

# Thermodynamic consistency of autocatalytic cycles

Thomas Kosc,<sup>1</sup> Denis Kuperberg,<sup>2</sup> Etienne Rajon,<sup>1</sup> and Sylvain Charlat<sup>1,\*</sup>

<sup>1</sup>*LBBE (Laboratoire de Biométrie & Biologie Evolutive) Université Lyon 1, CNRS, Villeurbanne, France*

<sup>2</sup>*LIP (Laboratoire de l'Informatique du Parallélisme), ENS Lyon, CNRS, France*

Autocatalysis is seen as a potential key player in the origin of life, and perhaps more generally in the emergence of Darwinian dynamics. Building on recent formalizations of this phenomenon, we tackle the computational challenge of exhaustively detecting autocatalytic cycles in reactions networks, and further evaluate the impact of thermodynamic constraints on their realization under mass action kinetics. We first characterize the complexity of the detection problem by proving its NP-completeness. This justifies the use of constraint solvers to list all autocatalytic cycles in a given reaction network, and also to group them into compatible sets, composed of cycles whose stoichiometric requirements are not contradictory. Crucially, we show that the introduction of thermodynamic realism does constrain the composition of these sets. Compatibility relationships among cycles can indeed be disrupted when the reaction kinetics obey thermodynamic consistency throughout the network. On the contrary, these constraints have no impact on the realizability of isolated cycles, unless upper or lower bounds are imposed on the concentrations of the reactants. Overall, by better characterizing the conditions of autocatalysis in complex reaction systems, this work brings us a step closer to assessing the contribution of this collective chemical behavior to the emergence of natural selection in the primordial soup.

**Significance Statement:** Describing the processes behind the origin of life requires us to better understand self-amplifying dynamics in complex chemical systems. Detecting autocatalytic cycles is a critical but difficult step in this endeavor. After characterizing the computational complexity of this problem, we investigate the impact of thermodynamic realism on autocatalysis. We demonstrate that individual cycles, regardless of thermodynamic parameters, can be activated as long as entities may occur at any required concentration. In contrast, two cycles can become mutually incompatible due to thermodynamic constraints, and will thus never run simultaneously. These results clarify the implications of physical realism on the realization of autocatalysis.

## I. INTRODUCTION

It is increasingly recognized that producing a consistent explanation for the origination of life will require us to explain how Darwinian evolution may have *gradually* emerged from a non-biological, purely physical world [1–7]. Gradually rather than suddenly, that is, without assuming that natural selection only came into play once chance alone had produced the first obvious “replicators”, displaying the same heritable variance as current organisms. Under this perspective of a smooth transition from physics to biology, natural selection is hypothesized to have been active already in the “prebiotic” soup as a driver to complexity; yet in a rudimentary and currently unrecognizable fashion.

To explore this path, autocatalysis is often taken as a plausible starting point (Box 1) [4, 8–12]. Here, more specifically, we envision autocatalytic cycles as the putative elementary components of higher level systems that may engage in “increasingly Darwinian” dynamics. In doing so, we aim at keeping the best of the two traditionally opposed approaches to the origin of life: physico-chemical realism of the metabolism-first

view, and evolvability of the gene-first perspective. Beyond the specifics of terrestrial life, progressing toward an articulation of Darwinian principles with physics appears as a prerequisite to assess their putative relevance to other physical systems [1, 3].

We build on recent theoretical and computational developments [8, 13, 14] to systematically search for autocatalytic cycles in reaction networks and then assess their thermodynamic consistency, i.e. the impact of thermodynamic constraints on their realization. We first prove that finding autocatalytic cycles in the network is an NP-complete problem – a question that was left open by earlier work [15, 16] – and converge with other authors in using constraint solvers as a technical solution [8, 14]. We then question whether such autocatalytic cycles, defined on the sole basis of the reaction network topology, can also be realized once thermodynamic constraints are introduced. To do so, we take into account the reaction kinetics that themselves depend on the Gibbs free energies and concentrations of the reactants, and the activation barriers of the reactions. We show that regardless of these physical quantities, any potential autocatalytic cycle may be instantiated in some region of the concentration space as long as one assumes this space is unbounded. In contrast, thermodynamic constraints do restrain compatibility relationships between autocatalytic cycles and will thereby impact the dynamics of complex chemical networks.

### Box 1: Related work on autocatalysis

The present study takes place within a flourishing body of literature taking autocatalysis as a plausible primary component of proto-biotic or proto-Darwinian systems. Our model contrasts from those based on the RAF framework [17, 18] in that it follows a bottom-up approach to autocatalysis: rather than setting catalytic relationships between components of the system and randomly picked reactions, we let the reaction network generate (or not) these relationships, as formalized by Blokhuis et al [13]. Catalysis and autocatalysis then simply emerge as pathways in the reaction network involving elements that act both as reactants and products. For example, in the reactions  $A+C \rightarrow AC$ ,  $AC+B \rightarrow ABC$ ,  $ABC \rightarrow AB +$

\* sylvain.charlat@univ-lyon1.fr

C, the C element can be simply described as a catalyst of the  $A+B \rightarrow AB$  reaction.

In taking such a bottom-up angle, our framework is much related to that of several recent studies [8–10, 14, 19, 20]. Some of these have considered the implications of thermodynamic constraints and mass action kinetics on specific autocatalytic motifs [9, 10, 19]. Others have implemented tools for the exhaustive detection of autocatalysis [8, 14, 20]. Here we jointly consider these two components of the problem, i.e. exhaustive detection and thermodynamic realism.

On a more conceptual ground, we share with Baum et al [4] the view that collections of autocatalytic cycles, rather than cycles alone, might constitute the scale at which incipient heritable variations may occur.

## II. FRAMEWORK AND DEFINITIONS

We analyze networks of bidirectional reactions governed by mass action kinetics. This is typically the case in reactions that simply consist in the association of two entities and the reciprocal dissociation (e.g.  $A+B \rightleftharpoons AB$ ). The entities are fully defined by their composition (e.g.  $A_2B_2$  is not distinct from  $B_2A_2$ ). Given a list of entities, this rule sets the list of all possible reactions, only some of which are assumed to exist to generate a particular reaction network – this is equivalent to assuming that some reactions have an infinite activation barrier and thus do not take place.

We can then apply the formalism of Blokhuis et al [13] to identify autocatalytic motifs in such reaction networks. Here, these motifs are more specifically referred to as *potential* autocatalytic cycles (PAC), to emphasize that they are defined on the sole basis of the reaction network topology, so that their realizability under thermodynamic constraints remains to be assessed. Intuitively, a PAC can be conceived as a cyclic sub-network admitting a regime where each entity has a positive net production rate.

**Definition 1.** A PAC is defined as a set of entities  $E_C$  and reactions  $R_C$  such that:

- For each reaction  $R \in R_C$ , at least one entity on each side is in  $E_C$ .
- There exists a vector  $\vec{v}$  of flows for reactions from  $R_C$  defining a regime where the total contribution of reactions from  $R_C$  is positive for each entity of  $E_C$ .

Consider for instance a reaction  $A+B \rightleftharpoons AB$ , with reactants  $A, B$  and product  $AB$ . Then the stoichiometry  $\sigma_A^R$  of  $A$  in  $R$  is  $-1$  (or  $-2$  if  $A=B$ ), and  $\sigma_{AB}^R$  is  $1$ . If an entity  $e$  does not appear in a reaction  $R$ , we define  $\sigma_e^R = 0$ .

We will note  $v_R$  the flow of the reaction at a given instant, that will be positive if the association rate is larger than the dissociation rate.

Given a PAC candidate  $C$  formed of entities  $E_C$  and reactions  $R_C$ , and an entity  $e$  in  $E_C$ , we define the variation of  $e$ 's concentration **due to**  $C$  as:

$$\Delta_C(e) = \sum_{R \in R_C} \sigma_e^R \cdot v_R \quad (1)$$

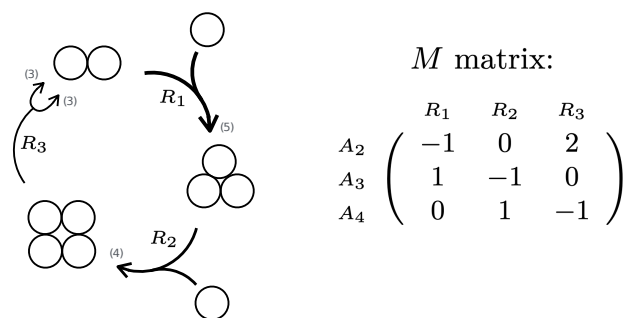


Figure 1. Schematic view of a formose-like potential autocatalytic cycle. Entities  $A_2$ ,  $A_3$  and  $A_4$  are part of the PAC while entity  $A_1$  serves as food. The arrows indicate the net direction of each reaction, while the line width indicates their respective flows, that must be decreasing from reactions  $R_1$  to  $R_3$  for the cycle to run. As an example, the flows values (indicated in brackets) would produce a net increase of 1 of each entity. The right panel shows the corresponding stoichiometric matrix  $M$ . Given the represented flow vector  $\vec{v} = (5, 4, 3)$ , we obtain  $M \cdot \vec{v} = (1, 1, 1)$ , showing that  $\vec{v}$  is indeed a PAC witness.

We define a PAC witness as a choice of  $v_R$  for each  $R \in R_C$ , such that for each  $e \in E_C$  we have  $\Delta_C(e) > 0$ . This can be formalized using linear algebra, following Blokhuis et al [13]. Indeed, if  $M$  is the stoichiometric matrix restricted to  $E_C$  and  $R_C$ , then the candidate  $C$  is a PAC if and only if there exists a witness vector  $\vec{v}$  such that all coordinates of  $M\vec{v}$  are strictly positive.

**Example 1.** We illustrate the PAC definition in Figure 1 with a simple formose-like cycle comprising three entities  $A_2$ ,  $A_3$  and  $A_4$  and using entity  $A_1$  as food, with detailed explanations provided in the caption.

Entities appearing in  $R_C$  that are not part of  $E_C$  will be called either “food” if they are consumed or “waste” if they are produced by a reaction of  $R_C$ , taking into account the sign of the witness vector that indicates the direction of reactions. Notice that an entity may simultaneously appear as food and waste in a PAC.

A PAC is said to *contain* another one if it includes all its entities and reactions. A PAC is *minimal* if it does not contain any other, in which case it corresponds to the “autocatalytic core” from Blokhuis et al [13] and to the stoichiometric autocatalysis of Gagrani et al [14]. In the following, we will focus on such minimal PACs, often omitting the “minimal” adjective for simplicity.

Notably, it is shown in Blokhuis et al [13] that in a minimal PAC, each entity is reactant of a unique reaction, and each reaction has a unique entity of the PAC as reactant (other reactants being food). This implies that the direction of reactions is consistent across all witnesses of a given PAC: flipping the direction of one reaction would force to flip all the others. Therefore, each reaction of the PAC has a unique possible net direction, that will be compatible with all its witnesses flow vectors.

### III. DETECTING POTENTIAL AUTOCATALYTIC CYCLES

Our goal is to enumerate all PACs in a reaction network. To this end, we first assess the complexity of this problem, in order to determine which computational tools are required for its resolution.

#### A. NP-completeness proof

Enumerating all PACs in a reaction network involves sequentially solving problems of the type: “is there a PAC in the system besides those previously found?”. We will show that a particular case of this question is already NP-complete, which justifies the use of an SMT solver. Namely, we will prove the NP-completeness of deciding whether a PAC exists that contains an entity  $A$  and takes food from a given subset  $F$ .

In this section, for simplicity, we will relax any compositionality constraint on reactions so that letters like  $A, B, E \dots$  will be shorthand for any kind of entity. Yet it would be straightforward (but less readable) to extend the construction to a strictly compositional framework.

Notably, the complexity of the autocatalysis detection problem has previously been considered by Andersen et al [15] but from a different angle. These authors have specifically shown that the following problem is NP-complete: considering a reaction network that contains a known autocatalytic cycle, can its resources be produced by the network? The difficulty of finding all autocatalytic cycles in a reaction network, that we tackle here, has thus not yet been addressed.

In the framework of Blokhuis et al [13], it is easy to check whether a proposed set of entities and reactions constitutes an autocatalytic cycle. Indeed, thanks to the linear algebra formulation summarized in Section II this problem is solved in polynomial time by Linear Programming. As will be shown, the difficulty rather lies in finding an autocatalytic cycle in a reaction network, among exponentially many possible candidates.

Formally, let PAC-DETECTION be the following algorithmic problem:

**Definition 2** (PAC-DETECTION problem).

**INPUT:** A reaction system defined by entities ( $E$ ) and reactions ( $R$ ), a target  $A \in E$ , and a set of allowed foods  $F \subset E$ .

**OUTPUT:** Is there a PAC containing  $A$  and using only foods from  $F$ ?

**Theorem 1.** PAC-DETECTION is NP-complete.

The detailed proof can be found in Appendix 1. We give here a brief description of the framework of the proof. Because a PAC candidate can be tested in polynomial time, PAC-DETECTION is in NP. It remains to be shown that it is NP-hard. To this end, we reduce from the well-known NP-complete problem SAT [21]. An instance of SAT asks whether an input formula on  $n$  boolean variables  $x_1, \dots, x_n$  is satisfiable, via a suitable assignation of variables with true/false values.

To perform the reduction, we associate to each such formula  $\varphi$  a reaction system  $S_\varphi$ , of size polynomial in  $\varphi$ , with a specified target entity  $A$  and a food set  $F$ . Reactions in  $S_\varphi$  are designed to mirror the structure of  $\varphi$ , ensuring that a PAC of the wanted form exists if and only if the formula is satisfiable. The only possible such PACs will actually directly encode satisfying assignments for  $\varphi$ .

This shows that PAC-DETECTION is NP-hard: a polynomial-time algorithm for PAC-DETECTION would yield a polynomial-time algorithm for SAT, via this reduction. We can conclude that PAC-DETECTION is NP-complete, since it is also in NP.

It would be interesting to investigate whether an unconstrained version of the PAC detection problem, i.e. without restricting the allowed foods, is also NP-complete. We leave this problem open. Since we will be interested in PAC enumeration, we must in any case be able to solve the constrained version.

#### B. Implementation

The above NP-completeness result justifies the use of an SMT Solver such as Z3 to enumerate all minimal PACs. We thus implemented this approach in C++ in the EmergeNS software [22] (more generally designed to simulate the dynamics of complex physicochemical systems and down the line to trace the physical emergence of natural selection). In practice, in a reaction system defined in EmergeNS, we ask the Z3 solver to find PAC candidates and to assess, for each candidate, the existence of a PAC witness, i.e. a reaction flow vector  $\vec{v}$  yielding only positive rates for entities of the candidate. Following our definition of PACs as minimal, i.e. equivalent to cores from [13], PAC candidates must use all of their entities exactly once as reactants. We then exclude the list of already found PACs from the search space, and repeat the process until no new PAC is found.

As remarked in Section III A, verifying whether a given candidate is indeed a PAC (by assessing the existence of a PAC witness) is a linear programming problem, and can thus be achieved efficiently, i.e. in polynomial time. The need for Z3 comes from the search for PAC candidates in an exponential space of possible subsets of entities and reactions.

### IV. PAC CONSISTENCY UNDER THERMODYNAMIC CONSTRAINTS

The kinetics of a reaction network must obey the second law of thermodynamics, a constraint that is not considered in the PAC definition. Indeed, this definition solely relies on the existence of a witness  $\vec{v}$  of reaction flows, that may or may not be compatible with thermodynamic constraints.

More precisely, once association and dissociation constants are derived from free energies and activation barriers, they cannot be freely chosen. Let us give more details on this link between energies and reaction rates. First, recall that each entity  $e$  is associated with a chemical potential  $\mu$ , which is the sum of their molar Gibbs free energy (or standard potential)  $G$ , and their activity (that we identify here to their

concentration  $[e]$ , hence placing ourselves in the ideal solution regime)<sup>1</sup>:

$$\mu = G + \ln [e] \quad (2)$$

In addition, the flow of a reaction  $R$ , denoted  $v_R$ , depends on the concentrations (i.e. activities) of its entities, under mass action kinetics:

$$v_R = k_R^{(+)} \cdot \prod_{e_i \text{ reactant}} [e_i]^{-\sigma_{e_i}^R} - k_R^{(-)} \cdot \prod_{e_j \text{ product}} [e_j]^{\sigma_{e_j}^R}, \quad (3)$$

where  $k_R^{(+)}$  and  $k_R^{(-)}$  are respectively the forward and backward kinetic rate constants and  $\sigma_{e_i}^R$  is the  $e_i$  stoichiometry in reaction  $R$ . The key point, to enforce the second law, is to relate the kinetic rate constants to Gibbs free energies using the local detailed balance condition (also known as Eyring's formula):

$$\begin{aligned} k_R^{(+)} &= \exp(-G_R^\ddagger + \sum_{e_i \text{ reactant}} (-\sigma_{e_i}^R \cdot G_i)) \\ k_R^{(-)} &= \exp(-G_R^\ddagger + \sum_{e_j \text{ product}} \sigma_{e_j}^R \cdot G_j) \end{aligned} \quad (4)$$

Notice that we have also introduced here an intermediate state energy  $G_R^\ddagger$  resulting in an activation barrier.

Injecting Equation (4) in Equation (3) leads to the following formulation that will be further used below:

$$v_R = e^{-G_R^\ddagger} \cdot \left( \prod_{e_i \text{ reactant}} e^{-\sigma_{e_i}^R \cdot \mu_i} - \prod_{e_j \text{ product}} e^{\sigma_{e_j}^R \cdot \mu_j} \right) \quad (5)$$

The question of addressing the thermodynamic consistency of PACs can thus be reformulated as follows: can a PAC flow witness  $\vec{v}$  be realized through a concentration vector  $\vec{c}$  of all entities of the system? If such a vector  $\vec{c}$  exists, the PAC will be considered as a thermodynamically Consistent Autocatalytic Cycle (CAC). The concentration vector  $\vec{c}$  will then be called a *CAC witness*. Here we reach the second key result of the present study, namely, that one can always find such a CAC witness – i.e. that a single PAC is always thermodynamically consistent – as long as the concentration space is unbounded.

**Theorem 2.** *Let  $\vec{v}$  be a PAC witness (or any flow vector compatible the directions of the PAC reactions). Then there exists  $\lambda > 0$  and a concentration vector  $\vec{c}$  such that  $\lambda \vec{v}$  is the flow vector induced by  $\vec{c}$ .*

The general proof of this theorem is given in Appendix 2, with an example provided in Appendix 2a. It should be noted that Theorem 2 actually has a broader scope than the problem specifically addressed in this section, since it demonstrates that *any* flow vector  $\vec{v}$  matching directions of the PAC

reactions can be realized in the concentration space up to some proportionality factor  $\lambda > 0$ . Notably, this holds for any values of activation barriers and Gibbs free energies, and even when food and waste concentrations are fixed to any arbitrary values.

This allows us to deduce the following corollary:

**Corollary 1.** *Any PAC is a CAC.*

*Proof of Corollary 1.* Consider a PAC formed of entities  $(e_1, \dots, e_n)$  and reactions  $(R_1, \dots, R_n)$  (recall that according to Blokhuis et al [13], a minimal PAC contains as many entities as reactions). Let  $\vec{v} = (v_1, \dots, v_n)$  be a PAC witness. Since the inequalities to be satisfied by a PAC witness are all linear with respect to the coordinates of  $\vec{v}$  (see Equation (1)), for all  $\lambda > 0$ ,  $\lambda \vec{v}$  is a PAC witness as well. Applying Theorem 2 gives us the existence of some  $\lambda > 0$  and a concentration vector  $\vec{c}$  yielding flows  $\lambda \vec{v}$ , which is a PAC witness. Thus  $\vec{c}$  is a CAC witness, and the arbitrary minimal PAC we started with is indeed a CAC.  $\square$

## V. COMPATIBILITY AMONG AUTOCATALYTIC CYCLES

In this section, we investigate whether thermodynamic constraints affect compatibility relationships among cycles. To do so, we first analyze compatibility among PACs, that is, we identify sets of cycles that are found compatible on the basis of the reaction network topology alone, hereafter called “multiPACs”. A set of PACs is a multiPAC if there exists a common witness  $\vec{v}$  of reaction flows allowing all the PACs of the set to run simultaneously. In the framework of Gagrani et al [14], this corresponds to a nonempty intersection of the flow-productive cones for the different autocatalytic cores considered. Here, instead of computing explicit intersections, we will use the Z3 solver to identify nonempty intersections, and directly ask for a single vector  $\vec{v}$  witnessing the different PACs simultaneously. To do so, we simply concatenate the requirements already defined for each individual PAC.

Interestingly, we note that incompatibility between PACs may occur for various reasons. The simplest case, illustrated in Figure 2, is when a reaction is shared between two PACs, but in opposite directions. This obviously prevents the existence of a common flow vector witness  $\vec{v}$ . Yet more subtle cases were also obtained from our software using randomly generated reaction networks. One example is described in Figure 3. Generally speaking, incompatibilities among PACs occur because of contradictory requirements in terms of flows, that is, when one PAC requires a reaction  $R_1$  to run faster than a reaction  $R_2$ , while a second PAC requires the opposite.

To investigate the impact of thermodynamic constraints on compatibility relationships among cycles, we now assess whether pairs of compatible PACs, witnessed by a common flow vector  $\vec{v}$ , also constitute pairs of compatible CACs, witnessed by a common concentration vector  $\vec{c}$ . As before, addressing this question using the SMT solver is straightforward: it suffices to concatenate the lists of constraints required for each of the CACs under study, and to ask if a  $\vec{c}$  vector exists that simultaneously satisfies them all.

This analysis reveals that compatibility among cycles inferred solely from the reaction network topology may overlook

<sup>1</sup> In this formulation, we set  $RT = 1$ , which can be done without loss of generality, as this is simply equivalent to expressing Gibbs free energies in  $RT$  units. In the following, the notation  $R$  will stand for a reaction, and the gas constant will not be referred to again



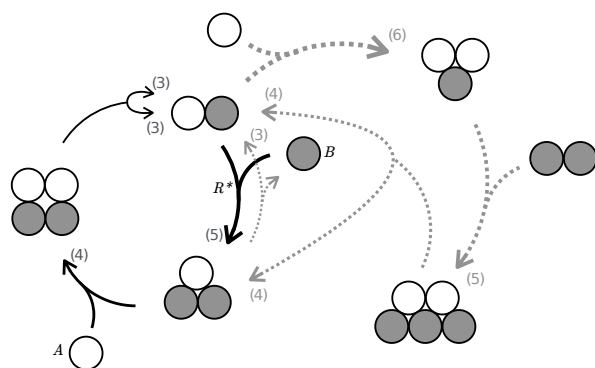


Figure 2. A simple example of two incompatible PACs. The PAC with full arrows is  $AB \rightarrow AB_2 \rightarrow A_2B_2 \rightarrow AB + AB$  and the gray dotted arrows PAC is  $AB_2 \rightarrow AB \rightarrow A_2B \rightarrow A_2B_3 \rightarrow AB + AB_2$ . Numbers in brackets indicate the flows of PAC reactions allowing for a (local) net production of 1 of all their entities. The two PACs share the reaction labeled  $R^*$ , but require it to run in opposite directions.

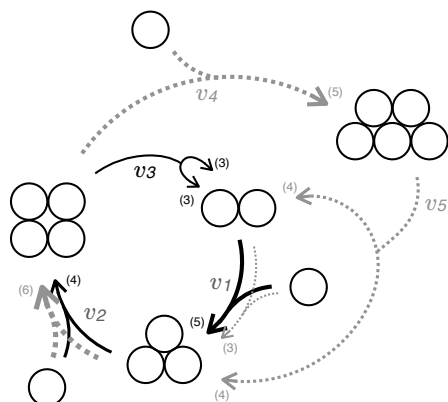


Figure 3. A more subtle example of two incompatible PACs, sharing two reactions indexed with flows  $v_1$  and  $v_2$ . Flows allowing for unitary production of PAC entities are shown in brackets for both PACs. The inner PAC depicted with solid black arrows ( $A_2 \rightarrow A_3 \rightarrow A_4 \rightarrow A_2 + A_2$ ) imposes the flow inequalities  $v_1 > v_2 > v_3 > v_1/2$ , while the gray dotted outer PAC ( $A_2 \rightarrow A_3 \rightarrow A_4 \rightarrow A_5 \rightarrow A_2 + A_3$ ) imposes inequalities  $v_1 + v_5 > v_2 > v_4 > v_5 > v_1$ . This yields contradictory requirements:  $v_1 > v_2$  for the first PAC and  $v_1 < v_2$  for the second.

thermodynamic inconsistencies. Indeed, we find several instances of two compatible PACs making incompatibles CACs. We give an example of such a behavior in Figure 4 and provide a formal proof of the CAC incompatibility in Box 2.

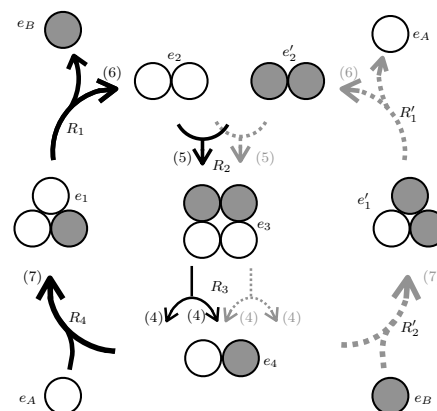


Figure 4. Schematic view of two autocatalytic cycles ( $\{R_1, R_2, R_3, R_4\}$  and  $\{R'_1, R'_2, R'_3, R'_4\}$ ) sharing two reactions. One can check that both cycles run simultaneously with flows  $v_1 = v'_1 = 6$ ,  $v_2 = 5$ ,  $v_3 = 4$  and  $v_4 = v'_4 = 7$  which allows for a unity production of all entities belonging to a PAC. However, it can be shown that these cycles can not be instantiated in the concentration space, i.e. they form a multiPAC but not a multiCAC.

## Box 2: Thermodynamic incompatibility between CACs

On the example shown in Figure 4, we can first show that the two PACs are compatible, since we can choose a set of reaction flows that witnesses both of them simultaneously. For instance setting  $v_1 = v'_1 = 6$ ,  $v_2 = 5$ ,  $v_3 = 4$ ,  $v_4 = v'_4 = 7$  leads to a positive production rate of each entity. For the two PACs to run simultaneously, the following inequalities must be satisfied:  $v_3 < v_2 < v_1 < v_4 < 2v_3$  and  $v_3 < v_2 < v'_1 < v'_4 < 2v_3$ .

We will show that these inequalities are not thermodynamically achievable, i.e. that they lead to contradictory requirements in the concentration space.

Notice that the inequalities imply that all  $v_i$  are strictly positive, because  $v_3 < 2v_3$  entails  $v_3 > 0$ , and all other  $v_i$  are larger than  $v_3$ . As proven below, this sign constraint alone is not satisfiable: not all reactions can flow in the wanted direction.

We denote  $x_i = e^{\mu_i}$  the exponential of the chemical potential of entity  $e_i$  (indexed as in Fig. 4), and  $b_i = e^{G_i^\ddagger}$  where  $G_i^\ddagger$  is with respect to reaction  $R_i$  (similarly for  $R'_i$ ). Reaction flows can thus be written as follows:

$$\begin{aligned} v_1 &= b_1(x_1 - x_B x_2) & b_4 &= b_4(x_A x_4 - x_1) \\ v'_1 &= b'_1(x'_1 - x_A x'_2) & v'_4 &= b'_4(x_B x_4 - x'_1) \\ v_3 &= b_3(x_3 - (x_4)^2) & v_2 &= b_2(x_2 x'_2 - x_3) \end{aligned}$$

From the positivity of all flows, we get:

$$\begin{cases} x_B x_2 < x_1 < x_A x_4 \\ x_A x'_2 < x'_1 < x_B x_4 \\ (x_4)^2 < x_3 < x_2 x'_2 \end{cases}$$

From the first line we deduce  $x_B x_2 / x_4 < x_A$ , and re-injecting in the second line we obtain  $x_B x_2 x'_2 / x_4 < x_A x'_2 < x_B x_4$ . This simplifies to  $x_2 x'_2 < (x_4)^2$ , which contradicts the condition of the third line. It should be noted that this proof is valid regardless of the activation barriers values, since the contradiction stems from the signs of the flows, while activation barriers only affect their amplitudes.

## VI. DISCUSSION

Working toward the long-term goal of an explicit grounding of Darwinian dynamics into physical processes, we addressed in this study the implications of thermodynamic constraints on the existence and detection of autocatalytic cycles given a reaction network. Our analysis builds on recent theoretical progresses made on the formalization of autocatalysis on the sole basis of the reaction network topology [13]. Under this definition, the exhaustive detection of autocatalysis proved here to be an NP-complete problem. This finding fully justifies the use of constraint solvers (e.g. SMT, Integer Programming) toward which we converge with others [8, 14].

We found that constraints imposed by free energies and activation barriers can always be compensated by adjusting concentrations, thereby allowing any minimal autocatalytic cycle to also be thermodynamically consistent. In other words, the list of autocatalytic cycles in a reaction network remains unaffected by these physical constraints, as long as concentrations are not limited by upper or lower bounds. However, as shown in Appendix 3, it should be noted that heterogeneity in free energies and activation barriers do restrict the volume of the concentration space where a cycle runs.

These conclusions on isolated cycles do not readily apply on combinations of cycles. Indeed, thermodynamic realism does restrict the list of mutually compatible cycles, even in an unlimited concentration space, such that topologically com-

patible cycles can turn out incompatible. Incompatibilities between two autocatalytic cycles can therefore stem from two distinct sources, namely the topology of the reaction network (PAC-incompatibility) and irreconcilable demands on concentrations (CAC-incompatibility).

A stimulating next step will be to investigate the implications of autocatalysis on the system’s dynamics through time and space. Among the many autocatalytic cycles that are thermodynamically achievable in a given system, which ones are actually encountered from a given starting point in the concentration space? Which ones are running in the long term, that is, once the system has reached a steady state? Which ones are running in the vicinity of this steady state, and perhaps contribute to drive the system in its direction? And finally, could autocatalysis contribute to generate more than one steady state in the concentration space? We anticipate that such multistability could enable a primordial form of heritable variation, paving the way to nascent Darwinian dynamics.

## ACKNOWLEDGMENTS

We are very grateful to Nicolas Lartillot for his insightful contribution in setting up our modeling approach, and to Benjamin Kuperberg for developing his open source GUI library OrganicUI and helping us use it. We also thank Olivier Rivoire and Yann Sakref for fruitful early discussions.

## APPENDIX

### 1. Proof of Theorem 1: NP-completeness

This section is devoted to the proof of Theorem 1.

First, by the above remark, if a candidate PAC is given then it can be checked in polynomial time; thus PAC-DETECTION is in NP. To achieve the proof, it remains to be shown that it is NP-hard. We do so by reducing from the classical NP-complete problem SAT, that is, by translating the SAT problem into PAC-DETECTION.

- [1] N. Goldenfeld and C. Woese, Annual Review of Condensed Matter Physics **2**, 375 (2011).
- [2] S. Charlat, T. Heams, and O. Rivoire, in *Evolutionary Thinking Across Disciplines*, Vol. 478, edited by A. Du Crest, M. Valković, A. Ariew, H. Desmond, P. Huneman, and T. A. C. Reydon (Springer International Publishing, Cham, 2023) pp. 287–296, series Title: Synthese Library.
- [3] S. Charlat, A. Ariew, P. Bourrat, M. Ferreira Ruiz, T. Heams, P. Huneman, S. Krishna, M. Lachmann, N. Lartillot, L. Le Sergeant d’Hendecourt, C. Malaterre, P. Nghe, E. Rajon, O. Rivoire, M. Smerlak, and Z. Zeravcic, Life **11**, 1051 (2021).
- [4] D. A. Baum, Z. Peng, E. Dolson, E. Smith, A. M. Plum, and P. Gagrani, Journal of The Royal Society Interface **20**, 20230346 (2023).
- [5] A. Sharma, D. Czégel, M. Lachmann, C. P. Kempes, S. I. Walker, and L. Cronin, Nature **622**, 321 (2023).
- [6] L. L. J. Schoenmakers, T. A. C. Reydon, and A. Kirschning, Life **14**, 175 (2024).
- [7] M. Kalambokidis and M. Travisano, Evolution **78**, 1 (2024).
- [8] Z. Peng, J. Linderth, and D. A. Baum, PLOS Computational Biology **18**, e1010498 (2022).
- [9] Y. Liu and D. J. Sumpter, Journal of Biological Chemistry **293**, 18854 (2018).
- [10] S. Sarkar and J. L. England, Physical Review E **100**, 022414 (2019).
- [11] M. Eigen and P. Schuster, Die Naturwissenschaften **64**, 541 (1977).
- [12] D. A. Baum, BioScience **65**, 678 (2015).
- [13] A. Blokhuis, D. Lacoste, and P. Nghe, Proceedings of the National Academy of Sciences **117**, 25230 (2020).
- [14] P. Gagrani, V. Blanco, E. Smith, and D. Baum, Journal of Mathematical Chemistry **62**, 1012 (2024).
- [15] J. L. Andersen, C. Flamm, D. Merkle, and P. F. Stadler, Journal of Systems Chemistry **3**, 1 (2012).
- [16] O. Weller-Davies, M. Steel, and J. Hein, Mathematical Biosciences **325**, 108365 (2020).
- [17] W. Hordijk and M. Steel, Biosystems **152**, 1 (2017).
- [18] S. A. Kauffman, Journal of Theoretical Biology **119**, 1 (1986).
- [19] A. Despons, Y. De Decker, and D. Lacoste, Communications Physics **7**, 224 (2024).
- [20] A. Arya, J. Ray, S. Sharma, R. Cruz Simbron, A. Lozano, H. B. Smith, J. L. Andersen, H. Chen, M. Meringer, and H. J. Cleaves, Chemical Science **13**, 4838 (2022).
- [21] S. A. Cook, in *Proceedings of the Third Annual ACM Symposium on Theory of Computing*, STOC ’71 (Association for Computing Machinery, New York, NY, USA, 1971) p. 151–158.
- [22] D. Kuperberg, Emergens, <https://github.com/deniskup/EmergeNS> (2024).

An instance of SAT is a boolean formula  $\varphi$  on  $n$  variables  $x_1, x_2, \dots, x_n$ . We will call literal a variable  $x_i$  or its negation  $\bar{x}_i$ . The formula  $\varphi$  is a conjunction of  $k$  clauses, i.e. of the form  $\varphi = \bigwedge_{1 \leq j \leq k} c_j$ , where each clause  $c_j$  is a disjunction of literals. The formula  $\varphi$  is satisfiable if there is a valuation setting a boolean value for each  $x_i$  (and the opposite for  $\bar{x}_i$ ), that allows  $\varphi$  to evaluate to *true*. The SAT problem, that is, asking whether an input formula  $\varphi$  is satisfiable, is known to be NP-complete [21].

Let us now encode SAT into PAC-DETECTION. Given an instance of SAT, i.e. a formula  $\varphi$  as above, we want to design a reaction system  $S_\varphi = (E, R)$ , together with a target entity  $A \in E$  and a food set  $F \subset E$ , such that there is a PAC satisfying this instance of PAC-DETECTION if and only if  $\varphi$  is satisfiable. This would mean that a (hypothetical) polynomial-time algorithm for PAC-DETECTION would yield a polynomial-time algorithm for SAT, provided the reduction can be done in polynomial time, and produces a reaction system of polynomial size, which will be ensured here.

We choose as entities the set:

$$E = \{A, A_2, N\} \cup \{E_i, \bar{E}_i \mid 1 \leq i \leq n\} \cup \{C_j \mid 1 \leq j \leq k\}.$$

This means we have  $3 + 2n + k$  entities in the system. Note that we use capital letters to distinguish entities from the variables and clauses from  $\varphi$  to which they refer.

For reactions, we choose the following set, where  $L_i$  always stands for either  $E_i$  or  $\bar{E}_i$ , e.g. the first line means that both reactions  $A \rightleftharpoons E_1 + N$  and  $A \rightleftharpoons \bar{E}_1 + N$  are present:

$$R : \begin{cases} A \rightleftharpoons L_1 + N \\ E_i + \bar{E}_i \rightleftharpoons L_{i+1} & \text{for } 1 \leq i < n \\ E_n + \bar{E}_n \rightleftharpoons C_1 \\ C_j + L \rightleftharpoons C_{j+1} & \text{for } 1 \leq j < k \\ & \text{and } L \text{ literal of } c_j \\ C_k + L \rightleftharpoons A_2 & \text{for } L \text{ literal of } c_n \\ A_2 \rightleftharpoons A + A \end{cases}$$

Finally, the target is  $A$  and the allowed foods will be  $F = \{E_i, \bar{E}_i \mid 1 \leq i \leq n\}$ .

The idea is that a valuation witnessing satisfiability of  $\varphi$  will correspond to a PAC of the form:

$$A \rightarrow L_1 \rightarrow \dots \rightarrow L_n \rightarrow C_1 \rightarrow \dots \rightarrow C_k \rightarrow A_2 \rightarrow A + A.$$

The crux of the construction is that reactions using  $C_j$  as reactant need as food one of the literals of  $c_j$ , that will be interpreted as the literal validating the clause  $c_j$ . In consequence, this literal must be a food, and cannot appear as one of the entities of the PAC. So the segment  $L_1 \rightarrow \dots \rightarrow L_n$  appearing in the PAC has to be formed exactly of the literals that are put to false in the valuation.

The entity  $N$  is used to force the direction of the reaction  $A \rightarrow L_1$  in the PAC: since it can only be used as waste and not as food, the reaction cannot go in the opposite direction. From there, any PAC containing  $A$  and using only foods from  $F$  must be of the above form, and witnesses a valuation satisfying  $\varphi$ .

**Example 2.** <sup>2</sup> Consider  $\varphi$  given by  $c_1 = x_1 \vee x_2$ ,  $c_2 = x_1 \vee \bar{x}_2$ , and  $c_3 = \bar{x}_1 \vee \bar{x}_2$ . Then the only correct valuation is the one setting  $x_1$  to true and  $x_2$  to false. This is witnessed in the system  $S_\varphi$  by the PAC  $A \rightarrow \bar{E}_1 \rightarrow E_2 \rightarrow C_1 \rightarrow C_2 \rightarrow C_3 \rightarrow A_2 \rightarrow A + A$ , using  $E_1$  and  $\bar{E}_2$  as foods (and  $N$  as waste).

If we add a clause  $c_4 = \bar{x}_1 \vee x_2$  to  $\varphi$ , the formula becomes unsatisfiable and accordingly, the corresponding system  $S_\varphi$  does not contain a PAC satisfying the constraints.

This achieves the proof that PAC-DETECTION is NP-complete.

## 2. Proof that all PACs are CACs

In the following, to simplify notations, reactions will be oriented in accordance to the PAC, i.e. all coordinates of the PAC witness will be positive. A “positive flow vector” is a vector where all coordinates are positive.

a. An example

**Example 3.** Let us first instantiate the proof with the PAC mentioned earlier, given by the following reactions (ignoring foods)  $2e_1 \rightarrow e_2 \rightarrow e_3 \rightarrow e_1 + e_2$ , depicted in Figure 5 ignoring food entities.

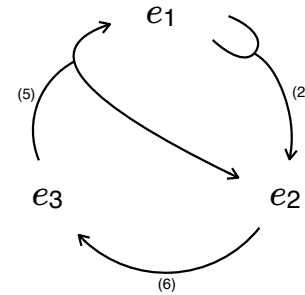


Figure 5. A PAC example of size 3 to give a first intuition of the proof.

A PAC witness is  $\vec{v} = (2, 6, 5)$ , ensuring a production flow of 1 for each entity. Using Equation (5), we now look for a concentration vector to realize this PAC witness:

- $v_1 = e^{-G_1^\dagger} (e^{2\mu_1} - e^{\mu_2})$ ,
- $v_2 = e^{-G_2^\dagger} (e^{\mu_2} - e^{\mu_3})$ ,
- $v_3 = e^{-G_3^\dagger} (e^{\mu_3} - e^{\mu_1 + \mu_2})$ .

<sup>2</sup> This is actually an instance of 2SAT which is a simpler problem, but it is just used here to illustrate the construction

We will instead aim at realizing the flow vector  $\lambda \vec{v}$  for some  $\lambda > 0$ . Letting  $x_i = e^{\mu_i}$  and  $w_i = v_i e^{G_i^\dagger}$  (so  $w_i > 0$ ), the equation system can be written conveniently:

- $x_1 = \sqrt{\lambda w_1 + x_2}$
- $x_2 = \lambda w_2 + x_3$
- $x_3 = \lambda w_3 + x_1 x_2$

Substituting  $x_2$  we obtain  $x_1 = \sqrt{\lambda w_1 + \lambda w_2 + x_3}$ , and finally substituting  $x_1$  as well we obtain:

$$x_3 = \lambda w_3 + (\sqrt{\lambda w_1 + \lambda w_2 + x_3})(\lambda w_2 + x_3).$$

Let us note  $g(\lambda, x_3)$  the above right-hand side expression, so that the equation becomes  $x_3 = g(\lambda, x_3)$ . Let us define  $h(\lambda, x_3) = g(\lambda, x_3) - x_3$ , so that we aim at  $h(\lambda, x_3) = 0$ . We will show that such a solution exists using the intermediate value theorem.

For  $x_3 = 0$  and any  $\lambda > 0$ , we have  $h(\lambda, 0) = \lambda w_3 + (\sqrt{\lambda w_1 + \lambda w_2})\lambda w_2 > 0$ . Let us choose any  $a \in (0, 1)$ . For  $x_3 = a$  and  $\lambda = 0$ , we have  $h(0, a) = a^{3/2} - a < 0$ . Because  $h$  is continuous there exists  $\varepsilon > 0$  such that for all  $\lambda \in (0, \varepsilon)$ ,  $h(\lambda, a) < 0$ . Let us choose  $\lambda \in (0, \varepsilon)$ , we know thanks to the intermediate value theorem (on  $x_3$  with this fixed value of  $\lambda$ ) that there exists  $x_3 \in (0, a)$  such that  $h(\lambda, x_3) = 0$ . From this  $x_3$  we can compute  $x_2 = \lambda w_2 + x_3$ , and  $x_1 = \sqrt{\lambda w_1 + x_2}$ . We have obtained a solution to our system by an appropriate choice of  $\lambda, x_1, x_2, x_3$ , thereby witnessing that the PAC is a CAC via Corollary 1.

#### b. General proof

We aim to prove Theorem 2 from Section IV by generalizing the proof exposed in the previous example to any PAC.

#### Notations.

Let  $\vec{v}$  be a positive flow vector, we want to show that there exists  $\lambda > 0$  and a concentration vector  $\vec{c}$  such that  $\lambda \vec{v}$  is the flow vector induced by  $\vec{c}$ .

We will even show that this can be attained for any fixed concentration values of food and waste entities.

Let  $\vec{v} = (v_1, \dots, v_n)$  be the target flow vector, where all coordinates are strictly positive. Here  $n$  is the size of the PAC, i.e. both the number of entities  $e_1, \dots, e_n$ , and reactions  $R_1, \dots, R_n$ , where for each  $i$ ,  $e_i$  is the sole reactant of  $R_i$  among entities of the PAC, see discussion in [13].

For an entity  $e_i$  ( $1 \leq i \leq n$ ) with Gibbs free energy  $G_i$ , recall that its chemical potential is given by  $\mu_i = G_i + \ln([e_i])$ . We will be interested here in the exponential of this potential, a variable that we denote  $x_i = e^{\mu_i}$ .

Notice that  $x_i$  can be freely adjusted to any strictly positive value by choosing the appropriate concentration  $[e_i]$ , so we can turn the problem into that of finding a solution vector  $\vec{x} = (x_1, \dots, x_n)$ . Analogs of components of  $\vec{x}$  associated to food and waste entities are fixed to 1 for simplicity. The proof can easily be adapted to any given values, but this will save us some notations, as food and waste entities can now be ignored in the computation of the reaction flows. One can

refer to Appendix 3 for an analysis of the impact of food and waste potentials on the feasibility of PACs.

Let us recall Equation (5) giving the flow of a reaction in terms of chemical potentials:

$$v_R = e^{-G_R^\dagger} \cdot \left( \prod_{e_i \text{ reactant}} e^{-\sigma_{e_i}^R \cdot \mu_i} - \prod_{e_j \text{ product}} e^{\sigma_{e_j}^R \cdot \mu_j} \right)$$

In our case, given that  $e_i$  is the sole reactant of  $R_i$  (apart from possible foods), this can be rewritten for any  $i \in [1, n]$ :

$$v_{R_i} = e^{-G_{R_i}^\dagger} (x_i^{\sigma_{i,i}} - \prod_j x_j^{\sigma_{i,j}})$$

Where  $\sigma_{i,i}$  is the stoichiometry of the reactant  $e_i$  in  $R_i$ , and  $\sigma_{i,j}$  is the stoichiometry of product  $e_j$  in  $R_i$ . It should be noted here that for simplicity we change the convention regarding the stoichiometry of reactants, so that all parameters are positive, including  $\sigma_{i,i}$ .

Let us note  $w_i = v_i e^{G_{R_i}^\dagger}$ . The fact that we aim for  $\lambda \vec{v}$  to be the flow vector induced by  $\vec{x}$  means that for all  $i \in [1, n]$ , we aim for  $v_{R_i} = \lambda v_i$ , so we must have  $\lambda w_i = x_i^{\sigma_{i,i}} - \prod_j x_j^{\sigma_{i,j}}$ , i.e.

$$x_i = \left( \lambda w_i + \prod_j x_j^{\sigma_{i,j}} \right)^{1/\sigma_{i,i}} \quad \text{Equation (i)}$$

Notice that  $\sigma_{i,i}$  corresponds to “inverted forks”, where several  $e_i$  must be used as reactants. This can indeed occur in a PAC, as witnessed by Example 3.

#### Variable elimination.

Our goal is to find  $\lambda > 0$  such that this system has a solution  $(x_1, \dots, x_n)$  with each  $x_i > 0$ .

We will do so by iteratively reducing the number of variables and equations. Firstly, we perform the following operation: as long as there exists an equation of the form  $x_i = A_i$ , where  $A_i$  is an expression not depending on  $x_i$ , we replace  $x_i$  by  $A_i$  in all other equations, and remove equation (i).

We will end up with a list of equations of the form  $x_i = g_i(\lambda, x_1, \dots, x_n)$  where  $i$  ranges over some subset of  $[1, n]$ , and where  $g_i$  is an expression depending on some of its variables including  $x_i$ , and generated by the following grammar:

$$g := x_j \mid gg \mid (\lambda w_j + g)^{1/k}$$

where  $j$  ranges over  $[1, n]$ , and  $k$  is a strictly positive integer. There are as many remaining equations as remaining variables. Without loss of generality, we assume that the remaining variables are  $x_1, \dots, x_p$ , where  $p \leq n$  is the number of remaining equations.

#### System rewriting as graph transformations.

Before tackling the resolution of the system per se, we relate some properties of the remaining equations to the topology of the PAC we started with. We associate to the PAC a graph  $\mathcal{G}$ , whose vertices are entities  $e_1, \dots, e_n$ , and there is an edge  $e_i \rightarrow e_j$  if  $e_j$  is among the products of the reaction using  $e_i$  as reactant.

Substituting variable  $x_i$  in the system according to Equation (i) amounts to removing vertex  $e_i$  and contracting all



edges originating in this vertex, as shown in Figure 6 with  $i = 5$ :

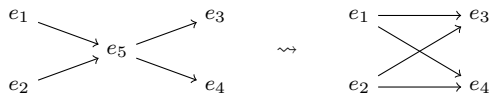


Figure 6. The effect of substituting the equation  $x_5 = \lambda w_5 + x_3 x_4$  on the graph  $\mathcal{G}$ .

This operation can only be performed if there is no self-loop on  $e_i$  in the current graph. Therefore, the case where we can end up with only one equation on an entity  $x_i$  (i.e. the case  $p = 1$ ) corresponds to the fact that the graph  $\mathcal{G} \setminus \{x_1\}$  is acyclic.

We give in Figure 7 an example where this does not happen, no matter in which order the substitutions are done. This example corresponds to the Type V pattern from [13].

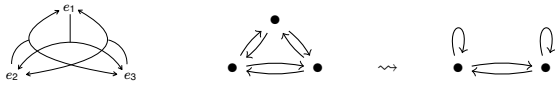


Figure 7. A 3-entity PAC, and its graph  $\mathcal{G}$ . After performing one substitution step, all nodes have self-loops and the substitution sequence must end, with equations of the form e.g.  $x_1 = \lambda w_1 + (\lambda w_3 + x_1 x_2) x_2$ .

**Shape of the reduced system.** Let us now turn to the following lemma, which gives some properties of the expressions  $g_i$  reached by this process.

**Lemma 1.** For all  $1 \leq i \leq p$ ,

- $g_i$  is always of the form  $(\lambda w_i + g)^{1/k_i}$ ,
- for  $\lambda = 0$ ,  $g_i = \prod_j x_j^{m_{i,j}}$  with  $m_{i,j} \geq 0$  for all  $j$ , and
  - if only one equation remains ( $p = 1$ ),  $m_{1,1} > 1$ ,
  - otherwise if  $p > 1$ , then for all  $i \in [1, p]$ , we have  $m_{i,i} = 1$  and  $m_{i,j} > 0$  for some  $j \neq i$ .

*Proof.* The first item is a direct consequence of the operation we performed, starting from Equation (i).

The general shape of the second item is guaranteed by construction.

We now show item 2.1. Let us first give an interpretation for exponents  $m_{i,j}$ . Consider Equation (i) when  $\lambda = 0$ :

$$x_i = \prod_j x_j^{\sigma_{i,j}/\sigma_{i,i}}.$$

Recall that we have defined  $\sigma_{i,i} > 0$  even though entity  $e_i$  is a reactant. Exponents in the equation system for  $\lambda = 0$  simply express the fact that in reaction  $R_i$  and starting from one unit of entity  $e_i$ , for all  $j$  one can produce  $\sigma_{i,j}/\sigma_{i,i}$  units of entity  $e_j$ . The substitution algorithm is equivalent to consuming all available  $e_j$  via reaction  $R_j$ , with exponents of the products keeping track of their quantity. By iteration, if one started with entity  $e_i$  and could remove all reactions from the system

(in our notation that would mean  $p = 1$ ), one should recover that  $e_i$  can self-amplify through the reactions of the cycle since it allows for autocatalysis of  $e_i$ . This intuitively explains why we expect  $m_{1,1} > 1$ . For illustration in Example 3, the cycle allows to produce  $3/2$  units of entity  $e_3$  starting from one unit of it.

Let us now give a formal proof of the fact that  $m_{1,1} > 1$  when  $p = 1$ , and let us note  $m = m_{1,1}$  for concision. The sequence of substitutions leading to the expression  $g_1$  amounts to describing a positive flow vector  $\vec{\tau} = (t_1, \dots, t_n)$  in the following way: the flow of reaction  $R_1$  is set to  $t_1 = 1$ , and then all available entities  $e_j$  other than  $x_1$  are entirely consumed and substituted with the products of reaction  $R_j$ , until only  $m$  units of  $e_1$  remain. This means that the net production of each entity is given by  $M\vec{\tau} = (m - 1, 0, 0, \dots, 0)$ , where  $M$  is the reaction matrix associated to the PAC. Let us now consider a PAC witness  $\vec{v}' = (v'_1, \dots, v'_n)$ , with  $v'_i > 0$  for each  $i$ , and scaled such that  $v'_1 = 1$  (this is always possible since scaling preserves the PAC witness inequalities). We will show that the existence of  $\vec{v}'$  implies  $m > 1$ . Notice that  $M \cdot \vec{\tau} = (m - 1, 0, \dots, 0)$ . Our goal is to show that for all  $i \in [2, n]$ ,  $v'_i < t_i$ .

We use the graph  $\mathcal{G}$  defined in the previous paragraph, and more precisely the fact that the case  $p = 1$  corresponds to acyclicity of  $\mathcal{G}' = \mathcal{G} \setminus \{e_1\}$ . Without loss of generality, we can assume that a topological order of  $\mathcal{G}'$  is given by  $e_2, e_3, \dots, e_n$ , that is to say there is no edge  $e_i \rightarrow e_j$  with  $j < i$ . In particular,  $e_2$  has no incoming edge in  $\mathcal{G}'$ . We compare the inflow and outflow of  $e_2$  according to flow vectors  $\vec{\tau}$  and  $\vec{v}'$ .

- the inflow is the same in both cases, as  $t_1 = v'_1 = 1$  and no other incoming edge exists,
- the outflows are  $t_2$  and  $v'_2$  respectively,
- the net production rate with respect to  $\vec{\tau}$  is 0, by construction, while it is strictly positive for  $\vec{v}'$  as  $\vec{v}'$  is a PAC witness.

These three items together imply that  $v'_2 < t_2$ . We can continue by induction with  $e_3, e_4, \dots$ : at step  $i$  for entity  $e_i$ , the inflow is smaller or equal with flows from  $\vec{v}'$  than from  $\vec{\tau}$  (using induction hypothesis  $v'_j < t_j$  for  $j < i$ ), but the net production is positive for  $\vec{v}'$  while it is 0 for  $\vec{\tau}$ . We can conclude the inequalities on outflows:  $v'_i < t_i$ . This allows us to conclude that for all  $i \in [2, n]$ ,  $v'_i < t_i$ . Let us call  $m'$  the inflow of  $e_1$  according to  $\vec{v}'$ , and recall that the corresponding inflow with respect to  $\vec{\tau}$  is  $m$ . Since every reaction producing  $e_1$  is strictly smaller in  $\vec{v}'$  than in  $\vec{\tau}$ , we have  $m' < m$ . However, to ensure net production of  $e_1$  according to  $\vec{v}'$ , and since the outflow is  $v_1 = 1$ , we must have  $1 < m'$ . We can conclude  $m > 1$  as desired.

Now if more than one equations remain, i.e.  $p > 1$  (item 2.2 of the Lemma) and  $m_{i,i} \neq 1$ , it means that one found a subset of reactions of the cycle allowing for the net production (or consumption) of entity  $e_i$ , in contradiction with the assumption that the cycle is minimal. An illustration of this phenomenon is given in Example 4. The existence of some  $j \neq i$  such that  $m_{i,j} > 0$  follows from the condition that at least one other equation than the  $i^{th}$  remains after the substitution algorithm ends. Indeed, if we could reach an equation

where  $x_i$  is the only remaining variable, without having substituting all the other ones, it would mean that the PAC is not strongly connected, contradicting its minimality.  $\square$

**Example 4.** We illustrate a case with  $m_{i,i} \neq 1$  while two equations remain in Figure 8.A. The cycle is not minimal, even though all entities appear only once as reactant, because it contains the cycle  $e_2 \rightarrow e_3 \rightarrow e_1 \rightarrow e_1 + e_2$  (depicted in black arrows). Notice that the forked reaction  $e_2 \rightarrow e_3 + e_4$  is not entirely in full line, because for this embedded minimal PAC,  $e_4$  corresponds to a waste. The equations corresponding to the full system are:

$$\begin{cases} x_1 = \lambda w_1 + x_2 x_3 \\ x_2 = \lambda w_2 + x_3 x_4 \\ x_3 = \lambda w_3 + x_1 \\ x_4 = \lambda w_4 + x_5 \\ x_5 = \lambda w_5 + x_3 x_4. \end{cases}$$

One can substitute  $x_3$ , then  $x_4$  and  $x_2$  to arrive to a system of two equations on  $x_1$  and  $x_5$ :

$$\begin{cases} x_1 = \lambda w_1 + [\lambda w_2 + (\lambda w_3 + x_1)(\lambda w_4 + x_5)](\lambda w_3 + x_1) \\ x_5 = \lambda w_5 + (\lambda w_3 + x_1)(\lambda w_4 + x_5). \end{cases}$$

Taking  $\lambda = 0$ , one obtains  $x_1 = x_1^2 x_5$  and  $x_5 = x_1 x_5$ , thus we have  $m_{1,1} = 2$  and  $m_{5,5} = 1$ . The appearance of  $x_1^2$  in the first equation betrays the existence of a smaller PAC where entity  $e_5$  is ignored. A visual interpretation is given in Figure 8.B, which shows the reduced cycle with the two entities  $e_1$  and  $e_5$  obtained at the end of the sequence of substitutions. Notice that entity  $e_5$ , in order to self-amplify, must produce  $e_1$ . However there exists a reaction path not involving entity  $e_5$  allowing for the production of 2 units of  $e_1$  starting with 1 unit of it (full black arrows in panel B), which is why  $m_{1,1} = 2$ . This indeed corresponds to the minimal PAC shown in full black arrows in panel A.

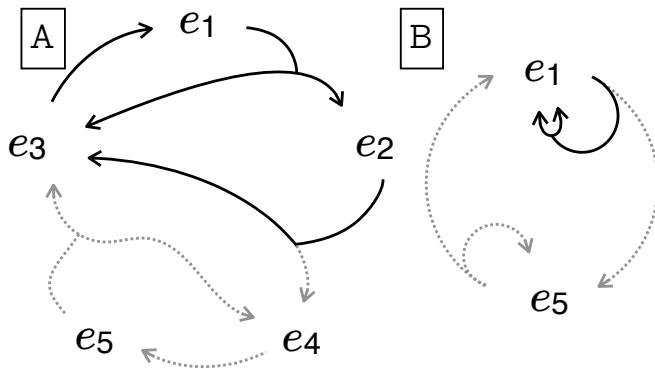


Figure 8. Non-minimal autocatalytic cycle.

**Solving the reduced system.** Turning back to the general proof of Theorem 2, it suffices to solve the remaining system

of equations, which is of the form  $x_i = g_i(\lambda, x_1, \dots, x_p)$  for all  $i \in [1, p]$ . We can then recover the missing  $x_j$  (for  $j \in [p+1, n]$ ) using the previous substitutions  $x_j = A_j$ .

**Case  $p = 1$ .** Let us first treat the easier case where  $p = 1$ , i.e. we have a single variable  $x_1$  and a single equation  $x_1 = g_1(\lambda, x_1)$ .

**Lemma 2.** There exists  $\varepsilon > 0$  such that for all  $\lambda \in (0, \varepsilon)$ , there exists  $x_1 \in (0, 1)$  verifying  $x_1 = g_1(\lambda, x_1)$ .

*Proof.* Let  $h_1(\lambda, x_1) = g_1(\lambda, x_1) - x_1$ . We will use the fact that  $h_1$  is continuous, and look at the values of  $h_1$  for a suitable  $\lambda$ ,  $x_1 = 0$  and  $x_1 = a$  for some  $a \in (0, 1)$  to conclude via the intermediate value theorem.

First of all for any  $\lambda > 0$  and  $x_1 = 0$ , we have  $h_1(\lambda, 0) > 0$ , since by construction  $g_1$  is built from positive terms using functions preserving positivity, and furthermore contains a term of the form  $\lambda w_1$  not depending on  $x_1$ .

Now, let us consider the case  $x_1 = a$  for any  $a < 1$ . We have  $h_1(\lambda, a) = g_1(\lambda, a) - a$ . By Lemma 1, for  $\lambda = 0$  we have  $g_1(0, a) = a^m$  with  $m > 1$ . Since  $a < 1$  and  $m > 1$ , we have  $a^m < a$ . Thus,  $h_1(0, a) < 0$ . Since the function  $\lambda \mapsto h_1(\lambda, a)$  is continuous, there exists  $\varepsilon > 0$  such that for all  $\lambda \in (0, \varepsilon)$ , we have  $h_1(\lambda, a) < 0$ . We can conclude (intermediate value theorem with respect to  $x_1$ ) that for all  $\lambda \in (0, \varepsilon)$ , there exists  $x_1 \in (0, a)$  such that  $h_1(\lambda, x_1) = 0$ .  $\square$

Therefore, choosing any  $\lambda < \varepsilon$  and applying this lemma gives a solution to the system.

**Case  $p > 1$ .** If  $p > 1$ , we will solve the system by constructing a sequence of partial continuous functions  $f_i : \mathbb{R}^i \rightarrow \mathbb{R}$  ( $1 \leq i \leq p$ ) and  $\varepsilon > 0$  such that:

- if  $x_1, \dots, x_{i-1} \in (0, 1)$  (resp.  $[0, 1]$ ) and  $\lambda \in (0, \varepsilon)$  (resp.  $[0, \varepsilon]$ ), then  $f_i(\lambda, x_1, \dots, x_{i-1})$  is defined and its image lies in  $(0, 1)$  (resp.  $[0, 1]$ ),
- if there exists  $\lambda, x_1, \dots, x_p$  such that for all  $i$ ,  $x_i = f_i(\lambda, x_1, \dots, x_{i-1})$ , then we have a solution of the wanted equations.

We define  $f_i$  by induction, starting with  $f_p$ .

Let us consider the equation  $x_p = g_p(\lambda, x_1, \dots, x_p)$ . Let  $h_p(\lambda, x_1, \dots, x_p) = g_p(\lambda, x_1, \dots, x_p) - x_p$ .

Our goal is to define  $f_p$  as a value of  $x_p$  yielding  $h_p = 0$ . Thus we have to show that such a root of  $h_p$  exists.

**Lemma 3.** For any  $x_1, \dots, x_{p-1} \in (0, 1)$ , there exists  $\varepsilon_p$  such that for all  $\lambda \in (0, \varepsilon_p)$ , there exists  $x_p \in (0, 1)$  verifying  $h_p(\lambda, x_1, \dots, x_{p-1}, x_p) = 0$ .

*Proof.* As before, we will use the fact that  $h_p$  is continuous (by construction of  $g_p$ ), and look at the values of  $h_p$  for  $x_p = 0$  and  $x_p = a < 1$ .

First of all for  $x_p = 0$ , for any  $\lambda, x_1, \dots, x_{p-1} > 0$ , we have  $h_p(\lambda, x_1, \dots, x_{p-1}, 0) > 0$ , by construction of  $g_p$ .

Now, let us consider the case  $x_p = a$  for some arbitrary  $a < 1$ . We have  $h_p(\lambda, x_1, \dots, x_{p-1}, a) = g_p(\lambda, x_1, \dots, x_{p-1}, a) - a$ . For  $\lambda = 0$ , we have  $g_p(0, x_1, \dots, x_{p-1}, a) = a \prod_{j \neq p} x_j^{m_{p,j}}$  by Lemma 1.

By Lemma 1, we have in addition  $m_{p,j} > 0$  for some  $j \neq p$ . Thus, for  $x_1, \dots, x_{p-1} < 1$ , we have  $h_p(0, x_1, \dots, x_{p-1}, 1) = a(\prod_{j \neq p} x_j^{m_{p,j}} - 1) < 0$ . Since  $h_p$  is continuous, there exists  $\varepsilon_p$  such that for all  $\lambda \in (0, \varepsilon_p)$ , we have  $h_p(\lambda, x_1, \dots, x_{p-1}, a) < 0$ . By the intermediate value theorem, for all  $\lambda \in (0, \varepsilon_p)$ , there exists  $x_p \in (0, a)$  such that  $h_p(\lambda, x_1, \dots, x_{p-1}, x_p) = 0$ . Here we could actually have taken  $a = 1$ , but we give the scheme that will be used throughout the rest of the induction.  $\square$

We can thus define  $f_p(\lambda, x_1, \dots, x_{p-1})$ , as the least value of  $x_p$  such that  $h_p(\lambda, x_1, \dots, x_{p-1}, x_p) = 0$ , as soon as the conditions of the lemma are met. Moreover,  $f_p$  is continuous as well, since it is defined as the first root of a continuous function with non-zero boundary conditions.

We now get rid of equation  $x_p = g_p$  by replacing  $x_p$  by  $f_p(\lambda, x_1, \dots, x_{p-1})$  in all other equations. We then have a system of the form  $x_i = g_i(\lambda, x_1, \dots, x_{p-1}, f_p(\lambda, x_1, \dots, x_{p-1}))$  for all  $i \in [1, p-1]$ .

Notice that replacing  $x_p$  by  $f_p$  allows to carry out the previous construction for  $f_{p-1}$ , as the only property asked of the various  $x_j$  other than the current  $x_{p-1}$  under consideration is to be in  $(0, 1)$ , and that is the case for  $f_p$ . Notice that  $x_p$  is now a function of  $\lambda, x_1, \dots, x_{p-1}$ , but this is not a problem, as long as it remains in  $(0, 1)$  (allowing for extremal cases if some arguments are extremal as well), the construction can be carried on.

We can thus iterate this construction, and obtain a sequence of functions  $f_i$  as wanted. At each step, the  $\varepsilon_i$  can be chosen to be the minimum of the  $\varepsilon_{i+1}$  obtained at the previous step and the  $\varepsilon$  needed at the current step. This guarantees that all  $f_j$  for  $j \geq i$  are defined for all  $\lambda < \varepsilon_i$ .

Continuing this construction, we will end up with a sequence of functions  $f_i$  such that for all  $\lambda < \varepsilon$ , there exists  $x_1, \dots, x_p$  such that  $x_i = f_i(\lambda, x_1, \dots, x_{i-1})$  for all  $i \in [1, p]$ . This will give us a solution to the system of equations, and thus a solution to the original problem. Indeed, it suffices to choose some  $\lambda < \varepsilon_1$ ,  $x_1 = f_1(\lambda)$ ,  $x_2 = f_2(\lambda, x_1)$ , etc. to obtain a solution to the system.

As described earlier, we can infer from  $x_1, \dots, x_p$  the values of  $x_{p+1}, \dots, x_n$  by using the previous substitutions  $x_j = A_j$ . From this, we can finally compute the concentration vector  $\vec{c}$  such that  $\lambda \vec{v}$  is the flow vector induced by  $\vec{c}$ .

As explained in Corollary 1, this is sufficient to show that any PAC is CAC, as it suffices to start with a PAC witness  $\vec{v}$  and apply this construction to find a CAC witness  $\vec{c}$ .

### 3. On the degree of feasibility of PACs

In this section we propose to explore how thermodynamic constraints affect the degree to which a PAC can be a CAC, in the sense of exploring quantitatively and analytically how such constraints, though not being able to prevent the feasibility of (minimal) PACs (see Corollary 1), nonetheless affect the limitations toward instantiating a PAC in the concentration space. We will focus on the specific topology of type I PACs, following the terminology of [13]. Such a PAC of size  $n$  has the following form (where only the forked reaction  $e_n \rightarrow e_1 + e_1$  is left without food and waste):

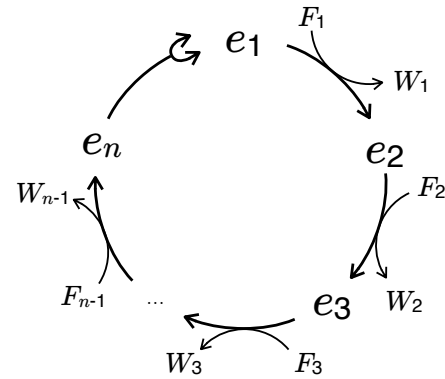


Figure 9. Illustration of a PAC of size  $n$  of type I, using the typology of [13].

Following the notations of the proof explained in Appendix 2, one would end up considering the following system of equations:

$$\begin{cases} F_1 x_1 &= \lambda w_1 + W_1 x_2 \\ F_2 x_2 &= \lambda w_2 + W_2 x_3 \\ &\vdots \\ F_{n-1} x_{n-1} &= \lambda w_{n-1} + W_{n-1} x_n \\ x_n &= \lambda w_n + x_1^2. \end{cases}$$

$F_i$  (respectively  $W_i$ ) is the exponential of chemical potential of food (respectively waste)  $i$ . The substitution algorithm from Section 2 allows us to inject the  $(i+1)^{\text{th}}$  line into the  $i^{\text{th}}$  until one finally obtains a single equation on  $x_1$ :

$$x_1 = \lambda \left( \frac{1}{F_1} w_1 + \frac{W_1}{F_1 F_2} w_2 + \dots + \frac{W_1 \dots W_{n-1}}{F_1 \dots F_{n-1}} w_n \right) + \frac{W_1 \dots W_{n-1}}{F_1 \dots F_{n-1}} x_1^2$$

rewritten as:  $x_1 = g(\lambda, x_1)$ .

For this particular PAC topology, the dependence of  $\lambda$  with respect to other parameters can be made explicit:

$$\lambda = \frac{\left(1 - \frac{1}{\beta} x_1\right) x_1}{\left(\frac{1}{F_1} w_1 + \dots + \frac{W_1 \dots W_{n-2}}{F_1 \dots F_{n-1}} w_{n-1} + \frac{W_1 \dots W_{n-1}}{F_1 \dots F_{n-1}} w_n\right)}. \quad (6)$$

It should be noted that the condition  $\lambda > 0$  imposes  $x_1 \in (0, \beta)$  with  $\beta = \frac{F_1 \dots F_{n-1}}{W_1 \dots W_{n-1}}$ . Notice that if all food and waste variables are set to 1, one recovers the interval of the general proof of Appendix 2 which was  $(0, 1)$ . Regardless, the PAC can be realized for any value of food and waste potentials.

Equation (6) shows that the solution  $\lambda$  has non-trivial dependencies on food, waste and PAC species chemical potentials, activation barriers of reactions and the witness  $\vec{v}$ . To facilitate the discussion, let us assume that the witness  $\vec{v}$  is the one allowing for a unitary production of all PAC species

(such a witness always exist because the stoichiometry matrix of a minimal PAC is square and invertible, see discussion in [13]). Hence, because the witness  $\lambda \vec{v}$  can be instantiated in the concentration space,  $\lambda$  can be interpreted as the amplitude of the PAC, i.e. the level up to which it can effectively produce its entities, or its degree of feasibility, relatively to the case of a unitary production of PAC entities.

Notice that both the numerator and denominator of Equation (6) depends on terms of the form  $\prod_i W_i / \prod_j F_j$ . The denominator becomes smaller as the food variables grow larger relative to the waste variables. As far as the numerator is concerned, we plot in Figure 10 the function  $\theta_\beta(x_1) = (1 - x_1/\beta)x_1$  for two values of  $\beta$  (namely  $\beta = 1$  and  $\beta = 2$ ).  $\theta_\beta$  is an inverted parabola which is positive for  $x_1 \in [0, \beta]$  and maximizes at  $x_1 = \beta/2$ , the maximum being  $\beta/4$ . Hence the numerator of Equation (6) (taken as a function of  $x_1$ ) can get larger as  $\beta$  gets larger, i.e. as food variables get large in comparison to waste variables. We thus conclude that  $\lambda$  can reach higher values when the food variables are in overall large in comparison to waste variables, which from now on is referred to as a favorable environment.

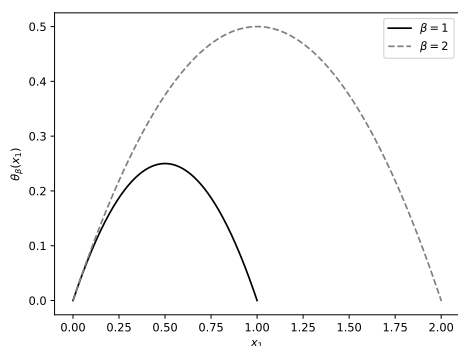


Figure 10. The numerator of Equation (6),  $\theta_\beta(x_1) = (1 - x_1/\beta)x_1$  as a function of  $x_1$  for  $\beta = 1$  and  $\beta = 2$ .

A favorable environment can always be reached by choosing appropriate food and waste concentrations. However it should

be recalled that the variables appearing in Equation (6) are related to the free energy and concentration of entities through the relation  $F_j = e^{G_j} [f_j]$  (illustrated for the  $j^{\text{th}}$  food species).

We have seen that the PAC realization is easier when  $F_j$ 's are large relative to the  $W_i$ 's. Recall that  $W_j = [w_j]e^{G_{w_j}}$ . If one waste free energy is increased by  $\delta G$ , then this corresponds to a multiplication of the corresponding  $W_j$  by  $e^{\delta G}$ , which will have to be compensated by similar variations in concentrations of food and/or waste entities. Thus, slight changes in the free energies values are exponentially passed on the concentration space. This highlights how PACs whose food and waste free energies are unfavorably distributed can be in practice hard to instantiate in the concentration space to the point that in given dynamics this rarely occurs.

Similar conclusions can be reached regarding the PAC entities concentrations. Notice that in Equation (6)  $\lambda$  only depends on  $x_1$  at the numerator that we illustrated in Figure 10. It shows that  $x_1$  is upper bounded, and as a consequence the concentration space of entity  $e_1$  associated to  $\lambda > 0$  (i.e. the PAC running) increases as its free energy is decreases. Notice once again that a favorable environment (i.e. large values of  $\beta$ ) is associated to a wider interval of  $x_1$  allowing  $\lambda > 0$ .

We conclude that regions of the concentration space where the PAC efficiently runs (i.e. runs with large values of  $\lambda$ ) always exist if all concentrations can be chosen freely, but that free energies unfavorably distributed among food, waste and PAC entities will significantly restrict their volume.

On a similar ground, one can study the impact of activation barriers on the feasibility of PACs. Suppose that all reactions of the PAC get penalized by increasing their activation barriers, such that  $G_k^\dagger \leftarrow G_k^\dagger + \Delta$  with  $\Delta > 0$  for all  $k \in [1, n]$ . This will in turn multiply the denominator of Equation (6) by a global factor  $e^\Delta$ , and as a consequence multiply  $\lambda$  by  $e^{-\Delta} < 1$ . This can be compensated by choosing a more favorable environment, i.e. by increasing food concentrations and diminishing waste concentrations such that  $\beta$  increases. Thus, the regions of the concentration space associated to large values of  $\lambda$  are tightened.