneous fine-tuning to environment in many-species chemical reaction networks," *Proc. Natl. Acad. Sci. U. S. A.* **114**, 7565 (2017).