



## Autopoiesis of the artificial: from systems to cognition

Francesco Bianchini

University of Bologna, Italy

### ARTICLE INFO

Handling Editor: Dr. A.U. Igamberdiev

#### Keywords:

Autopoiesis  
Biological systems  
Artificial intelligence  
Cognition  
Life and cognition modeling

### ABSTRACT

In the seminal work on autopoiesis by Varela, Maturana, and Uribe, they start by addressing the confusion between processes that are history dependent and processes that are history independent in the biological world. The former is particularly linked to evolution and ontogenesis, while the latter pertains to the organizational features of biological individuals. Varela, Maturana, and Uribe reject this framework and propose their original theory of autopoietic organization, which emphasizes the strong complementarity of temporal and non-temporal phenomena. They argue that the dichotomy between structure and organization lies at the core of the unity of living systems. By opposing history-dependent and history-independent processes, methodological challenges arise in explaining phenomena related to living systems and cognition. Consequently, Maturana and Varela reject this approach in defining autopoietic organization. I argue, however, that this relationship presents an issue that can be found in recent developments of the science of artificial intelligence (AI) in different ways, giving rise to related concerns. While highly capable AI systems exist that can perform cognitive tasks, their internal workings and the specific contributions of their components to the overall system behavior, understood as a unified whole, remain largely uninterpretable. This article explores the connection between biological systems, cognition, and recent developments in AI systems that could potentially be linked to autopoiesis and related concepts such as autonomy and organization. The aim is to assess the advantages and disadvantages of employing autopoiesis in the synthetic (artificial) explanation for biological cognitive systems and to determine if and how the notion of autopoiesis can still be fruitful in this perspective.

### 1. Introduction

In the seminal work on autopoiesis by Varela, Maturana, and Uribe, published almost fifty years ago (Varela et al., 1974), the authors emphasize in the introduction the distinction between history-dependent processes and history-independent processes in the biological world. The former is particularly associated with evolution and ontogenesis, while the latter is connected to the organizational characteristics of biological individuals. The framework in which the concept of autopoiesis is developed is based on an opposition that Varela, Maturana, and Uribe claim to reject. They make an even more significant statement, asserting that the organization of biological systems and their unity precede other aspects of their living nature, such as evolution and reproduction. This implies that something non-temporal takes precedence over a set of temporal phenomena.

This opposition should not mislead us about the core concept of autopoiesis, which underlies the core of living systems. The core revolves around the dichotomy between structure and organization. As we delve further into this article, we will discover that the autopoietic

organization, as defined by Maturana and Varela, is a network of functional relationships involving the synthesis and destruction of components, resulting in the creation of the very components that constitute it. In essence, they describe a dynamic network that produces and destroys its own components, akin to a metabolic network. They emphasize that this dynamic network, as the organizational aspect of living systems, remains constant and unchanged throughout their ontogenetic dynamics. It is the common denominator shared by all living systems and remains unaffected by temporal fluctuations. We can thus speak of complementarity when considering the temporal and non-temporal factors associated with the autopoietic organization and its relationship to life. In living organisms, what remains temporally invariant and what changes over time are intertwined. The structure of living systems is in a constant state of variation, and this perpetual structural variation is what regenerates their organization - a network of relationships between components. Hence, the organization represents the unchanging aspect of life. Simultaneously, the organization, as a dynamic network of functional relations between components, perpetually triggers structural variations and regenerate the structure. Consequently, the

E-mail address: [francesco.bianchini@unibo.it](mailto:francesco.bianchini@unibo.it).

<https://doi.org/10.1016/j.biosystems.2023.104936>

Received 13 February 2023; Received in revised form 28 May 2023; Accepted 29 May 2023

Available online 4 June 2023

0303-2647/© 2023 Elsevier B.V. All rights reserved.

structure represents the variant aspect of life.

This represents the most distinctive aspect of autopoiesis, which can be summarized by considering the organizational aspect as stable and the structural aspect as unstable. The nature of living systems, by definition, requires continuous change in their structural components to maintain stability over time. However, this organizational stability appears to hold greater significance compared to other biological phenomena associated with temporal changes, such as reproduction and evolution. The dichotomy between history-dependent processes and history-independent processes, explicitly highlighted in this article and not revisited in other works by Maturana and Varela (e.g., [Maturana and Varela, 1980](#)), nonetheless carries significance for the framework in which the discourse on autopoiesis is presented (particularly considering the time period in which the article was published) and for the reflections it inspires regarding subsequent attempts in the field of artificial sciences to build synthetic artifacts encompassing life and cognition features. Indeed, the opposition between history-dependent and history-independent processes raises methodological challenges when explaining phenomena related to living systems and cognition. This is why Maturana and Varela reject this approach when defining autopoietic organization. I argue that the issues highlighted by Maturana, Varela, and Uribe in the 1974 article, such as the relationship between history-dependent and history-independent processes, can also be observed in recent developments in the field of artificial sciences, where similar issues arise within major or minor research trends. For instance, the use of cognitive computational architectures has been a long-standing methodology in artificial intelligence (AI) and cognitive science for modeling human cognition ([Lieto, 2021](#)). Computational architectures can be considered as an organizational approach that is highly useful in explaining cognition, despite not fully addressing the coupling between cognitive systems and the environment. On the other hand, the significant advancements in machine learning have brought about the problem of Explainable AI, which entails explaining the behavior and outcomes of neural networks and deep learning systems ([Miller et al., 2022](#)). In this case, these systems demonstrate exceptional performance due to their unified (mathematical and computational) behavior, but the contribution of individual components to the overall system behavior, seen as a unified entity, remains less clear. Here, we encounter another instance in which the system exhibits autonomy, akin to autopoiesis, and possesses an autonomous "life" as an organizational unity with emergent capabilities. But what about cognition? Further examples can be found in the domains of AI and artificial life approaches from the past two decades.

In the well-known book by Maturana and Varela on autopoiesis ([Maturana and Varela, 1980](#)), they establish a close relationship between autopoiesis and cognition, and consequently between life and cognition. However, the two parts of the book, focusing on the biology of cognition and the organization of living systems, also reveal a symbolic gap between these two domains. This gap, present in the bipartite structure of the volume,<sup>1</sup> continues to pose both a problem to explain and a research program. Nevertheless, the underlying assumption that studying life and cognition can benefit from their interconnectedness remains fruitful. It suggests that these domains are not as separate as they may initially appear, even though the study of life and cognition has historically emerged from different fields, each with its own perspectives and aims.

In my article, I aim to explore the relationship between biological systems and cognition by considering recent developments in artificial

<sup>1</sup> It is worth mentioning that the book by Maturana and Varela is composed of two distinct parts. The first part, titled "Biology of Cognition," can be considered a preliminary version of the autopoiesis concept and was written by Maturana in 1969. The second part of the book is the English translation of "Máquinas Y Seres Vivos" (originally written between 1971 and 1972, and published in 1973).

systems and AI trends that can be connected to autopoiesis and related concepts, such as autonomy and organization. The objective is to assess the advantages and limitations of using autopoiesis to explain cognitive biological systems in a synthetic way and to determine whether the notion of autopoiesis remains fruitful in this context. I will attempt to demonstrate that, in most cases, the fundamental features of autopoiesis are not present in artificial systems. While these recent advancements in AI are moving towards more embodied approaches to cognition, they cannot be fully recognized as autopoietic approaches. This limitation hinders their ability to achieve their desired goal of functioning as cognitive systems with autonomous "life." In other words, the challenge of explaining life and cognition through an artificial, synthetic approach is still far from being realized, primarily because autopoietic characteristics are either partially or completely neglected. Furthermore, I will endeavor to show that this limitation is not inherent in principle, and that autopoiesis could still serve as a valuable conceptual framework for the specific goals of the artificial sciences.

A further observation needs to be made regarding Maturana and Varela's conception of cognition, which differs significantly from that of AI systems, especially during the period when the autopoietic theory was formulated. While it is beyond the scope of this article to provide a detailed analysis of Maturana and Varela's notion of cognition, it can be summarized as the idea that living itself is cognition ([Maturana and Varela, 1980](#); [Thompson, 2004](#)). In other words, there is no separation between life and cognition—the aforementioned gap that the autopoietic theory seeks to bridge. Recent advancements in cognitive science have incorporated these conceptual aspects into certain embodied perspectives, particularly enactive approaches that emphasize cognition as the outcome of continuous interaction between a system and its environment, characterized by mutual dynamic coupling. However, in the field of artificial sciences, these developments have not progressed at the same pace as cognitive science approaches. Thus, the aim of this article is to analyze recent developments in AI and identify those approaches that, from an autopoietic perspective, appear most promising in terms of producing synthetic models aligned with an autopoietic vision of cognition. Additionally, the article aims to highlight the limitations of these approaches from an autopoietic standpoint, while also emphasizing the potential contributions of autopoietic theory to future advancements in the realm of artificial sciences.

The structure of the article is as follows. In Section 2, I examine the concept of organization in relation to cognitive systems and autopoiesis, providing an overview of the introduction of autopoiesis by Varela, Maturana, and Uribe. Section 3 explores contemporary approaches in AI to identify potential connections between recent AI methodologies and autopoietic systems. In Section 4, I discuss the "chemical" aspect of AI as the most promising avenue for bridging the gap between AI systems and autopoiesis. Finally, in Section 5, I present the conclusions drawn from this analysis and propose potential future directions for research in this field.

## 2. Organization, cognition, and autopoiesis

The relationship between organizations and cognitive systems has been a longstanding topic in cognitive science ([Bemudez, 2020](#)). From a functionalist perspective, the mind is typically organized into interconnected parts within a well-defined structure or architecture. This traditional view in cognitive science is associated with modularity, which suggests that cognition is the outcome of an organized system of modules that process information through structured interactions. Within this framework, a cognitive system is examined as being structured into different levels of organization that correspond to various cognitive abilities such as perception, learning, decision-making, problem-solving, language, and so on. However, organization can also be

observed in the study of the brain: neurons exhibit organized patterns, and the brain is comprised of distinct neuronal parts or regions. In this case as well, a modular structure exists that links function and matter.<sup>2</sup> Therefore, organization plays a central role in cognitive science and neuroscience modeling, serving as an integral part of the explanation for the functions and capabilities under investigation.

On the other hand, computational cognitive models rely precisely on the use of organization as an explanatory framework. These models aim to capture the general aspects of AI (such as the ongoing efforts to achieve General Artificial Intelligence) or at least some generalizable features. For instance, they focus on the flexible utilization of a conceptual knowledge base or frameworks that explain how information being processed can be cognitively integrated with each other, as seen in global workspace models.<sup>3</sup> In this context, the concept of organization or an organized system becomes an integral part of the mechanistic standpoint, from which mechanistic explanations arise as a consequence. For example, Bechtel and Abrahamsen (2005: 423) provide a broad definition of a mechanism as follows: “A mechanism is a structure performing a function in virtue of its component parts, component operations, and their organization. The orchestrated functioning of the mechanism is responsible for one or more phenomena.” This definition clearly highlights the connection with the notion of organization, which involves the integrated action of the structure’s components. Overall, the use of organization in computational cognitive models aligns with the mechanistic perspective, where the organized functioning of the system’s components plays a crucial role in explaining various phenomena.

A more refined concept of organization is employed to explain systems that exhibit self-organization. In this context, the crucial point lies in the specific structure, coupled with a set of actions, required to achieve such a phenomenon. Self-organizing systems have gained significant prominence in the field of AI, particularly in recent decades. In various domains, the notions of emergence and self-organization are used to elucidate the functioning and origins of systems. For example, in the realm of evolutionary computing, techniques such as evolutionary programming, evolutionary algorithms, and genetic algorithms are employed (Eiben and Smith, 2015). Self-organizing systems are also commonly encountered in neural networks and robotics. In fact, the concept of self-organization was already proposed by Alan Turing in his early writings on AI (Turing, 1950) as a goal to be achieved to develop thinking machines. Self-organizing systems, as evidenced by their applications in AI-related fields, exhibit a strong connection to natural systems, where aspects of life and cognition are intertwined (Holland, 1992). However, the primary focus of AI researchers studying these systems has not been to explain life itself but rather to understand certain aspects of cognition or cognition itself as a whole, while also exploring numerous practical applications (including solving optimization problems).

With the concept of self-organizing systems, AI, including its connections with Artificial Life, approaches the notion of autopoiesis to a significant extent. Autopoiesis, however, possesses distinct characteristics and, perhaps, a greater potential for capturing the relationship between the living and cognition through explanatory and synthetic modeling approaches. The term “autopoiesis” was introduced and thoroughly analyzed in the 1980 book by Humberto Maturana and Francisco Varela, *Autopoiesis and Cognition: The Realization of the Living*. The title itself highlights the authors’ intention to establish an intimate connection between cognition and the living. This development is notable in an era when AI, as an integrated discipline within cognitive science, was predominantly focused on logical-symbolic methodologies for constructing computational models of cognition. Its main applied outcomes revolved around knowledge representation and natural

language processing. The living was not a central focus within AI during those years. However, in the subsequent decade, there was a significant shift, marked by the emergence of neural networks, artificial life, and the study of complex systems in computational terms. It was during this period that the living began to be increasingly integrated into AI’s purview. Furthermore, autopoiesis and self-organization exhibit only a few shared features and are not interchangeable concepts. Autopoiesis primarily focuses on the characteristics of a system that dynamically creates and maintains itself over time, possessing an autonomous self-maintenance mechanism. Autopoiesis, therefore, can be described as the ability of certain systems, such as cells foremostly, to sustain and reproduce themselves by continually generating and reproducing their own constituent parts. While cells exemplify this notion, it can be extended to other types of even more complex systems.

Maturana and Varela establish a connection between the concept of autopoiesis and the notion of a machine, emphasizing how it is closely intertwined with the idea of organization. They describe an autopoietic machine as a machine that is organized and defined as a unity through a network of processes involved in the production, transformation, and destruction of its components. In their own words: “An autopoietic machine is a machine *organized* (defined as a unity) as a *network of processes* of production (transformation and destruction) of components which: (i) through their interactions and transformations continuously regenerate and realize the network of processes (relationships) that produced them; and (ii) constitute it (the machine) as a concrete *unity in space* in which they (the components) exist by specifying the *topological domain of its realization as such a network*” (Maturana and Varela, 1980: 78–79 [emphasis added]).

The autopoietic machine is not merely organized; it is organized according to a principle of unity in space that enables it to operate autonomously. Notably, the principle of autopoietic unity does not necessitate any pre-determined commitment to the material or substance from which the autopoietic machine is constructed. What matters are the relationships within the network that constitute the structure of the machine’s topological domain. Thus, an autopoietic machine can be constructed from various substrates, including artificial ones, as long as its operational and relational characteristics remain unchanged. This aligns with the essential nature of living systems, which continuously change the material they are composed of, while maintaining their unity over time. The definition of an autopoietic machine aims to capture this aspect, as well as the mechanical organizational constitution of a living entity (Damiano and Stano, 2020).

In their 1974 article published in *Biosystems*, authors present the concept of autopoiesis and its strong relationship with organization by defining autopoietic organization as follows:

- “(i) a unity by a network of productions of components which participates *recursively* in the same network of production of components that produced these components, and
- (ii) *realize the networks* of productions as a unity in the space in which the components exist” (Varela et al., 1974: 188 [emphasis added]).

The network of production processes within the autopoietic machine is actively involved in *recursively* producing its own components. This self-application of the machine on itself, guaranteeing unity, autonomy, and persistence over time, is achieved through the concept of recursion. Recursion is fundamental in computation theory as it relates to the definition of computable functions. Interestingly, during the same period when the article for *Biosystems* was written, the connection between recursion, computation, and AI was extensively investigated and emphasized (Hofstadter, 1979). In the case of autopoiesis, recursion serves to establish a principle of self-referentiality in an operational and empirical sense within the network of production processes leading to the production of components. This principle is crucial for maintaining the autopoietic machine’s integrity. Without recursion, the production

<sup>2</sup> On this topic and the related debate see for example Pessoa (2014).

<sup>3</sup> Related to the Global Workspace Theory by Baars (1988).

of external components that are not part of the machine itself would occur, as observed in biological reproduction, which extends beyond the intrinsic homeostatic nature of the machine (Maturana and Varela, 1980: 78). Once again therefore, the connection between autopoiesis and AI has been deeply rooted since the early stages, fostering an almost underground parallelism between the two fields, particularly in biologically inspired approaches to AI.

The primary example of an autopoietic machine can be found in the biological cell, which can be understood as a network of chemical reactions producing molecules, some of which constitute its components. However, it is also possible to consider other entities as autopoietic machines based on a condition that appears necessary, although it is uncertain whether it can be considered sufficient: the close connection between network and organization. In other words, it seems inevitable for an autopoietic machine to rely on a network that can be described in organizational terms, where its parts possess *well-defined functionalities*. This aspect is particularly relevant in network cognition models, including neural networks, where not all organizational aspects are fully explainable. Consequently, this poses an increasing challenge in the field of explainable AI and interpretable machine learning (Murdoch et al., 2019).

Within the autopoiesis framework, Varela, Maturana, and Uribe view living systems as a subset of mechanistic systems. This perspective diverges from a line of research that gained significant importance in subsequent decades, particularly in cognitive neuroscience: the mechanistic explanation with a focus on decomposition and localization (Bechtel and Richardson, 1993; Craver, 2001), where the mechanisms are seen both in their contextualization within a system and from the point of view of their constituent parts. In this explanatory approach, the organization of the context in which a mechanism works and of its constituent parts plays a crucial role. The aim is to overcome the classical formulation of the multiple realizability thesis (Bechtel and Mundale, 1999) and address the challenge of establishing the connection between cognitive states and neural states.

While the decomposition and localization approach moves toward resolving the aforementioned problem, the autopoietic organization remains aligned with the notion of multiple realizability. This stems from the specific nature of autopoietic organization, where “the *same* organization may be realized in *different* systems with different kinds of components as long as these components possess the properties that realize the *required relations*. It is *obvious* that with respect to their *organization* such systems are members of the *same class*, even though with respect of the *nature* of their components they may be *distinct*” (Varela et al., 1974: 188 [emphasis added]). This is interesting for two reasons. Firstly, even from a mechanistic explanation perspective, the link between the cognitive and the living hinges on the organizational nature of the system. Secondly, the concept of autopoiesis preserves the essential dichotomy of “same organization/different systems”, which appears to be a fundamental requirement for creating artificial systems that can be deemed as having life and cognition. Autopoiesis, therefore, retains a quality that many subsequent attempts in cognitive science and cognitive neuroscience seem to have to relinquish.

The emphasis on organization is significant as it allows for the preservation of multiple realizability regardless of the material composition of a system. This requirement aligns with the standard approach in AI and cognitive science, particularly the symbolic manipulation perspective. However, the outcome differs substantially in the context of autopoiesis, which explicitly relies on organizational and (partially) self-organizational aspects, not typically present in the traditional AI approach. Another crucial point is made in the beginning of the 1974 *Biosystems* article. The authors distinguish between history-dependent and history-independent processes, arguing that the overlap between these two types of processes has led to confuse phenomena which instead must be kept distinct. Specifically, reproductive and evolutionary processes are not considered constitutive features of living organization which instead “can only be characterized unambiguously by

specifying the networks of interactions of components which constitute a living system as a whole, that is, as a unity” (Varela et al. 1974: 187). In this framework, reproductive and evolutionary features, like other biological phenomena, are considered secondary to the adequate organization of living, which is deemed a necessary and sufficient condition for life itself. The complementary relationship between organizational and structural aspects of the living unity condenses invariant and variant features of a system within the autopoietic perspective, where temporality and change contribute to the possibility of a living system. Other historical phenomena of life take a subordinate role to the understanding of temporality within the dynamics of structure and organization. For instance, the reproduction of a living system preserves the invariance of its relational unity across generations, ensuring stability (Damiano and Stano, 2023). Thanks to this dynamics, temporal processes do not affect the stability of the living. Consequently, temporal transformations that living beings inevitably undergo do not impose an excessively rigid constraint on achieving an alleged “objective” stability, nor do they serve as a reason for system dissipation or destruction.

This balance represents one of the most original contributions of the autopoietic theory, which remains unique even today in comparison to the directions taken by AI and cognitive AI in recent decades. In fact, it is precisely this aspect that sets autopoiesis apart, at least to some extent, from certain more recent bio-inspired approaches to AI. Only when considering autopoiesis in terms of organizational aspects, specifically autopoietic organization, can it be seen as a precursor to Artificial Life. A question that arises is whether this particular aspect can make Artificial Life or other AI fields more suitable for achieving the goal of creating artificial systems that can be considered alive and cognitive, or if it instead poses a hindrance. Undoubtedly, this characteristic of the autopoiesis concept brings AI systems, in which autopoiesis can be identified, closer to Synthetic Biology (SB), another discipline that can be fully classified as one of the sciences of the artificial (Bianchini, 2021). Therefore, what are the AI systems, if any, in which autopoiesis can be discussed in terms of Varela, Maturana, and Uribe?

### 3. Autopoiesis, organization, and AI

The field of SB, as a discipline within the realm of the artificial, is intriguing within the emerging framework because it tackles the problem of artificiality from a living perspective. Its goal is to create something that can be defined as living, without relying heavily on evolutionary techniques. This characteristic brings SB closer to autopoiesis. For instance, one of SB’s aims has been to construct fundamental living organisms like cells (Buddingh’ and van Hest, 2017; Gaut and Adamala, 2021). While achieving this goal has proven to be quite challenging, SB’s focus appears to be more oriented towards the realm of the living compared to Artificial Life. It goes beyond simply creating artificial systems that fulfill a set of modeling requirements and constraints to be considered living. Instead, SB aims to start with living systems and progress towards artificial living systems, utilizing both materials and organizational structures characteristic of living organisms (such as genetic and biochemical matter). Particularly in terms of the latter, SB aligns closely with the concept of autopoiesis in its conventional sense.

If it is possible to transition from living systems to artificial living systems, as pursued by SB, one might contemplate whether it is feasible to take a step further and move from living systems to cognitive systems using the same techniques. While SB still seems distant from providing an answer to this question, it could be valuable to consider the notion of autopoiesis and its conceptual foundations as useful elements for attaining artificial cognitive systems within the broader field of cognitive AI. However, this line of research does not appear to be widely explored or fully pursued in relation to the most recent approaches to AI (Waser, 2014; Damiano and Stano, 2021).

The development of AI over the decades has resulted in the emergence of diverse approaches, leading to a multitude of methodological

and epistemological perspectives within the fields of AI and cognitive AI. These approaches encompass various methodologies, including the traditional logical-symbolic approach, dynamic systems, traditional machine learning, neural networks, complex systems, and biology-inspired methods. Each of these approaches exhibits its own internal variations and is influenced by several factors, such as the conceptual frameworks that define them. These approaches, characterized by their internal diversity, now coexist within the field of AI, contingent upon several factors, including the conceptualizations that shape their theoretical foundations. They are influenced, thus, by factors such as the notion of intelligence, available technologies, the concept of a system, the notion of autonomy, the aims of AI, different intelligent capabilities (such as language, reasoning, creativity, sensorimotor abilities, emotions, morality, etc.), the role of the body, the role of the brain, and the *notion of cognition*. In exploring the potential role of autopoiesis in relation to cognitive AI, two contemporary approaches that appear relevant from this perspective will be considered.

A potential initial statement to define the role of autopoiesis in relation to cognitive AI is as follows: if autopoiesis emphasizes organization rather than history-dependent processes such as evolution and reproduction, then organizational AI systems that are *specifically related to cognition* should incorporate *autopoietic features or mechanisms*. Is this claim valid? Does it offer promising prospects? To address these inquiries, we will examine the first of the two relevant approaches."

The first approach is the cognitive architectures AI approach. Cognitive architectures are modular systems in which different parts are joined and organized to exhibit various behaviors, aiming to replicate the general aspects of intelligence incrementally by adding blocks/modules. One of the most well-known examples is SOAR (State, Operator, and Result), which is based on Newell's unified cognitive theory. This theory emphasizes the Problem Space Hypothesis (Newell, 1990), a cornerstone of classical AI that considers all intelligent behavior as goal-seeking within a state space. SOAR, as a cognitive architecture (Laird, 2012), consists of modules that process symbolic information and are designed to implement cognitive abilities such as memory, perception, learning, decision making, and motor activities. Another widely recognized cognitive architecture is ACT-R (Adaptive Control of Thought - Rational), which also processes symbolic information and features an organized modular structure (Anderson, 2007). Various versions of ACT-R have been developed for different modeled tasks, and there have also been proposals for connectionist extensions based on neural networks (Labièrre and Anderson, 1993). A key and significant aspect of this framework for constructing cognitive models is its direct inspiration from the modular organization of brain functions (such as memory, control, motor execution) and the corresponding brain areas. Like other cognitive architectures, this approach relies on a modular structure, with each module serving specific functions. The choice of whether to use a symbolic approach or neural networks in building the modules depends on the particular architecture.<sup>4</sup>

A common feature shared by many cognitive architectures is the assumption of modularity, which is utilized in various ways. This approach involves implementing different functions within the architecture that are identified as cognitive or relevant to cognitive activities, corresponding to specific parts of the architecture. A direct consequence of this modular construction principle is that these parts are organized in a mutually interconnected manner. In other words, a cognitive architecture can be viewed as a structured system in which organization arises from the integration and exchange of information among its modules, i.e. a *module network integration*. The architecture is composed of modules that are recognized as cognitive or relevant to specific activities, and they interact by exchanging symbolic information *at their respective levels*, with the purpose of forming a comprehensive cognitive

model. This organizational characteristic seems to bear similarities to the properties of autopoietic entities, albeit with a fundamental difference. In the case of cognitive architectures, it follows an allopoietic rather than an autopoietic approach. This means that the system as a whole, functioning as a unified unit, is built by a programmer with explicit reference to cognitive parts. Even though these parts/modules may implement self-organization processes in information processing, the overall system's construction and maintenance are mostly external to the system itself. These external processes may be rooted in human or evolutionary programming, or of another type altogether. Consequently, attempts to explain or model cognition within this framework differ somewhat from those applicable to autopoietic systems. In this case, cognition does not result from the autopoietic organizational nature of the system, and therefore, an explanation or construction of the cognitive dimension *in such terms* is absent.

A more promising approach that bears resemblance to autopoietic systems appears to be that of neural networks. In neural networks, the biological inspiration derived from certain elements of brain neurons brings cognitive performance closer to the biological system that produces it, at least in principle. In artificial neural networks, the organization, characterized by layers of nodes/neurons connected by weighted arcs, plays a crucial role. The connections can involve all the neurons between two layers or only a subset of them. The system's topology is essential in achieving specific performances after the training phase. For instance, in contrast to traditional feedforward neural networks where information flows only in one direction, recurrent neural networks (Schmidhuber, 2015) allow the output of certain nodes to become input for the nodes themselves, creating cycles within the network. This dynamic behavior enables the network to exhibit understanding and recognition of temporally determined events. In cases like these, as well as in other types of neural networks (e.g., convolutional networks), and more broadly in the field of deep learning with organized systems across multiple layers of neurons, the system's topology plays a critical role in achieving cognitive performances. These performances include tasks such as categorization, pattern recognition, sensorimotor action, and others, which directly arise from the network's organization. Therefore, the organization of the network is crucial for cognitive performance.

Even in the case of neural networks, or specific types of neural networks, a close resemblance to autopoietic systems from a cognitive standpoint is not present. Neural networks are not a form of autopoiesis because while the system does self-organize, cognitive performance emerges as a result of self-organization following a training phase within a pre-existing system built with specific features. This pre-existing system is characterized by the organization of nodes and connections into layers and follows a predetermined topology. The autopoietic proposal extensively explores the relationship between cognition and biology (as seen in a significant portion of the book *Autopoiesis and Cognition*, for example). The nervous system is described as the foundation for major cognitive abilities such as learning and memory. It is based on the principles of unity, self-regulation, and continuous self-construction, which underlie the internal and external interactive nature of the nervous system's behavioral dynamics. This third feature, continuous self-construction, not only enables cognition at the system's inception but at every moment throughout its lifespan. As stated by Maturana and Varela (1980:119), "for any autopoietic system, its cognitive domain is necessarily relative to the particular way in which autopoiesis is realized". The autopoietic nature is absent in neural networks, which, despite their good or excellent cognitive performance, face a challenge when it comes to explaining the processes and outcomes they produce. This challenge is known as the problem of explainable AI (Miller et al., 2022), which is common in systems where understanding the individual components' contributions is difficult. This problem has helped highlight the gap between neural networks and explanations of cognitive phenomena. The autopoietic dimension, specifically the self-referentiality and self-observation features of autopoietic systems, could shed light on aspects of cognitive explanation and contribute to

<sup>4</sup> A detailed discussion of the main features of these and other cognitive architectures is in Lieto (2021).

bridging the explainability gap. However, the exploration of neural networks evolving in this direction is still largely unexplored.

#### 4. A “chemical” AI?

The analyzed cognitive approaches to AI, although fundamentally different in their principles and underlying assumptions, share two common features: 1) they are organizational, and 2) they are non-autopoietic. The fact that these two features work so well to build systems implementing – even in very different ways – performances which are recognizable as cognitive, divide living systems from cognitive ones from the standpoint of autopoiesis. Computational cognitive models perform well and can be integrated within the synthetic methodology framework (Datteri and Tamburrini, 2007; Webb, 2001) for modeling and simulating cognition or its various aspects. On the other hand, living systems are characterized by the realization of some form of autopoietic principles as their fundamental feature. In essence, autopoiesis is the foundation of life, while non-autopoiesis (in the form of allopoietic systems) serves as the basis for cognitive systems. The gap between living systems and cognitive systems can be observed in the different synthetic methodologies used to construct artificial systems. Even in cases where biological inspiration is strong, as in evolutionary computation, the living organisms that serve as inspiration do not adhere to the autopoietic principles that precede and underlie evolutionary processes, as explained in the seminal 1974 article on autopoiesis. Evolutionary processes (an example of history-dependent processes) can exist without autopoiesis. Consequently, the resulting systems cannot be considered, for all (autopoietic) intent and purposes, living systems in an autopoietic sense. As a result, if autopoiesis is at the core of the living, artificial systems built in such a manner are not properly living and do not belong to true artificial life.

A possible solution to this gap could involve considering the chemical or biochemical aspects of living organisms as inspiration for computational modeling, given their significance in autopoietic processes. The relational and organizational nature of autopoietic systems heavily relies on the underlying biochemical processes, which enable the system and its components to exhibit self-creation and self-maintenance capabilities. Therefore, the question arises: should we take into account the “chemical” features of AI systems to bridge this gap? What are the chemical characteristics underlying living organisms and the cognition we aim to model? Recently, there has been a new line of research known as Chemical Artificial Intelligence that explores the construction of computational systems based on chemical and biochemical aspects. This approach aims to chemically implement binary or multi-valued logic, fuzzy logic, artificial neuron models, and chemical robots (Gentili, 2013, 2022). Chemical AI strives to develop intelligent chemical systems in wetware, emulating for example basic aspects of human intelligence related to logic. The use of hardware and software to reproduce the same phenomena is still little explored, but in principle there seems to be no hindrance to implement the relationships and organizational processes involved in this type of phenomena on standard computational substrates. This could potentially lead to the creation of artificial autopoietic systems capable of exhibiting intelligent or cognitive behavior.<sup>5</sup>

The utilization of computational methods to investigate biochemical phenomena in living organisms dates back to the early stages of AI. In a 1952 article, Turing discusses the concept of morphogenesis, proposing that “a system of chemical substances, called morphogens, reacting together and diffusing through a tissue, is adequate to account for the main phenomena of morphogenesis” (Turing, 2004: 519). These reaction-diffusion systems are studied in relation to an isolated ring of cells, aiming to explore a possible mechanism by which the genes of a

zygote can establish the anatomical structure of the resulting organism. Turing emphasizes that this pattern formation arises from well-known physics laws and specifies that the development of an organism typically involves transitioning from one pattern to another, rather than from homogeneity to a pattern. While this general process can be mathematically explained, Turing acknowledges that there is no comprehensive theory encompassing all these processes, aside from stating the equations. Turing’s proposal is to employ computational methods (specifically “digital computers”) as “this method has the advantage that it is not so necessary to make simplifying assumptions as it is when doing a more theoretical type of analysis. It might even be possible to take the *mechanical aspects* of the problem into account as well as the *chemical*, when applying this type of method” (Turing, 2004, 561 [emphasis added]). This highlights why computational methods are particularly valuable for addressing certain chemical aspects *in process* at the basis of living as an organic unity. In section K of the Turing Digital Archive, there are a couple of sheets where Turing has calculated the function of a morphogenesis system, revealing the spatial structure in which the pattern takes shape.<sup>6</sup> This represents the “chemical” computation that simulates the *inner* structure of the system.

Another well-known attempt to understand the fundamental aspects of life from a computational perspective during the early days of AI is the Von Neumann approach to self-replicative systems. In the 1948 Hixon Symposium, he introduced the concept of a logical theory of self-replication and explored the role of replication errors in the process of evolution (von Neumann, 1951). The purpose was to develop a logical theory of automata, later known as cellular automata, by uncovering the logical principles underlying evolutionary phenomena. Through the application of specific rules, individual entities are capable of self-reproduction, leading to the emergence of complex effects and structures. These repeating patterns are formed by the assembly of multiple individuals, following simple logical laws of self-replication. Cellular automata represent computational spaces where cells undergo changes based on predefined rules, mirroring similar phenomena observed in living organisms. Von Neumann focused on simulating self-reproduction through logical models that encompassed both single cells and organic clusters composed of multiple cells. His research explored two dimensions of self-replication for simulating self-reproduction:

- 1) the self-replication of a single entity: a cell.
- 2) the self-replication of systems consisting of cells replicating themselves.

In the second case, the focus lies on the study of replication and self-replication of complex systems, which reproduce themselves at an emergent level, not necessarily the highest (von Neumann, 1966). Subsequent developments in this line of research yielded intriguing mathematical results that underpin the study of living systems and laid the foundations for the emergence of Artificial Life (Langton, 1986). Within the framework of a chemical AI implementation, what is particularly noteworthy is that the logical laws of self-replication simulate cellular biochemical processes at a certain level of abstraction. In this specific case, unlike Turing’s morphogenesis, chemical computation can be seen as a simulation of *outer* processes of living entities. Moreover, it is intriguing to note that if these processes are considered external to the cell, they cannot be solely classified as internal or external in the context of multicellular organisms, the actual cellular automata. The “chemical” dimension collapses the dichotomy of inside/outside, manifesting a series of processes that are constitutive of the system’s unity. This convergence closely aligns these systems with autopoietic systems, where the internal organizational nature plays a

<sup>5</sup> On this topic and the notion of a “chemical explorative AI” see also Damiano and Stano (2023).

<sup>6</sup> The files are, in particular, two with the following references: AMT/K/3/11 and AMT/K/3/12.

constitutive role in shaping the dynamic and organic processual structure of the system.

Even the complex adaptive systems approach has some aspects in common with autopoietic systems. For example, AI genetic algorithms (Holland, 1992; Mitchell, 1998) utilize certain features of natural selection as a metaheuristic to find solutions for search and optimization problems. Their primary focus is on performance efficiency rather than the synthetic modeling of natural processes. The same principle can be applied to complex adaptive systems in general, which consist of numerous components interacting with each other and capable of adapting to the context and learning from it (Holland, 2006). The field of complex adaptive systems is vast, but several common features can be identified. These systems are based on dynamic networks of interactions, exhibit unpredictable behavior, and are self-organizing. They exist as far-from-equilibrium systems that continuously update themselves based on the context to which they need to adapt. While they share self-organizing characteristics and are built upon process networks of nonlinear dynamics, determining their boundaries or structural aspects of organization is often challenging, unlike autopoietic systems.

One of the goals of complex adaptive systems is to study the emergent macroscopic properties observed in complex systems, such as living systems. However, they can be considered as implementing a computation to simulate inner and outer processes of the system, with a focus on the desired outcome. The “biochemical” features of their component interactions are in the background. What matters is the system as a whole and its emergent behavior. In this sense, they do not constitute a strong attempt to simulate or explain the living and cognitive phenomena to the same extent as autopoietic systems can be. Complex adaptive systems primarily aim to logically define and computationally leverage the laws of formation, replication, and evolution present in biological systems. These phenomena are history-dependent, as described by Maturana and Varela, and they do not exhibit autopoietic organizational properties. The explicit consideration of self-maintenance properties, related to the complementary relationship between structure and organization, is not addressed.

A final opportunity to establish connections with autopoiesis in AI and Artificial Life systems, which simulate the chemical features of biological systems, can be found in systems inspired by specific biological entities: superorganisms (Hölldobler and Wilson, 2009), such as ant colonies (Bonabeau et al., 1998). These systems fall under the umbrella of complex adaptive systems but specifically focus on the collective behavior of swarm systems (Bonabeau et al., 1999) associated with intelligence, cognition, or robotics (according to the related field of modeling). These bio-inspired artificial systems adopt organizational features and relational networks of processes from living organisms, modeling their interaction methods based on biochemical signals. An ant colony can be considered an autopoietic organization in many respects, including unity (even without distinct boundaries), network structure, and chemical relationships. The specific mechanism of stigmergy, based on pheromones, enables widespread coordination among the agents forming the colony in response to the environmental context. This exemplifies an agent/environment relationship based on chemical factors. In a superorganism, life and cognition, i.e., the parts and their goal-oriented, cognitively understandable functions, are closely intertwined. However, they are the life and cognition of the superorganism itself, rather than the individual organisms that comprise it. Consequently, life and cognition emerge as properties of the collective behavior of the colony-forming agents, which can be computationally modeled. The connection with the life and cognition of the constituents assumes a secondary role, as observed in relation to complex adaptive systems in general. This seems to violate the principle of unity that makes autopoietic machines what they are, by making superorganisms and the systems inspired by them a collection of autopoietic machines, rather than a unitary system. Nevertheless, the intertwined aspects of life and cognition remain intriguing from the perspective of simulating and explaining life and cognition through computational models

inspired by superorganisms. A potential solution to reconcile these two perspectives could involve redefining the notion of biological unity within the framework of autopoiesis. By adopting a more flexible understanding of system unity, the relationship between life and cognition can be better characterized within an autopoietic organizational view. However, it is crucial to determine to what extent a deflationary notion of system unity is acceptable within the autopoietic framework without compromising its explanatory power regarding the characteristics of living systems it supports.

## 5. Conclusion

In the final part of their 1974 *Biosystems* article, Varela, Maturana, and Uribe present a model that serves as a simple embodiment of autopoietic organization. This model depicts a universe composed of only a few elements capable of composing, concatenating, or disintegrating. The rules governing the interactions are explicitly formulated, allowing for a step-by-step execution of the resulting computational model. The actions described by these rules primarily involve chemical processes, such as catalyzing reactions formalized within the considered universe. Autopoietic organization is the specific process generated by the properties of the components, enabling the creation of the system’s dynamic unity. Importantly, the autopoietic organization is not preexisting in the initial state of the elements, whether in representational or embodied forms. It emerges as a result of the inherent properties of the system. The ultimate outcome is a network of processes involving the components, which becomes a recognizable entity within the universe and co-produces itself along with its components. Prior to the beginning of the overall process, there is no pre-existing network or its components. As stated by authors, “The properties of an autopoietic system [...] are determined by the constitution of this unity, and are, in fact, the properties of the network created by, and creating, its components. Therefore, to ascribe a determinant value to any component, or to any of its properties, because they seem to be ‘essential’, is a semantic artifice” (Varela et al., 1974: 192). While components are necessary for the production of the network, none of them is essential in the process. Each component contributes to the network’s formation, but no individual component holds inherent importance.

This leads to a continuous and radical form of emergentism, wherein every interaction and step is explicitly accounted for, even though it is undoubtedly impossible to predict the evolution of the system in the most complex living systems. This aligns with the computational irreducibility often encountered in attempts to model such complex systems (Zwirn and Delahaye, 2013). However, the contribution of individual components remains clear and transparent. Their nature as components is a consequence of the process that makes the interaction among processes within a dynamic network the crucial element of the autopoietic mechanism. There is no external creation, but rather a self-construction that inherently enables self-maintenance as a constitutive property of the system itself. The macro and micro levels are distinct yet interconnected without abrupt transitions. In this sense, life and cognition can be traced back to the same foundation, the same set of principles governing the formation and persistence of the system as an entity. From a cognitive standpoint, these principles gradually give rise to increasingly complex modes of interaction with the external environment, transitioning from being inherent properties of living systems to forms of cognition.

It is beyond the scope of this discussion to delve into other aspects that connect autopoiesis and cognition, which Maturana and Varela also explore and establish from the beginning. For instance, the autopoietic aspects of neural systems and processes, which form the basis of the study of cognition in autopoietic terms, have not been addressed here. The primary objective of this article is to emphasize and explicitly highlight the role of organization in autopoietic systems and to investigate the relationship between the properties of living systems and their

cognitive capabilities. To this end, various contemporary approaches in AI, where the organization of the system is crucial and which exhibit bio-inspired characteristics, have been considered. However, these approaches are not fully realized autopoietic systems to varying degrees. The hypothesis is that bio-inspired “chemical” AI approaches (along with their models and simulations) serve as promising candidates for bridging the gap between life and cognition within the autopoietic framework of artificial modeling. Nonetheless, this line of research is still far from yielding conclusive results. The concept of applying computational AI to chemistry (as opposed to utilizing AI for the analysis of chemical and biochemical data to make new discoveries) remains an underexplored field, despite the aforementioned foundations present in the evolution of AI.

Nonetheless, if the goal is to situate the explanation and synthetic modeling of cognition within the framework of the living, autopoiesis, in its original formulation, provides a valid conceptual framework, especially considering the development of computational theories on cognition in recent decades. Additionally, the fundamental assumption is that achieving a complete realization of cognition detached from the living is hard. Instead, significant progress has been made in modeling specific cognitive performances over the past two decades, inevitably bringing forth epistemological challenges for AI systems, such as black box systems, explainability, appropriate simulation definitions, modeling problems, and the general relationship between natural and artificial systems. The hypothesis supported in this article is that by utilizing autopoiesis theory and its specific notion of organization (a network of invariant processes that self-produce along with its components, independent of the substrate), it is possible to create self-emergent systems that serve as promising candidates for AI (cognitive) systems. Moreover, these systems should also exhibit self-maintenance at a lower level. Consequently, drawing stronger and deeper inspiration from chemical and biochemical processes in computational modeling is crucial to further develop complex autopoietic computational mechanisms.

### Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

### References

- Anderson, J.R., 2007. *How Can the Human Mind Occur in the Physical Universe?* Oxford University Press.
- Baars, B.J., 1988. *A Cognitive Theory of Consciousness.* Cambridge University Press, Cambridge, Mass.
- Bechtel, W., Abrahamsen, A., 2005. Explanation: a mechanist alternative. *Stud. Hist. Phil. Biol. & Biomed. Sci.* 421–441. <https://doi.org/10.1016/j.shpsc.2005.03.010>.
- Bechtel, W., Mundale, J., 1999. Multiple realizability revisited: linking cognitive and neural states. *Philos. Sci.* 66, 175–207.
- Bechtel, W., Richardson, R.C., 1993. *Discovering Complexity: Decomposition and Localization as Strategies in Scientific Research.* Princeton University Press, Princeton (NJ).
- Bemudez, J.L., 2020. *Cognitive Science. An Introduction to the Science of the Mind*, third ed. Cambridge University Press, Cambridge.
- Bonabeau, E., Theraulaz, G., Deneubourg, J.L., 1998. The synchronization of recruitment-based activities of ants. *Biosystems* 45 (3), 195–211. [https://doi.org/10.1016/S0303-2647\(98\)00004-5](https://doi.org/10.1016/S0303-2647(98)00004-5).
- Bonabeau, E., Dorigo, M., Theraulaz, G., 1999. *Swarm Intelligence from Natural to Artificial Systems.* Oxford University Press, Oxford.
- Bianchini, F., 2021. A New Definition of “Artificial” for Two Artificial Sciences. *Found Sci.* <https://doi.org/10.1007/s10699-021-09799-w>.

- Buddingh', B.C., van Hest, J.C.M., 2017. Artificial cells: synthetic compartments with life-like functionality and adaptivity. *Acc. Chem. Res.* 50 (4), 769–777. <https://doi.org/10.1021/acs.accounts.6b00512>.
- Craver, C.F., 2001. Role functions, mechanisms and hierarchy. *Philos. Sci.* 68, 53–74.
- Damiano, L., Stano, P., 2020. On the “life-likeness” of synthetic cells. *Front. Bioeng. Biotechnol.* 8, 953. <https://doi.org/10.3389/fbioe.2020.00953>.
- Damiano, L., Stano, P., 2021. A wetware embodied AI? Towards an autopoietic organizational approach grounded in synthetic biology. *Front. Bioeng. Biotechnol.* 9 <https://doi.org/10.3389/fbioe.2021.724023>.
- Damiano, L., Stano, P., 2023. Explorative Synthetic Biology in AI. Criteria of relevance and a taxonomy for synthetic models of living and cognitive processes. *Artif. Life* 29 (3).
- Datteri, E., Tamburrini, G., 2007. Biorobotic experiments for the discovery of biological mechanisms. *Philos. Sci.* 74 (3), 409–430. <https://doi.org/10.1086/522095>.
- Eiben, A.E., Smith, J.E., 2015. *Introduction to Evolutionary Computing.* Springer Berlin, Heidelberg. <https://doi.org/10.1007/978-3-662-44874-8>.
- Gaut, N.J., Adamala, K.P., 2021. Reconstituting natural cell elements in synthetic cells. *Advanced Biology* 5, 2000188. <https://doi.org/10.1002/adbi.202000188>.
- Gentili, P.L., 2013. Small steps towards the development of chemical artificial intelligent systems. *RSC Adv.* 3, 25523–25549.
- Gentili, P.L., 2022. Photochromic and luminescent materials for the development of chemical artificial intelligence. *Dyes Pigments* 205, 110547. <https://doi.org/10.1016/j.dyepig.2022.110547>.
- Hofstadter, D.R., 1979. *Gödel, Escher, Bach: an Eternal Golden Braid.* Basic Books, New York.
- Holland, J.H., 1992. *Adaptation in Natural and Artificial Systems. An Introductory Analysis with Applications to Biology, Control, and Artificial Intelligence.* The MIT Press, Cambridge, Mass.
- Holland, J.H., 2006. *Studying Complex Adaptive Systems*, vol. 19. *Jrl Syst Sci & Complex*, pp. 1–8. <https://doi.org/10.1007/s11424-006-0001-z>, 2006.
- Hölldobler, B., Wilson, E.O., 2009. *The Superorganism. The Beauty, Elegance, and Strangeness of Insect Societies.* W. W. Norton & Company, New York.
- Lebiere, C., Anderson, J.R., 1993. A connectionist implementation of the ACT-R production system. In: *Proceedings of the Fifteenth Annual Conference of the Cognitive Science Society* Lawrence Erlbaum Associates, New York, pp. 635–640.
- Laird, E.J., 2012. *The Soar Cognitive Architecture.* The MIT Press, Cambridge, Mass.
- Langton, C.G., 1986. Studying artificial life with cellular automata. *Phys. Nonlinear Phenom.* 22, 120–149.
- Lieto, A., 2021. *Cognitive Design for Artificial Minds.* Routledge.
- Maturana, H., Varela, F., 1980. *Autopoiesis and Cognition. The Realization of Living.* Reidel Publishing Company, Dordrecht.
- Miller, T., Hoffman, R., Amir, O., Holzinger, A., 2022. Special issue on explainable artificial intelligence (XAI). *Artif. Intell.* 307 <https://doi.org/10.1016/j.artint.2022.103705>.
- Mitchell, M., 1998. *An Introduction to Genetic Algorithms.* The MIT Press, Cambridge, Mass.
- Murdoch, W.J., Singh, C., Kumbier, K., Abbasi-Asl, R., Yu, B., 2019. Definitions, methods, and applications in interpretable machine learning. *Proc. Natl. Acad. Sci. U. S. A.* 116 (44), 22071–22080. <https://doi.org/10.1073/pnas.1900654116>.
- Newell, A., 1990. *Unified Theories of Cognition.* Harvard University Press, Cambridge, Mass.
- Pessoa, L., 2014. Understanding brain networks and brain organization. *Phys. Life Rev.* 11 (3), 400–435. <https://doi.org/10.1016/j.plrev.2014.03.005>.
- Schmidhuber, J., 2015. Deep learning in neural networks: an overview. *Neural Network* 61, 85–117.
- Thompson, E., 2004. Life and mind: from autopoiesis to neurophenomenology. A tribute to Francisco Varela. *Phenomenol. Cognit. Sci.* 3, 381–398.
- Turing, A.M., 1950. Computing machinery and intelligence. *Mind* 59 (236), 433–460.
- Turing, A.M., 2004, 1952–1954. *The Chemical Basis of Morphogenesis.* Philosophical Transactions of the Royal Society of London, Series B, vol. 237, pp. 37–1954. reprinted in J. Copeland (ed.), *The Essential Turing.* Clarendon Press, Oxford, pp. 519–1954.
- Varela, F.G., Maturana, H.R., Uribe, R., 1974. Autopoiesis: the organization of living systems, its characterization and a model. *Biosystems* 5, 187–196. [https://doi.org/10.1016/0303-2647\(74\)90031-8](https://doi.org/10.1016/0303-2647(74)90031-8).
- von Neumann, J., 1951. The general and logical theory of automata. In: Jeffress, L.A. (Ed.), *Cerebral Mechanisms in Behavior: the Hixon Symposium. Proceedings of a Meeting Held in Pasadena, 1948.* Wiley, New York, pp. 1–41.
- von Neumann, J., 1966. Theory of self-reproducing automata. In: *And Completed by A. W. Burks.* University of Illinois Press, Urbana and London.
- Waser, M.R., 2014. Bootstrapping a structured self-improving & safe autopoietic self. *Proc. Comput. Sci.* 41, 134–139.
- Webb, B., 2001. Can robots make good models of biological behaviour? *Behav. Brain Sci.* 24, 1033–1050.
- Zwrin, H., Delahaye, J.-P., 2013. Unpredictability and computational irreducibility. In: Zenil, H. (Ed.), *Irreducibility and Computational Equivalence: Emergence, Complexity and Computation*, vol. 2. Springer-Verlag, Berlin, Heidelberg, pp. 273–295. [https://doi.org/10.1007/978-3-642-35482-3\\_19](https://doi.org/10.1007/978-3-642-35482-3_19).