

# A MODEL OF THE TRANSITION TO BEHAVIORAL AND COGNITIVE MODERNITY USING REFLEXIVELY AUTOCATALYTIC NETWORKS

LIANE GABORA AND MIKE STEEL

**ABSTRACT.** This paper proposes a model of the cognitive mechanisms underlying the transition to behavioral and cognitive modernity in the Upper Paolelithic using Reflexively Autocatalytic and Food set generated (RAF) networks. Autocatalytic networks have been used to model life's origins. More recently, they have been applied to the emergence of *cognitive* structure capable of undergoing *cultural* evolution. Mental representations of knowledge and experiences play the role of catalytic molecules, the interactions among them (e.g., the forging of new associations or affordances) play the role of reactions, and thought processes are modeled as chains of these interactions. Building on an earlier autocatalytic model of the cognitive transition underlying the transition from Oldowan to Acheulian tool technology, we posit that a genetic mutation allowed thought to be spontaneously tailored to the situation by modulating the degree of (1) divergence (versus convergence), (2) abstractness (versus concreteness), and (3) context-specificity. This culminated in persistent, unified cognitive autocatalytic networks that bridged previously compartmentalized knowledge and experience. We explain the model using the example of the oldest-known uncontested example of figurative art: the carving of the Löwenmensch, or lion-man. The model suggests an explanation for the lag between anatomically modern *Homo sapiens* and behavioral/cognitive modernity.

*Keywords:* autocatalytic network, behavioral modernity, cognitive modernity, cultural evolution, Hohlenstein-Stadel Löwenmensch figurine, semantic network, Upper Paleolithic

Corresponding Author:

Liane Gabora

Department of Psychology, University of British Columbia

Okanagan Campus, Kelowna BC, Canada

[liane.gabora@ubc.ca](mailto:liane.gabora@ubc.ca)

Mike Steel:

Biomathematics Research Centre, University of Canterbury, Christchurch, New Zealand

[mike.steel@canterbury.ac.nz](mailto:mike.steel@canterbury.ac.nz)

## 1. INTRODUCTION

How did we become distinctively human? What enabled us to develop imagination, ingenuity, and complex belief systems? These questions are central to understanding who we are, how we got here, and where we are headed. Behavioral and cognitive modernity are thought to have come about in the Upper Paleolithic, between 100,000 and 30,000 years ago, as evidenced by the sudden proliferation in cultural artifacts of both utilitarian and aesthetic value.<sup>1</sup> Some attribute this transition to an enhanced ability to process social information [94, 98]. Cognitive explanations have been proposed; for example, it has been attributed to the onset of conceptual fluidity [74], dual modes of information processing [28, 79], or enhanced working memory [19]. We propose that each of these proposals holds merit and that they are not mutually exclusive but, in turn, reflect the onset a new kind of semantic network structure, which is modeled here.

Although evidence of human culture dates back millions of years, behavioral-cognitive modernity is associated with the transition to cultural change that is not just adaptive (new innovations that yield some benefit for their bearers tend to predominate), but also cumulative (later innovations build on earlier ones) and open-ended (the space of possible innovations is not finite, since each innovation can give rise to spin-offs). In other words, culture became an *evolutionary* process [9, 14, 17, 25, 34, 49, 71, 84]. By *culture*, we mean extrasomatic adaptations, including behavior and artifacts, that are socially rather than genetically transmitted. Although cultural *transmission*—in which one individual acquires elements of culture from another—is observed in many species, cultural *evolution* is much rarer (and perhaps, unique to our species).<sup>2</sup>

---

<sup>1</sup>Some researchers push behavioral modernity back further [102], and the concept itself has been called into question [22]. This paper does not delve into these discussions so as to focus squarely on the task of modeling the cognitive changes underlying this cultural transition.

<sup>2</sup>The term ‘cultural evolution’ is occasionally used in a less restricted sense to refer to the generation and transmission of novelty without the requirement of cumulative, adaptive, open-ended change (e.g., [99]).

Networks allow for a comprehensive understanding of the dynamics of complex entities and their relationships [45]. Network-based approaches to characterizing the kind of cognitive structure that could sustain cultural evolution enable us to address the question of how minds carry out the contextual, combinatorial, and hierarchically structured thought processes needed to generate cumulative, adaptive, and open-ended cultural novelty [30, 39]. Here, rather than a generic semantic or neural network, we use an autocatalytic network. Autocatalytic network theory grew out of studies of the statistical properties of *random graphs* consisting of nodes randomly connected by edges [26]. As the ratio of edges to nodes increases, the size of the largest cluster increases, and the probability of a phase transition resulting in a single giant connected cluster also increases. The recognition that connected graphs exhibit phase transitions led to their application to efforts to develop a formal model of the origin of life (OOL), namely, of how abiogenic catalytic molecules crossed the threshold to the kind of collectively self-sustaining, self-replicating structure we call ‘alive’ [61, 60]. In the application of graph theory to the OOL, the nodes represent catalytic molecules and the edges represent reactions. It is exceedingly improbable that any catalytic molecule present in the primordial soup of Earth’s early atmosphere catalyzed its own formation. However, reactions generate new molecules that catalyze new reactions, and as the variety of molecules increases, the variety of reactions increases faster. As the ratio of reactions to molecules increases, the probability increases that the system will undergo a phase transition. When, for each molecule in a set, there is a catalytic pathway to its formation, the set is said to be collectively *autocatalytic*, and the process by which this state is achieved has been referred to as *autocatalytic closure* [61]. The molecules thereby become a self-sustaining, self-replicating structure (i.e., a living protocell [53]). Thus, the theory of autocatalytic networks has provided a promising avenue for modeling the OOL and thereby understanding how biological evolution began [100]. The approach transitions is consistent claims for the centrality of transitions across the life sciences [87].

Autocatalytic networks have been developed mathematically and generalized for cross-disciplinary application in other settings in the theory of Reflexively Autocatalytic Food set generated (RAF) networks [54, 89]. The term *reflexively* is used in its mathematical sense, meaning that every element is related to the whole. The term *food set* refers to the reactants that are initially present, as opposed to those that are the products of catalytic reactions. RAFs have been used extensively to model the origins of biological evolution [54, 89, 95, 100]. Thus, one strength of the approach is that by adapting a formalism that has been used successfully to model one evolutionary process to model another, we pave the way for a broad conceptual framework that can shed light on both [30, 4]. This is in keeping with the suggestion that autocatalytic networks may hold the key to understanding the origins of *any* evolutionary process, including the origin of culture [30, 31, 34, 37, 41].<sup>3</sup> In application to culture, the products and reactants are not catalytic molecules but culturally transmittable *mental representations*<sup>4</sup> (MRs) of experiences, ideas, and chunks of knowledge, as well as more complex mental structures such as schemas and scripts (Tables 1 and 2).

<b>Graph Theory</b>	<b>Origin of Life (OOL)</b>	<b>Origin of Culture (OOC)</b>
node	catalytic molecule	mental representation (MR)
edge	reaction pathway	association
cluster	molecules connected via reactions	MRs connected via associations
connected graph	autocatalytic closure [60, 61]	conceptual closure <sup>5</sup> [30]

TABLE 1. Application of graph theoretic concepts to the origin of life (OOL) and origin of culture (OOC).

<sup>3</sup>For related approaches, see [4, 16, 77].

<sup>4</sup>Although we use the term ‘mental representation’, our model is consistent with the view (common amongst ecological psychologists and in the situated cognition and quantum cognition communities) that what we call mental representations do not ‘represent’, but instead act as contextually elicited bridges between mind and world.

Abbreviation	Meaning
OOL	Origin of Life
OOC	Origin of Culture
MR	Mental Representation
RR	Representational Redescription
RAF	Reflexively Autocatalytic and Food set generated (F-generated)
CCP	Cognitive Catalytic Process

TABLE 2. Abbreviations used throughout this paper.

Another strength of the approach is that because it distinguishes reactants that are external in origin—in our case, MRs that were acquired through social learning or individual learning of *existing* information—from those that are the products of internal reactions—in our case, MRs that come about through the creative generation of *new* information—MRs are tagged with their source. This enables us to model how networks emerge and to trace cumulative change in cultural lineages step by step.

In previous work, we used the RAF framework to model what is arguably the earliest significant transition in the archaeological record: the transition from Oldowan to Acheulean tool technology approximately 1.76 million years ago (mya) [41, 42]. We posited that this was precipitated by onset of the capacity for *representational redescription* (RR), in which the contents of working memory are recursively restructured by drawing upon similar or related ideas, or through concept combination. This enabled the forging of associations between MRs, and the emergence of hierarchically structured concepts, making it possible to shift between levels of abstraction as needed to carry out tasks composed of multiple subgoals. This culminated in what is referred to as a transient RAF, a critical step toward what has been referred to as conceptual closure [31] and the emergence of a persistent cognitive RAF.

The application of RAFs presented in this paper builds on that work to model what is perhaps the most spectacular cultural transition in human history: the cultural transition of the Upper Paleolithic, which has been referred to as the “origins of art, religion, and science” [72]. We propose that behavioral and cognitive modernity was brought about by the emergence of a semantic network that was autocatalytic. We first summarize the archaeological evidence for a transition to behavioral modernity in the Upper Paleolithic. We then present our RAF model of the underlying cognitive transition that brought it about. Finally, we compare and contrast our proposal with existing literature.

## 2. ARCHAEOLOGICAL EVIDENCE FOR BEHAVIORAL AND COGNITIVE MODERNITY

We begin with a brief summary of the evidence for a transition to behavioral and cognitive modernity in the Upper Paleolithic.<sup>6</sup> Although one can argue that the earliest stone tools marked the onset of a ‘proto’ form of cultural evolution, following the initial appearance of the Acheulian hand axe, the archaeological record exhibits over a million years of cultural stasis, and—with the exception of a more sophisticated knapping (Levallois) technique 200,000-400,000 years ago—little in the way of creative embellishment or improvement [92].

This changed dramatically in the Aurignacian period of the Upper Paleolithic, at which point there is evidence of recognizably human ways of living and thinking. The earliest evidence of behavioral and cognitive modernity comes from Africa less than 100,000 years ago, in Sub-Himalayan Asia and Australasia more than 50,000 years ago [76], and in Continental Europe until approximately 30,000 years ago [70]. This evidence consists of a proliferation of different complex, task-specific tools including effective cutting blades [3, 66]. It also marks the appearance of representational art [5, 78, 82], artifacts indicating personal symbolic ornamentation [24], complex living spaces [81], sophisticated ways of obtaining food, including aquatic resources [27], burial sites indicating ritual [55] and possibly religion [85]. The Upper

---

<sup>6</sup>A more detailed discussion can be found elsewhere [39, 40].

Paleolithic is also widely believed to have marked the onset of modern syntactically rich language [11] (though some argue that language arose more gradually; c.f. [64]). In short, this period witnessed an unprecedented dramatic increase in the variety, utility, and aesthetic value of cultural artifacts.

A celebrated example of Upper Paleolithic art to which we will devote considerable attention is the Löwenmensch or ‘lion-man’ figurine from the Hohlenstein-Stadel cave in Germany (Figure 1). This figurine, carbon-dated to the Interpleniglacial period between 35,000 and 40,000 years ago, is one of the oldest-known zoomorphic (animal-shaped) sculpture in the world, and one of the oldest-known examples of figurative art. It measures 31.1 cm, and was carved out of mammoth ivory using a flint stone knife.

### 3. AN UNDERLYING COGNITIVE TRANSITION

The model of the transition to cognitive and behavioral modernity in the Upper Paleolithic developed here grew out of the hypothesis that it was due to the onset of *contextual focus*: the capacity to, in a spontaneous and ongoing manner, shift between convergent and divergent modes of thought, thereby tailoring ones’ mode of thought to one’s situation [32, 33, 39].

Focused attention is conducive to *convergent thought* because the activation of cell assemblies is constrained enough to zero in on the most defining properties. In this compact form, the contents of thought are more readily amenable to deliberate executive level operations. In convergent thought, one can access only *close associates* of the current thought: items that are highly related to it with respect to the most conventional, default context. For example, FIG and PLUM are close associates because they are both fruits, and they are most commonly thought about with respect to their membership in the category FRUIT.

In contrast, defocused attention is conducive to *divergent thought* because it causes diffuse activation of cell assemblies in memory, such that obscure (but potentially relevant) properties come into play [31, 33, 35, 36]. This is useful for creative tasks, and when one is in need of a



FIGURE 1. Sketch of the Löwenmensch or 'lion-man' figurine from the Hohlenstein-Stadel cave in Germany. According to the Ulm Museum,  $^{14}\text{C}$  dates put it at an age of 35,000 to 40,000 years. (Obtained with permission from the artist, Cameron Smith).

new approach or innovative solution. Divergent thought may include both more details of the current subject of thought, or incorporate related items; one is simply taking in more of the situation and its associations. In divergent thought, one can access *remote associates*: words or concepts that not related to each other with respect to their most conventional, default context. Highly divergent thought may result in cross-domain thinking, in which ideas from different domains are combined, or a solution to a problem in one domain is borrowed from another domain.

Although divergent thought is useful for escaping local minima, it confers the risk of getting perpetually side-tracked, whereby irrelevant thoughts readily intrude, interfering with survival tasks. Unless the capacity for divergent thought is accompanied by the ability to reign it in, it would be counterproductive. Therefore, it seems likely that in the pre-modern mind, before the advent of the capacity to shift along the spectrum from convergent to divergent, all mental contents were processed convergently, such that each successive thought was a close associate of its predecessor, and remote associates were not accessed.

Contextual focus came about through the onset of the capacity to adjust the focus of attention to current constraints and affordances, making it more focused or diffuse, as needed, thereby stretching or shrinking conceptual space, and tailoring working memory to task demands (or lack thereof, as in mind wandering). Contextual focus made it possible to use convergent thought to modify the content of working memory on the basis of close associates when that proves sufficient, and divergent thought to usher in remote associates when ‘stuck in a rut’. The theory that contextual focus can have a transformative impact on cultural evolution was tested using an agent-based model [39]. Incorporating the ability to shift between convergent and divergent processing modes into neural network-based agents in the agent-based model resulted in an increase in the mean fitness of cultural outputs.

The model that follows builds on the hypothesis that behavioral modernity was due to the onset of contextual focus, but goes further in positing that thought acquired the capacity to shift along a multimodal spectrum if thought through spontaneous tuning of the following three variables.

**3.1. Divergence.** The first variable is the capacity to spontaneously shift between divergent and convergent thought, as discussed above.

**3.2. Level of abstraction.** The second variable is *degree of abstraction*. It has been shown that that there is what is called a *basic level* of abstraction (e.g., BIRD, as opposed to ANIMAL

or SWALLOW) that mirrors the correlational structure of properties in the object's real-world perception and use [86]. Categories form, and are first learned and perceived, at this basic level, before they are further discriminated at the subordinate level (e.g., SWALLOW), and abstracted at the superordinate level (e.g., ANIMAL).<sup>7</sup> Since basic-level categories contain the degree of abstraction most useful for carrying out daily activities [86], it seems reasonable that they precede other levels of abstraction not just developmentally but evolutionarily. We posit that the arrival of behavioral/cognitive modernity involved onset of the capacity to shift along the hierarchy from abstract to concrete, thereby identifying relatedness at different hierarchical levels, and incorporating these distinctions into one's mental model of the world. Abstraction provides another means of connecting MRs, but instead of forging a remote association between them, it makes explicit that they are both instances of some more general MR (e.g., LION and MAMMOTH are both instances of ANIMAL). Thus, the second variable entails a facility shift from basic-level categories to other levels of abstraction.

**3.3. Context-specificity.** To generate ideas and solutions that are not just new but also task-relevant may require thought that is not just divergent but also context-specific [35]. Thus, the third variable is *context-specificity*: the degree to which thought is biased by a specific motivating contextual factor such as a goal or desire [7]. Divergent thought need not *always* be context-specific (e.g., during mind-wandering or writing free verse). However, context-specific divergent thought allows one to access information that is related to the current contents of working memory in ways that may be unconventional yet precisely relevant to the current situation [69]. For example, thinking of lions in the context of wishing for an inspirational reminder of a lion's power might prompt one to modify one's concept of lion to incorporate the

---

<sup>7</sup>Note that abstract processing is not the same as convergent processing. An item at a particular level of abstraction, such as LION, would, in convergent thought, be held in working memory in a compact manner stripped of details, whereas during divergent thought, it would be rich in the characteristics of, and feelings evoked by, lions. One might speculate that richly detailed visions of religious deities occurs in a mode of thought that is abstract yet divergent.

possibility of carving a lion. This unusual context makes this remote yet feasible relationship ‘pop out’.

**3.4. Multimodality.** The spectrum of thought is multimodal, where by a ‘mode’ we refer to a particular combination of the three variables (e.g., divergent, abstract, and context-specific). In short, we posit that by using this multimodal spectrum to modify how one thought gives way to the next, cognitive processes could be carried out more effectively. Moreover, the fruits of one mode of thought could become ingredients for another mode, thereby facilitating the forging of a richly integrated network of understandings about the world and one’s place in it, sometimes referred to as a *worldview*. This, we posit, set the stage for behavioral modernity.<sup>8</sup>

Note that although it seems likely that we are ‘wired for culture’ [65], and the cognitive changes underlying this cultural transition were brought on by one or more genetic mutations [32, 40, 19, 20], we are not proposing that control over these variables came online instantaneously, nor that control over each of them arose simultaneously. The challenge may have been not so much to *possess* the capacity to change these variables as to *coordinate* them so as to continuously tune one’s mode of thought in response to changing task demands and effectively navigate semantic space. The evolutionary and developmental tinkering required to achieve this could explain the lag between anatomically modern *Homo sapiens* 200,000 to 100,000 years ago, and behavioral modernity 100,000 to 30,000 ago.

#### 4. AUTOCATALYTIC NETWORKS

We now summarize the key concepts of RAF theory. A *catalytic reaction system* (CRS) is a tuple  $\mathcal{Q} = (X, \mathcal{R}, C, F)$  consisting of a set  $X$  of molecule types, a set  $\mathcal{R}$  of reactions, a catalysis set  $C$  indicating which molecule types catalyze which reactions, and a subset  $F$  of  $X$  called the

---

<sup>8</sup>Note that, in this view, language enhanced not just the ability to communicate and collaborate (thereby accelerating the pace of cultural innovation), but also the ability to think ideas through for oneself and manipulate them in a manner that was controlled, deliberate, and multimodal.

food set. A *Reflexively Autocatalytic and F-generated* set—i.e., a RAF—is a non-empty subset  $\mathcal{R}' \subseteq \mathcal{R}$  of reactions that satisfies the following two properties:

- (1) *Reflexively autocatalytic*: each reaction  $r \in \mathcal{R}'$  is catalyzed by at least one molecule type that is either produced by  $\mathcal{R}'$  or is present in the food set  $F$ ; and
- (2) *F-generated*: all reactants in  $\mathcal{R}'$  can be generated from the food set  $F$  via a series of reactions only from  $\mathcal{R}'$  itself.

A set of reactions that forms a RAF is simultaneously self-sustaining (by the *F-generated* condition) and (collectively) autocatalytic (by the RA condition; as each of its reactions is catalyzed by a molecule associated with the RAF). A CRS need not have a RAF, but when it does there is a unique maximal one. Moreover, a CRS, may contain many possible RAFs, and it is this feature that allows RAFs to evolve as demonstrated (both in theory and in simulation studies) through selective proliferation and drift acting on possible subRAFs of the maxRAF [54, 95].

In the OOL context, a RAF emerges in systems of polymers (molecules consisting of repeated units called monomers) when the complexity of these polymers (as measured by their maximum length) reaches a certain threshold [61, 75]. The phase transition from no RAF to a RAF incorporating most or all of the molecules depends on (1) the probability of any one polymer catalyzing the reaction by which a given other polymer was formed, and (2) the maximum length (number of monomers) of polymers in the system. This transition has been formalized and analyzed (mathematically and via simulations), and applied to real biochemical systems [50, 51, 52, 54, 75], ecology [18], and cognition [41, 42]. RAF theory has proven useful for identifying how phase transitions might occur, and at what parameter values.

**4.1. Terminology.** We now introduce the mathematical framework and terminology that will be used to model the transition to cognitive modernity. All mental representations (MRs) in a given individual  $i$  are denoted  $X_i$ , and a particular MR  $x = x_i$  in  $X_i$  is denoted by writing

$x \in X_i$ . As in an OOL RAF, we have a *food set*; for individual  $i$ , this is denoted  $F_i$ . In the origin of culture (OOC) context,  $F_i$  encompasses MRs for individual  $i$  that are either innate, or that result from direct experience in the world, including natural, artificial, and social stimuli.  $F_i$  includes everything in the long-term memory of individual  $i$  that was not the direct result of individual  $i$  engaging in RR. This includes information obtained through social learning from *someone else* who may have obtained it by way of RR. For example, if individual  $i$  learns from individual  $j$  how to edge a blank flake through percussive action, this is an instance of social learning, and the concept EDGING is therefore a member of  $F_i$ .

$F_i$  also includes existing information obtained by  $i$  through individual learning (which, as stated earlier, involves learning from the environment by nonsocial means), so long as this information retains the form in which it was originally perceived (and does not undergo re-description or restructuring through abstract thought). The crucial distinction between food set and non-food set items is not whether another person was involved, nor whether the MR was originally obtained through abstract thought (by *someone*), but whether the abstract thought process originated in the mind of the individual  $i$  in question. Thus,  $F_i$  has two components:

- $\mathbb{S}_i$  denotes the set of MRs arising through direct stimulus experience that have been encoded in individual  $i$ 's memory. It includes MRs obtained through social learning from the communication of an MR  $x_j$  by another individual  $j$ , denoted  $\mathbb{S}_i[x_j]$ , and MRs obtained through individual learning, denoted  $\mathbb{S}_i[l]$ , as well as contents of memory arising through direct perception that does not involve learning, denoted  $\mathbb{S}_i[p]$ .
- $I_i$  denotes any *innate knowledge* with which individual  $i$  is born.

A particular catalytic event (i.e., a single instance of RR) in a stream of abstract thought in individual  $i$  is referred to as a *reaction*, and denoted  $r \in \mathcal{R}_i$ . A stream of abstract thought, involving the generation of representations that go beyond what has been directly observed, is modeled as a sequence of catalytic events. Following [41], we refer to this as a *cognitive catalytic*

*process* (CCP). The set of reactions that can be catalyzed by a given MR  $x$  in individual  $i$  is denoted  $C_i[x]$ . The entire set of MRs either *undergoing* or *resulting from*  $r$  is denoted  $A$  or  $B$ , respectively, and a member of the set of MRs undergoing or resulting from reaction  $r$  is denoted  $a \in A$  or  $b \in B$ .

The term *food set derived*, denoted  $\neg F_i$ , refers to mental contents that are *not* part of  $F_i$  (i.e.  $\neg F_i$  consists of all the products  $b \in B$  of all reactions  $r \in R_i$ ). In particular,  $\neg F_i$  includes the products of any reactions derived from  $F_i$  and encoded in individual  $i$ 's memory. Its contents come about through mental operations *by the individual in question* on the food set; in other words, food set derived items are the direct product of RR. Thus,  $\neg F_i$  includes everything in long-term memory that *was* the result of one's own CCPs.  $\neg F_i$  may include a MR in which social learning played a role, so long as the most recent modification to this MR was a catalytic event ( i.e., it involved RR).<sup>9</sup>

The set of *all* possible reactions in individual  $i$  is denoted  $\mathcal{R}_i$ . The mental contents of the mind, including all MRs and all RR events, is denoted  $X_i \oplus \mathcal{R}_i$ . This includes  $F_i$  and  $\neg F_i$ . Recall that the set of all MRs in individual  $i$ , including both the food set and elements derived from that food set, is denoted  $X_i$ .

$\mathcal{R}_i$  and  $C_i$  are not prescribed in advance; because  $C_i$  includes reminders and associations on the basis of one or more shared property, different CCPs can occur through interactions amongst MRs. Nevertheless, it makes perfect mathematical sense to talk about  $\mathcal{R}_i$  and  $C_i$  as sets. Table 3 summarizes the terminology and correspondences between the OOL and the OOC.

Our model includes elements of cognition that have no obvious parallel in the OOL. We denote the subject of attention at time  $t$  as  $w_t$ . It may be an external stimulus, or a MR retrieved from memory. Any other contents of  $X_i \oplus \mathcal{R}_i$  that are accessible to working memory,

---

<sup>9</sup>This distinction between food set and food set-derived may not be so black and white but for simplicity we avoid that subtlety for now.

Term	Origin of Life (OOL)	Origin of Culture (OOC)
$X_i$	all molecule types in protocell $i$	all mental representations (MRs) in individual $i$
$x \in X_i$	a molecule in $X_i$	a MR in $X_i$
$F_i$	food set for protocell $i$	innate or directly experienced MRs by $i$
$r \in \mathcal{R}_i$	a particular reaction in $i$	a particular representational redescription (RR) in $i$
$C_i[x]$	reactions catalyzed by $x$ in $i$	RR events ‘catalyzed’ by $x$ in $i$
$(x, r) \in C$	$x$ catalyzes $r$	$x$ ‘catalyzes’ redescription of $r$
$a \in A$	member of set of reactants in $r$	member of set of MRs undergoing $r$
$b \in B$	member of set of products of $r$	member of set of MRs resulting from $r$
$\neg F_i$	non food set for $i$ (i.e., all $B$ of $\mathcal{R}_i$ )	MRs resulting from $R_i$ (i.e., all $B$ of $R_i$ )

TABLE 3. Terminology and correspondences between the Origin of Life (OOL) and the Origin of Culture (OOC).

such as close associates of  $\hat{w}_t$ , or recently attended MRs, are denoted  $W_t$ , with  $W_t$  constituting a very small subset of  $X_i \oplus \mathcal{R}_i$ . The focus here is on how non-food set derived MRs (i.e., a non-empty  $\neg F$ ) emerge and connect giving rise to a semantic network that is self-organizing and autocatalytic.

## 5. RAF MODEL OF THE COGNITIVE TRANSITION

We now use the RAF formalism to model the transition to behavioral/cognitive modernity in the Upper Paleolithic. To address how the mind as a whole acquired autocatalytic structure, the model is, by necessity, abstract. It does not distinguish between semantic memory (memory of words, concepts, propositions, and world knowledge) and episodic memory (personal experiences); indeed, we are sympathetic to the view that these are not as distinct as once thought [63]. Nor does it address how MRs are obtained (i.e., whether through Hebbian learning versus probabilistic inference). Although MRs are represented simply as points in an  $N$ -dimensional

space (where  $N$  is the number of distinguishable differences, i.e., ways in which MRs could differ), our model is consistent with models that use convolution [57], random indexing [58], or other methods of representing MRs.

We assume that associations form between MRs but do not address whether they are due to similarity or co-occurrence, or whether they are learned through Bayesian inference [43] or other means. We view associations as probabilistic; when we say that an association was forged between two MRs we mean a spike in the probability of one MR evoking another, which we refer to as the ‘catalysis’ of one MR by the other. We view context as anything external (e.g., an object or person) or internal (e.g., other MRs) that influences the instantiation of a MR in working memory. Although our approach is influenced by how context is modeled in quantum approaches to concepts [2, 1], it is not committed to any formal approach to modeling context.

MRs are composed of one or more *concepts*: mental constructs such as CAT or FREEDOM that enable us to interpret new situations in terms of similar previous ones. The rationale for treating MRs as catalysts comes from the literature on concept combination, which provides extensive evidence that when concepts act as contexts for each other, their meanings change in ways that are often non-trivial and defy classical logic [1, 2, 47, 80]. The extent to which one MR modifies the meaning of another is referred to here as its *reactivity*. A given MR’s reactivity varies depending on the other MRs present in working memory.<sup>10</sup> Although we do not explicitly model the dimensionality of semantic space itself (i.e., the features or properties of MRs), we do so indirectly, by representing hierarchical structure in terms of reactivity, as explained below. Our model hinges on the fact that interactions between two or more MRs in working memory alter (however slightly) the network of association strengths [15, 67]. Conceptual closure is

---

<sup>10</sup>For example, in a study of the influence of context and mode of thought on the perceived meanings of concepts (as measured by property applicabilities and exemplar typicalities), the concept PYLON was rated low as an exemplar of HAT; however, in the context FUNNY (as in ‘worn to be funny’), it was rated high as an exemplar of HAT [96]. Thus, the degree to which PYLON qualified as an instance of a HAT changed depending on the context. The context FUNNY had an even greater effect on the rating of MEDICINE HAT (as in the name of the Canadian town) as an instance of HAT. We say that the *reactivity* was high here because the context exerted a dramatic influence on the perceived meaning.

achieved and a cognitive RAF network emerges when, for each MR, there is an associative pathway to its formation; in other words, any given concept can be explained using other concepts, and new ideas can be re-framed in terms of existing ones.

We now show how the RAF framework is used to model the emergence of a persistent and integrated cognitive RAF, through onset of the capacity to spontaneously control the ‘spectrum of thought’ variables introduced in Section Three, and summarized in Table 4. In this table, the variables  $\gamma_D$ ,  $\gamma_A$ , and  $\gamma_C$  quantify the three variables: divergence ( $D$ ), abstractness ( $A$ ) and context-specificity ( $C$ ), respectively.

Variable	Example	Symbol
divergence	LION $\rightarrow$ KILL $\rightarrow$ POWER	$\gamma_D$
abstractness	LION $\rightarrow$ ANIMAL $\rightarrow$ ANIMATE BEING	$\gamma_A$
context-specificity	LION (context: desire to possess lion’s power) $\rightarrow$ LION FIGURINE	$\gamma_C$

TABLE 4. Examples of the three variables of the spectrum of thought.

We model the capacity to shift between convergent and divergent thought by introducing a metric geometry. We let  $d$  denote the *semantic distance* between an item  $m$  in memory  $M_t$  and an item in working memory  $\hat{w}$ . In convergent thought, the semantic distance  $d$  between successive contents of working memory remains small, and only close associates catalyze RR reactions and participate in CCPs. In contrast, by spontaneously engaging in divergent thought when stymied, or as a form of mental exploration or mind-wandering, the modern mind gained access to remote associates (i.e., items for which the semantic distance  $d$  to the content of working memory was large). These remote associates catalyzed RR reactions, and participated in CCPs. Thus, divergent thought could (in our terminology) bring about reactions amongst previously unconnected MRs, including MRs from different knowledge domains. The variable  $\gamma_D$  determines how ‘remote’ an associate can be in order to catalyze an update (i.e., how ‘far

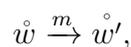
afield’ one looks for ingredients in one’s stream of thought). Thus,  $\gamma_D$  provides a threshold on  $d$  that increases as one shifts from convergent to divergent thought.

The more abstract a concept, the more associations it can have with other MRs. Therefore, we represent hierarchical semantic structure from concrete instances to increasingly abstract concepts in terms of reactivity. Consequently, the more abstract (as opposed to concrete)  $x$  in individual  $i$  is, the larger the value of  $C_i[x]$ . Thus, abstract concepts facilitate the navigation of semantic space through CCPs. For example, during the transition from thinking about a particular sharp axe to thinking about the abstract quality of sharpness, abstractness increase, and therefore so does the reactivity, potentially leading the CCP to something quite different from a sharp axe, such as a ‘claw’.

As mentioned earlier, the capacity for divergent thought could be made even more useful by broadening the sphere of associates in a context-specific manner, such that one’s current needs bias the retrieval of information from memory. This facilitates the forging of new connections between MRs that would be irrelevant in most contexts but relevant in the current one. We make the notion of context-specificity more precise by introducing a context-dependent association structure. As above, we let  $d$  denote the *semantic distance* between an item  $m$  in memory  $M_t$  and an item in working memory  $\hat{w}$ . A small value of  $d(m, \hat{w})$  means that in the current context,  $m$  is closely related to  $\hat{w}$  whereas a large value of  $d(m, \hat{w})$  means that with respect to the current context, they are distantly related.

Let  $\mathcal{C}_t$  denote a context at time  $t$  (which is determined by the goals and needs of the individual at time  $t$ ). We can represent the context-dependent associations explicitly by writing  $m \sim_{\mathcal{C}_t} m'$  if  $m$  and  $m'$  are related with respect to context  $\mathcal{C}_t$ .

For  $m$  to catalyze a cognitive updating reaction



$m$  should satisfy at least one of the following properties (where  $\gamma_D$  is as described above, and  $\gamma_C$  is the extent to which context can facilitate the catalysis of a particular RR reaction):

- (i)  $d(m, \dot{w}) \leq \gamma_D$ , or
- (ii)  $m \sim_{\mathcal{C}_t} \dot{w}$  and  $d(m, \dot{w}) \leq \gamma_D(1 + \gamma_C)$ .

In other words, for  $m$  to catalyze a RR reaction involving  $\dot{w}$ , either the default semantic distance to  $m$  must be sufficiently small that divergent thought makes it accessible, or it is pulled within reach because context-specificity warps semantic space in such a way as to make this particular association salient. In addition, a particular context  $\mathcal{C}_t$  at time  $t$  may mean that a stimulus  $s_t$  that is relevant to the current contents of working memory, catalyzes a RR reaction  $\dot{w} \xrightarrow{s} \dot{w}'$  that would not occur otherwise. For example, seeing an animal puncture a food source with its claw could be a source of ideas for how to make a tool sharper.

**5.1. Example: The Hohlenstein-Stadel figurine.** We now make the transition from pre-modern to modern mind more concrete using the example of the Löwenmensch or ‘lion-man’ figurine from the Hohlenstein-Stadel cave, discussed in Section 2. Although we cannot know exactly how the Hohlenstein-Stadel figurine was created, by reverse-engineering the process it is possible to infer what conceptual structure would, at a minimum, have had to be in place [38, 93, 97]. We carry this out using available evidence, such as our knowledge that the lion was the largest and most dangerous predator in the ecosystem of the Interpleniglacial [83, 62], and likely a source of fear and awe due to its power and aggression [46]. Since the word ‘representation’ is often used to refer to an internal, mental construct of something in the world, to avoid confusion, we use the term *iconic* to refer to an object that represents something else in a way that is not merely symbolic but captures its physical attributes.

We now consider the sequence of steps culminating in the creation of the lion man, summarized in Table 5 and depicted in Figures 2 and 3. Note that the steps culminating in the Hohlenstein-Stadel figurine, were preceded by, and dependent upon, the development of lithic

reduction (i.e., knapping and carving) techniques. (These are not discussed here, since it is the subject of another paper [42].)

Step	Description	Origin	Mode
1.	Carve from stone $\rightarrow$ carve (from something)	RR	abstract
2.	Transmission of lithic reduction techniques	social learning	convergent
3.	Carve functional tool $\rightarrow$ carve (something)	RR	divergent, abstract
4.	Transmission of CARVE (something)	social learning	convergent
5.	Carve (something) $\rightarrow$ carve iconic likeness	RR	divergent, concrete
6.	Combine human form and lion head $\rightarrow$ internalize lion's power	RR	divergent, cross-domain, context-specific
7.	Assimilate features of lion and man	individual learning	divergent
8.	Carve Löwenmensch	individual learning + RR	divergent

TABLE 5. The sequence of steps culminating in the creation of the Hohlenstein-Stadel Löwenmensch figurine.

- (1) **Form abstract concept, CARVE** (from any suitable material). This consisted of abstracting the general concept of lithic reduction with stone as the source material to lithic reduction using any suitable material (e.g., mammoth ivory). There is evidence that the capacity to abstract a general concept from particular instances dates back to at least 1.76 mya (well before the Paleolithic) [90].
- (2) **Social learning of first step.** Creative contributions to culture begin with a preparation stage involving thorough assimilation of relevant background knowledge [101]. Thus, the first step that took place in the Upper Paleolithic involved the social learning of existing knapping and carving techniques, including the abstract concept CARVE (from any suitable material).

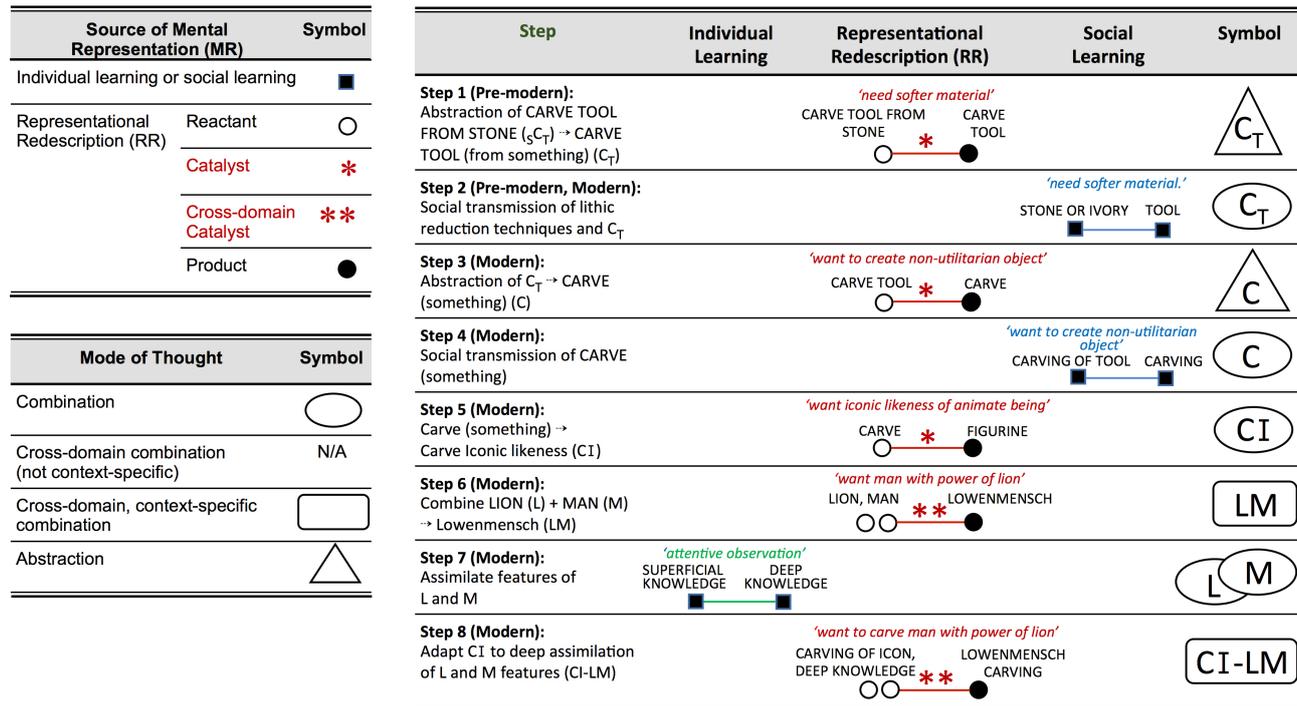


FIGURE 2. Top left: Sources of mental representations involved in the creation of the Hohlenstein-Stadel figurine, and the symbols used to depict them. Bottom left: Modes of thought and symbols used to depict them. Right: Steps involved in the creation of the figurine. Top two rows show the steps that occurred prior to the Upper Paleolithic; subsequent rows depict steps that took place during the Upper Paleolithic.

- (3) **Abstraction of CARVE TOOL to CARVE (something).** The next step was to extricate the concept CARVE TOOL from its conventional function of generating something utilitarian such as a hand axe. This resulted in a the abstract concept CARVE, which could now be applied in domains other than technology, such as art. This would have likely involved divergent, abstract thought. The existence of objects in bone, ochre, and ostrich eggshell with geometric engravings from southern Africa, dates this to at least 77,000 years ago [48].

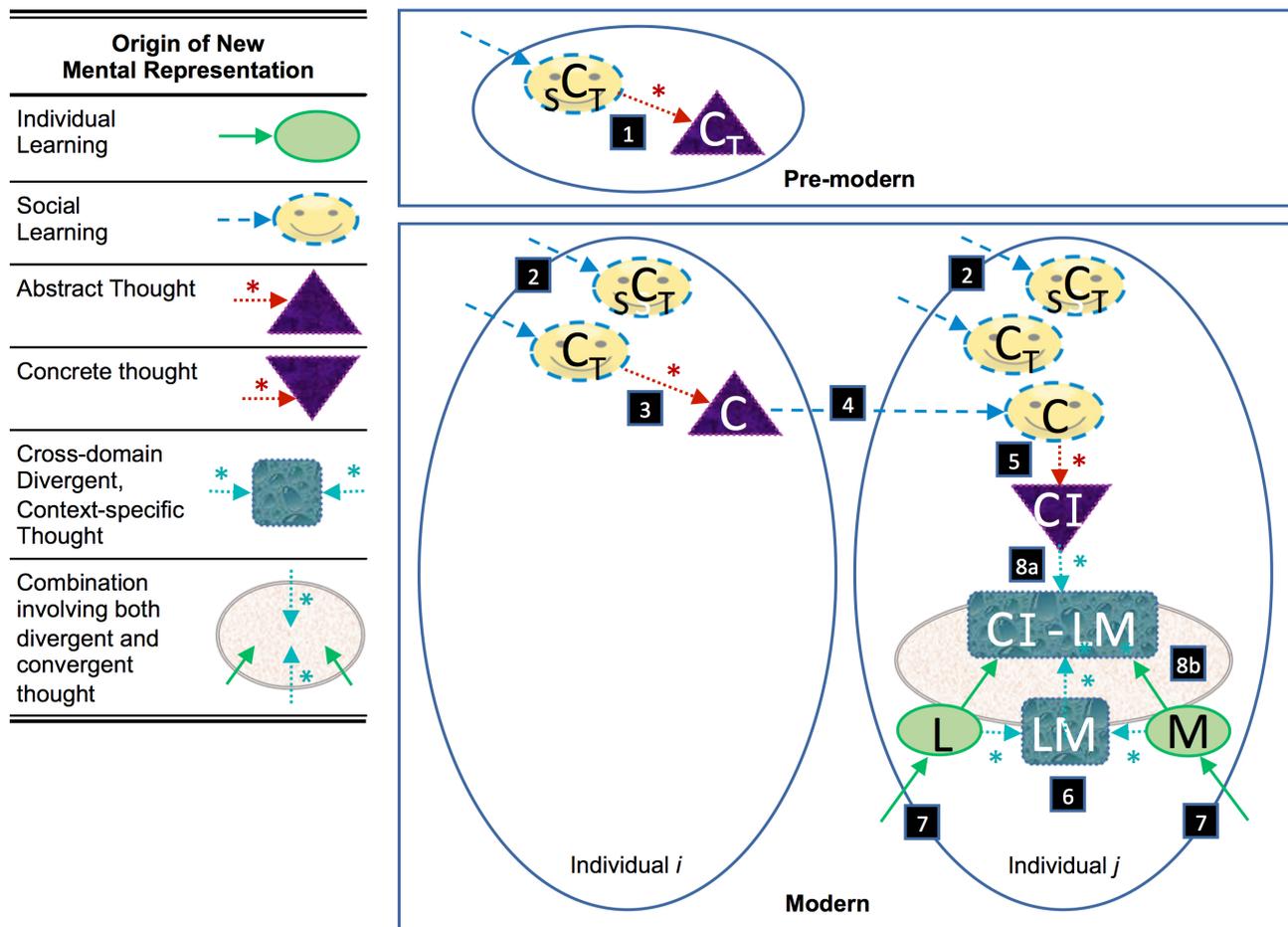


FIGURE 3. Steps culminating in creation of Hohlenstein-Stadel figurine. Meanings of symbols are defined in Figure 2 (see text for details).

- (4) **Social learning of the abstract concept CARVE (something)**. The carver of the Hohlenstein-Stadel figurine acquired the abstract concept CARVE (something) through social learning.<sup>11</sup>
- (5) **Apply CARVE (something) to the domain of figures (animal and human), yielding CARVE FIGURINE**. We may never know exactly what motivated the first artisan who took the step of carving an iconic likeness, a figurine. It may have been the

<sup>11</sup>We cannot know for certain that it was not invented independently (particularly given the distance between Hohlenstein-Stadel and southern Africa).

product of idle mind wandering. An alternative and perhaps more likely explanation is that it was shaped by a goal or desire, such as to (1) know the depicted subject more deeply, or (2) gain a sense of control or mastery over it, or (3) preserve a memory of it, or (4) have constant access to a feeling associated with it (such as the feeling of power associated with a lion). Whatever the motive, it involves taking the concept CARVE and applying it to a new domain, that of ANIMATE BEINGS.

- (6) **Combine LION HEAD with HUMAN BODY.** We also do not know what motivated this step. Like the previous step, it is possible that it was the product of idle mind wandering. It could be that by endowing a human body like ours with the head of a lion, the artisan hoped that those who held it would internalize the lion's power as their own. An alternative possibility is that it held some religious significance. Again, for the purpose of this model it is not essential to know which of these is correct, for whatever the underlying motive, this cross-domain combination would have required RR using divergent, context-specific thought.
- (7) **Assimilate features.** To carve an iconic figurine, knowledge of lithic reduction techniques is not sufficient; the artisan would have had to deeply absorb the physical characteristics of lions and humans through individual learning. We characterize this process as divergent because it involves assimilating the details and, potentially, any feelings they evoke. The artisan would then have used RR to creatively adapt this technical knowledge to the new task of rendering a figure in ivory.
- (8) **Carve figurine.** The actual carving of the figurine would have required RR in divergent mode to creatively adapt known carving techniques to the new task of rendering the detailed characteristics of the lion and human forms. Engagement in a tactile process meant that thought was concrete, ensuring that the features of the figurine were recognizably human or lion-like.

The entire mental trajectory through the spectrum of thought culminating in the creation of the Hohlenstein-Stadel figurine is depicted in Figure 4.

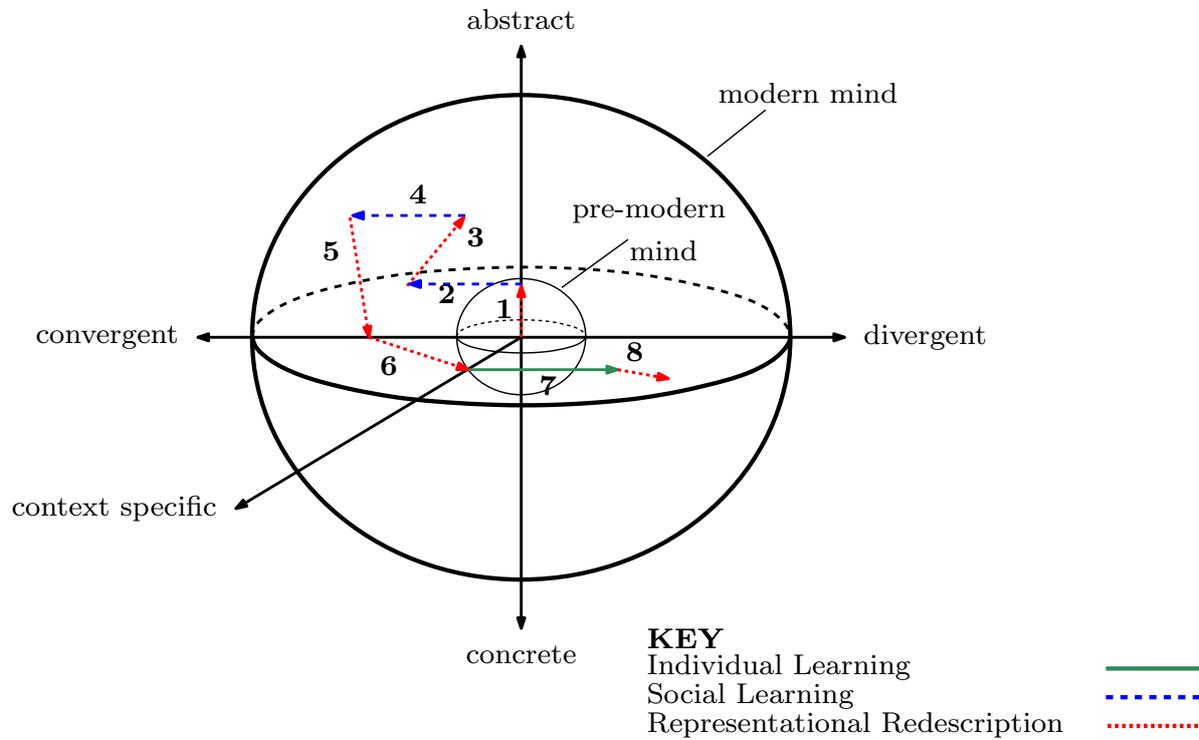


FIGURE 4. Trajectory through the spectrum of thought culminating in the Hohlenstein-Stadel figurine. Convergent-to-divergent is on the  $x$ -axis, abstract-to-concrete is on the  $y$ -axis, and degree of context-specificity is on the  $z$ -axis. Different combinations of these three variables comprise different ‘modes’ of thought (i.e., ways of navigating memory and processing information). Although the pre-modern mind could, to some degree, form abstractions, its thought trajectories used only a tiny portion of this space, as indicated by the small sphere. It was therefore restricted to a single mode of thought. The modern mind could engage in all combinations of these three variables, thereby engaging in many modes of thought, as indicated by the large sphere. The numbered arrows correspond to the eight steps listed in Table 5; thus, they depict how the mode of thought shifted over course of the figurine-making process.

## 6. THE TRANSITION TO COGNITIVE MODERNITY

In the pre-modern mind, information is thought to have been compartmentalized into domain-specific modules, and RR only operated within particular domains of human inquiry (e.g.,

‘tools’) [74]. Pre-modern cognition was largely (though not entirely) restricted to basic level categories—an intermediate level of abstraction—without the context-specific RR needed for cross-domain thinking. This was modeled by restricting RAFs to closed subsets of MRs that only reacted with other members of the same subset, resulting in RAF structure that was transient and fragmented [41, 42].

We now describe a simple mathematical model of the formation and persistence of cognitive RAFs culminating in behavioral/cognitive modernity and the cultural transition of the Upper Paleolithic. Let  $\mathcal{W} = \mathcal{W}(t)$  be a continuous measure of the scope of working memory of an individual at time  $t$  ( $t$  varies continuously over the lifetime of that individual). Cognitive processes make use of items obtained through individual learning, social learning, or RR, with items persisting in working memory for a short (but variable) time. As in [41], we model this as a non-deterministic process. If we let  $W = W(t) = \mathbb{E}[\mathcal{W}(t)]$  denote the expected (i.e., mean) value of  $\mathcal{W}(t)$  and assume a limit on the rate at which  $\mathcal{W}(t)$  can grow, this leads to the following nonlinear first-order equation:

$$(1) \quad \frac{dW}{dt} = -\mu W + \lambda \mathbb{E}[f(\mathcal{W})] + S.$$

Here  $S = S(t) \geq 0$  is a measure of information that is externally derived, either through individual learning or social learning at time  $t$ , the value  $1/\mu$  is the mean time that items remain in working memory, and  $\lambda$  describes the rate of RR reactions; and the function  $f$  is a concave increasing continuous function that satisfies  $f(0) = 0$ , and is asymptotically bounded (i.e.  $\lim_{x \rightarrow \infty} f(x) = K$  for some value  $K$ ). Simple default choices for  $f$  would be  $f(x) = \min\{x, K\}$  or  $f(x) = x(1 - x/C)$ , though we do not explicitly assume either of these here.

Crucially, the parameter  $\lambda$  also depends on total memory (a richer memory of knowledge and experiences allows more opportunity to catalyse RR reactions) and it is influenced by the three variables described above that we propose distinguish pre-modern from modern cognition.

More precisely, if  $\mathcal{M} = \mathcal{M}(t)$  denotes a continuous measure of the scope of total memory at time  $t$  and  $M = M(t) = \mathbb{E}[\mathcal{M}(t)]$  is the expected (mean) value of  $\mathcal{M}(t)$ , then  $\lambda$  is dependent on  $M$  (i.e.  $\lambda = \lambda(M)$ ). Thus, Equation (1) is coupled to the growth in  $M$ , and a simple model for the dynamics of expected total memory  $M$  is the first-order linear differential equation:

$$(2) \quad \frac{d(M - W)}{dt} = \nu W,$$

where  $\nu \in (0, 1)$  parameterises the extent to which items in working memory become encoded in long-term memory (which may also depend on time, as  $\nu$  may vary during the lifetime of the individual). The coupled (non-linear) system of Equations (1) and (2) leads to certain predictions. In particular, when  $\lambda$  lies below a critical threshold (dependent on the other parameters), CCPs do not form or persist, and thoughts are driven externally (individual learning or social learning). However, once  $\lambda$  passes this threshold, CCPs can form and persist indefinitely, even when the term  $S$  (in Eqn. (1)) drops to zero. The justification of these two claims and further mathematical details are provided in the Appendix.

The ability to shift between convergent and divergent thought, to consider the same item at multiple levels of abstraction, and to allow context to bias retrieval from memory by adjusting  $\gamma_D$ ,  $\gamma_A$  and  $\gamma_C$ , provide distinct and complementary mechanisms for  $\lambda$  to change. If the resulting new MRs are encoded in long-term memory, the positive dependence of  $\lambda$  on  $M$  provides more routes for catalysis of CCPs. This increases  $M - W$  from Eqn. (2), which, in turn, influences the dynamics of  $W$  by Eqn. (1).

The modern mind could carry out logical operations during convergent thought, and make new connections using divergent thought. Divergent thought could be biased toward a specific need by making thought more contextual. The modern mind could also shift up and down the hierarchy from concrete to abstract. By tuning the mode of thought along the three variables of the above multimodal spectrum to match the situation one is in, the modern mind acquired

the capacity to work out how elements of their world were interrelated, and where each element fitted with respect to the whole (i.e., the integrated internal model of the world, or *worldview*). The modern mind could now synthesize different domains of understanding into an integrated internal model of the world, using not only basic level concepts [86] but also higher or lower levels of abstraction, from fine-grained details to the ‘big picture’, as appropriate.

The worldview of the modern mind is a ‘metabolism’ in the sense that it has in place entropy-defying processes that maintain its organization. Like the protocell that constituted the earliest structure that could be said to be alive, the structure of the autocatalytic cognitive network as a whole is now maintained through the interactions amongst its parts. New experiences are interpreted, understood, and encoded in memory, in terms of existing cognitive structure already in place; thus, the structure of the memory becomes increasingly web-like.

## 7. COMPARISON WITH OTHER THEORIES

This model builds on the theory that the burst of creativity in the Paleolithic was due to the onset of contextual focus: the capacity to shift between divergent and convergent modes of thought [32]. That theory is superficially similar to the proposal that the distinguishing feature of human cognition is our capacity for dual processing [28, 79].<sup>12</sup> Our model builds on both the contextual focus and dual processing theories by positing that a single-variable spectrum of thought is insufficient to achieve an integrated internal model of the world.

---

<sup>12</sup>Dual processing posits that humans engage in not just a primitive implicit Type 1 mode for free association and fast ‘gut responses’, but also an explicit Type 2 mode for deliberate analysis. However, although dual processing makes the split between older, more automatic processes and newer, more deliberate processes, contextual focus theory posits that pre-modern thought was intermediate between two extremes (each valuable in different ways): a divergent mode based on relationships of correlation, and a convergent mode based on relationships of causation. Earlier hominids’ memories were coarser-grained, so there were fewer routes for meaningful associations, and less processing of previous experiences. Rather than convergent or divergent processing of previously assimilated material, there was greater tendency to focus on the here and now, so items in memory tended to remain in the same form as when they were originally assimilated. For a comparison of the divergent thought and dual processing theories see [88].

Our model is also consistent with Mithen's [73] theory that the transition was due to the connecting of domain-specific information processing modules, thereby enabling metaphorical thinking and cognitive fluidity: the capacity to combine ideas from different domains, fuse different knowledge processing techniques, or adapt a solution to one problem to a different problem. It is also consistent with Coolidge and Wynn's [19] theory that it was due to expanded working memory.<sup>13</sup> Conceptual fluidity and expanded working memory are underwritten by divergent thought but, as explained above, the capacity to engage in divergent thought without the capacity to control *how* divergent one's thinking is, would be perilous. Although it is not the focus of this paper, like [19], as well as [20] we are sympathetic with the view that genetic mutation was involved (see [40]).

Our proposal is consistent with the view that complex languages, symbolic representation, and myth lay at the heart of this transition [12, 13, 23]. However, we put the emergence of a persistent (i.e., stable) and integrated RAF network as central, with language both facilitating and being facilitated by this structure. Given evidence of recursive reasoning well before behavioral modernity, our framework is inconsistent with the hypothesis that the onset of recursive thought enabled mental time travel and cognitive modernity [21, 91]; nevertheless, the ability to shift through a multimodal spectrum of thought would have brought on the capacity to make vastly better use of it. The proposal that behavioral modernity arose due to onset of the capacity to model the contents of other minds, sometimes referred to as the 'Theory of Mind' [94], is somewhat underwritten by recursive RR, since the mechanism that allows for recursion is required for modeling the contents of other minds (though in this case the emphasis is on the social impact of recursion, rather than the capacity for recursion itself). Our proposal is also consistent with explanations for behavioral modernity that emphasize social-ecological factors [29, 98], but places these explanations in a broader framework by suggesting a mechanism that aided not just social skills but other skills (e.g., technological) as well.

---

<sup>13</sup>Working memory is just the part of memory that is, at any moment, working.

## 8. DISCUSSION AND CONCLUSIONS

Formal models exist of many aspects of human cognition, such as learning, memory, planning, and concept combination. However, there is little in the way of formal models of how they came to function together as an integrated whole, and how the unique cognitive abilities of *Homo sapiens* came about. RAF networks provide a means of addressing these questions. Building on earlier models of the cognitive transition underlying the earliest origins of human culture and the invention of the Acheulean hand axe, resulting in a transient autocatalytic structure, in this paper, we developed a model of the transition to a persistent, integrated RAF network. We proposed that rapid cultural change in the Middle-Upper Paleolithic required the ability to, not just recursively redescribe the contents of thought, but also tailor the ‘reactivity’ of thought to the current situation. This was accomplished through continuous, spontaneous tuning of three variables that concern not the content of thought *per se*, but how it is processed. The first involves shifting between convergent and divergent processing. The second involves shifting between concrete and abstract representations. The third involves biasing divergent processing according to a pressing need or context. Together, these enabled *Homo sapiens* to reflect on the contents of thought from different perspectives and at different levels of abstraction. This culminated in the crossing of a threshold to conceptual closure and the achievement of self-organizing autocatalytic semantic networks that spanned different knowledge domains, and routinely integrated new information by reframing it in terms of current understandings.

The model is highly simplified, and we do not know the precise details of the cognitive events modelled here took place (though the model does not hinge on these details). We hope that future versions will incorporate inhibition (in conjunction with the existing catalysis), as well as a more sophisticated representation of the interactions amongst MRs [1, 2] and a dynamic representation of context [56, 96]. There remains much work to be done on how cognitive RAFs

replicate and evolve (see [4] for informal suggestions in this regard) and on the developmental question of how persistent, integrated RAF networks emerge in the mind of a child.

We hope that future research will build on this direction by comparing the cultural RAF approach developed here with other standard semantic network approaches [8, 6, 10, 59, 68]. Although these standard semantic networks suffice for modeling semantic structure in individuals, we believe that the RAF approach will turn out to be superior because it distinguishes semantic structure arising through social or individual learning (modeled as food set items) from semantic structure *derived from* this pre-existing material (modeled as non-food set items generated through abstract thought processes that play the role of catalyzed reactions). This makes it feasible to model how cognitive structure emerges, and to trace lineages of cumulative cultural change step by step. It also frames this project within the overarching scientific enterprise of understanding how evolutionary processes (be they biological or cultural) begin, and unfold over time.

REFERENCES

- [1] D. AERTS, J. BROEKAERT, L. GABORA, AND S. SOZZO, *Generalizing prototype theory: A formal quantum framework*, *Frontiers in Psychology (Cognition)*, 7 (2016), p. 418, <https://doi.org/10.3389/fpsyg.2016.00418>.
- [2] D. AERTS, L. GABORA, AND S. SOZZO, *Concepts and their dynamics: A quantum theoretical model*, *Topics in Cognitive Science*, 5 (2013), pp. 737–772, <https://doi.org/10.1111/tops.12042>.
- [3] S. AMBROSE, *Paleolithic technology and human evolution*, *Science*, 291 (2001), pp. 1748–1753.
- [4] C. ANDERSSON AND P. TÖRNBERG, *Toward a macroevolutionary theory of human evolution: The social protocell*, *Biological Theory*, 14 (2019), pp. 86–102, <https://doi.org/10.1007/s13752-018-0313-y>.
- [5] M. AUBERT, P. SETIAWAN, A. OKTAVIANA, M. BRUMM, P. H. SULISTYARTO, E. W. SAPTOMO, AND ET AL., *Palaeolithic cave art in borneo*, *Nature*, 564 (2018), pp. 254–257.
- [6] A. BARONCHELLI, R. FERRER-I-CANCHO, R. PASTOR-SATORRAS, N. CHATER, AND M. H. CHRISTIANSEN, *Networks in cognitive science.*, *Trends in Cognitive Sciences*, 17 (2013), pp. 348–360, <https://doi.org/10.1016/j.tics.2013.04.010>.
- [7] L. W. BARSALOU, *Context-independent and context-dependent information in concepts*, *Memory & cognition*, 10 (1982), pp. 82–93.
- [8] R. E. BEATY, M. BENEDEK, P. J. SILVIA, AND D. L. SCHACTER, *Creative cognition and brain network dynamics*, *Trends in Cognitive Science*, 20 (2016), pp. 87–95.
- [9] R. A. BENTLEY, M. W. HAHN, AND S. J. SHENNAN, *Random drift and culture change*, *Proceedings of the Royal Society of London. Series B: Biological Sciences*, 271 (2004), pp. 1443–1450, <https://doi.org/10.1098/rspb.2004.2746>.
- [10] R. F. BETZEL AND D. S. BASSETT, *Generative models for network neuroscience: Prospects and promise*, *Journal of The Royal Society Interface*, 14 (2017), p. 20170623.
- [11] D. BICKERTON, *Language evolution: A brief guide for linguists*, *Lingua*, 117 (2007), pp. 510–526.
- [12] D. BICKERTON, *More than nature needs: Language, mind and evolution*, Harvard University Press, Cambridge, 2014.
- [13] D. BICKERTON AND E. SZATHMÁRY, *Biological foundations and origin of syntax*, MIT Press, Cambridge, 2009.

- [14] R. BOYD AND P. RICHERSON, *Culture and the evolutionary process*, University of Chicago Press, Chicago, 1988.
- [15] J. BROCKMEIER, *After the archive: Remapping memory*, *Culture and Psychology*, 16 (2010), pp. 5–35, <https://doi.org/10.1177/1354067X09353212>.
- [16] K. R. CABELL AND J. VALSINER, *The catalyzing mind: Beyond models of causality (Annals of Theoretical Psychology, Volume 11)*, Springer, Berlin, 2013, <https://doi.org/10.1007/978-1-4614-8821-7>.
- [17] L. L. CAVALLI-SFORZA AND M. W. FELDMAN, *Cultural transmission and evolution: A quantitative approach*, Princeton University Press, Princeton, NJ, 1981.
- [18] R. CAZZOLLA GATTI, B. FATH, W. HORDIJK, S. KAUFFMAN, AND R. ULANOWICZ, *Niche emergence as an autocatalytic process in the evolution of ecosystems*, *Journal of Theoretical Biology*, 454 (2018), pp. 110–117, <https://doi.org/10.1016/j.jtbi.2018.05.038>.
- [19] F. L. COOLIDGE AND T. WYNN, *Working memory, its executive functions, and the emergence of modern thinking*, *Cambridge Archaeological Journal*, 15 (2005), pp. 5–26.
- [20] M. CORBALLIS, *The origins of modernity: Was autonomous speech the critical factor?*, *Psychological Review*, 111 (2004), pp. 543–552, <https://doi.org/10.1037/0033-295X.111.2.543>.
- [21] M. CORBALLIS, *The recursive mind: The origins of human language, thought and civilization*, Princeton University Press, Princeton, NJ, 2011, <https://doi.org/10.2307/j.ctt6wpzjd>.
- [22] S. DAVIES, *Behavioral modernity in retrospect*, *Topoi*, (2019), pp. 1–12.
- [23] T. DEACON, *The symbolic species: The coevolution of language and the brain*, Norton, New York, 1997.
- [24] F. D'ERRICO, N. BARTOND, A. BOUZOUGGAR, H. MIENIS, D. RICHTER, AND P. . . . LOZOUET, *Additional evidence on the use of personal ornaments in the Middle Paleolithic of North Africa*, *Proceedings of the National Academy of Sciences, USA*, 106 (2009), pp. 16051–16056, <https://doi.org/10.1073/pnas.0903532106>.
- [25] M. ENQUIST, S. GHIRLANDA, AND K. ERIKSSON, *Modelling the evolution and diversity of cumulative culture*, *Philosophical Transactions of the Royal Society B: Biological Sciences*, 366 (2011), pp. 412–423, <https://doi.org/10.1098/rstb.2010.0132>.
- [26] P. ERDÖS AND A. RÉNYI, *On the evolution of random graphs*, *Publication of the Mathematical Institute of the Hungarian Academy of Sciences*, 5 (1960), pp. 17–61.
- [27] J. M. ERLANDSON, *The archaeology of aquatic adaptations: Paradigms for a new millennium*, *Journal of Archaeological Research*, 9 (2001), pp. 287–350.

- [28] J. EVANS, *Dual-process accounts of reasoning, judgment and social cognition*, Annual Review of Psychology, 59 (2008), pp. 255–278.
- [29] R. FOLEY AND C. GAMBLE, *The ecology of social transitions in human evolution*, Philosophical Transactions of the Royal Society B: Biological Sciences, 364 (2009), pp. 3267–3279, <https://doi.org/10.1098/rstb.2009.0136>.
- [30] L. GABORA, *Autocatalytic closure in a cognitive system: A tentative scenario for the origin of culture*, Psycology, 9 (1998), pp. [adap-org/9901002].
- [31] L. GABORA, *Conceptual closure: How memories are woven into an interconnected worldview.*, in Closure: Emergent Organizations and their Dynamics, G. Van de Vijver and J. Chandler, eds., no. 901 in Annual Review Series, Annals of the New York Academy of Sciences, 2000, pp. 42–53, <https://doi.org/10.1111/j.1749-6632.2000.tb06264.x>.
- [32] L. GABORA, *Contextual focus: A cognitive explanation for the cultural transition of the Middle/Upper Paleolithic.*, in Proceedings of the 25th Annual Meeting of the Cognitive Science Society, A. R. and H. D., eds., Lawrence Erlbaum Associates, Hillsdale, NJ, 2003, pp. 432–437.
- [33] L. GABORA, *Revenge of the 'neurds': Characterizing creative thought in terms of the structure and dynamics of human memory.*, Creativity Research Journal, 22 (2010), pp. 1–13.
- [34] L. GABORA, *An evolutionary framework for culture: Selectionism versus communal exchange*, Physics of Life Reviews, 10 (2013), pp. 117–145, <https://doi.org/10.1016/j.plrev.2013.03.006>.
- [35] L. GABORA, *How insight emerges in a distributed, content-addressable memory*, in The Cambridge handbook of the neuroscience of creativity, O. Vartanian and J. Jung, eds., MIT Press, Cambridge, MA, 2018, pp. 58–70.
- [36] L. GABORA, *Reframing convergent and divergent thought for the 21st century*, in Proceedings of the 2019 Annual Meeting of the Cognitive Science Society, A. Goel, C. Seifert, and C. Freska, eds., Cognitive Science Society, Austin, TX, 2019, pp. 1794–1800.
- [37] L. GABORA AND D. AERTS, *A model of the emergence and evolution of integrated worldviews*, Journal of Mathematical Psychology, 53 (2009), pp. 434–451, <https://doi.org/10.1016/j.jmp.2009.06.004>.
- [38] L. GABORA, S. LEIJNEN, T. VELOZ, AND C. LIPO, *A non-phylogenetic conceptual network architecture for organizing classes of material artifacts into cultural lineages*, in Proceedings of the 33rd annual meeting of the Cognitive Science Society, L. Carlson, C. Hölscher, and T. F. Shipley, eds., Cognitive Science Society, Psychology Press, 2011, pp. 2923–2928.

- [39] L. GABORA AND C. SMITH, *Two cognitive transitions underlying the capacity for cultural evolution*, Journal of Anthropological Science, 96 (2018), pp. 27–52, <https://doi.org/10.4436/jass.96008>.
- [40] L. GABORA AND C. SMITH, *Exploring the psychological basis for transitions in the archaeological record*, in Handbook of cognitive archaeology: Psychology in Prehistory, T. Henley, E. Kardas, and M. Rossano, eds., Routledge / Taylor and Francis, Abingdon, UK, 2019, ch. 12.
- [41] L. GABORA AND M. STEEL, *Autocatalytic networks in cognition and the origin of culture*, Journal of Theoretical Biology, 431 (2017), pp. 87–95, <https://doi.org/10.1016/j.jtbi.2017.07.022>.
- [42] L. GABORA AND M. STEEL, *Modeling a cognitive transition at the origin of cultural evolution using autocatalytic networks*, Cognitive Science, (in press).
- [43] T. L. GRIFFITHS, M. STEYVERS, AND J. B. TENENBAUM, *Topics in semantic representation*, Psychological Review, 114 (2007), pp. 211–244, <https://doi.org/10.1037/0033-295X.114.2.211>.
- [44] G. GRIMMETT AND D. STIRZAKER, *Probability and random processes (3rd ed.)*, Oxford University Press, 2001.
- [45] D. M. GYSI AND K. NOWICK, *Construction, comparison and evolution of networks in life sciences and other disciplines*, Journal of the Royal Society Interface, 17 (2020), p. 20190610, <https://doi.org/doi.org/10.1098/rsif.2019.0610>.
- [46] J. HAHN, *Kraft und aggression. Die botschaft der eiszeitkunst im Aurignacien Süddeutschlands?*, Archaeologica Venatoria, Tübingen, 1986.
- [47] J. A. HAMPTON, *Disjunction of natural concepts*, Memory and Cognition, 16 (1988), pp. 579–591, <https://doi.org/10.3758/BF03197059>.
- [48] T. HENLEY, M. J. ROSSANO, AND E. KARDAS, *Handbook of cognitive archaeology: A psychological framework*, Routledge / Taylor and Francis, Abingdon, UK, 2020, <https://doi.org/10.4324/9780429488818>.
- [49] C. J. HOLDEN AND R. MACE, *Spread of cattle led to the loss of matrilineal descent in africa: a coevolutionary analysis*, Proceedings of the Royal Society of London. Series B: Biological Sciences, 270 (2003), pp. 2425–2433, <https://doi.org/10.1098/rspb.2003.2535>.
- [50] W. HORDIJK, J. HEIN, AND M. STEEL, *Autocatalytic sets and the origin of life*, Entropy, 12 (2010), pp. 1733–1742, <https://doi.org/10.3390/e12071733>.
- [51] W. HORDIJK, S. A. KAUFFMAN, AND M. STEEL, *Required levels of catalysis for emergence of autocatalytic sets in models of chemical reaction systems*, International Journal of Molecular Science, 12 (2011), pp. 3085–3101, <https://doi.org/10.3390/ijms12053085>.

- [52] W. HORDIJK AND M. STEEL, *Detecting autocatalytic, self-sustaining sets in chemical reaction systems*, Journal of Theoretical Biology, 227 (2004), pp. 451–461, <https://doi.org/10.1016/j.jtbi.2003.11.020>.
- [53] W. HORDIJK AND M. STEEL, *Autocatalytic sets and boundaries*, J. Syst. Chem., 6:1 (2015).
- [54] W. HORDIJK AND M. STEEL, *Chasing the tail: The emergence of autocatalytic networks*, Biosystems, 152 (2016), pp. 1–10, <https://doi.org/10.1016/j.biosystems.2016.12.002>.
- [55] E. HOVERS, S. LANI, O. BAR-YOSEF, AND B. VANDERMEERSCH, *An early case of color symbolism: Ochre use by modern humans in Qafzeh cave*, Current Anthropology, 44 (2003), pp. 491–522.
- [56] M. W. HOWARD AND M. J. KAHANA, *A distributed representation of temporal context*, Journal of Mathematical Psychology, 46 (2002), pp. 269–299, <https://doi.org/10.1006/jmps.2001.1388>.
- [57] M. N. JONES AND D. J. K. MEWHORT, *Representing word meaning and order information in a composite holographic lexicon*, Psychological Review, 114 (2007), pp. 1–37, <https://doi.org/10.1037/0033-295X.114.1.1>.
- [58] P. KANERVA, *Hyperdimensional computing: An introduction to computing in distributed representations with high-dimensional random vectors*, Cognitive Computation, 1 (2009), pp. 139–159, <https://doi.org/10.1007/s12559-009-9009-8>.
- [59] E. A. KARUZA, S. L. THOMPSON-SCHILL, AND D. S. BASSETT, *Local patterns to global architectures: influences of network topology on human learning*, Trends in Cognitive Sciences, 20 (2016), pp. 629–640.
- [60] S. A. KAUFFMAN, *Autocatalytic sets of proteins*, Journal of Theoretical Biology, 119 (1986), pp. 1–24, <https://doi.org/10.3390/ijms12053085>.
- [61] S. A. KAUFFMAN, *The origins of order*, Oxford University Press, 1993.
- [62] C. KIND, N. EBINGER-RIST, S. WOLF, T. BEUTELSPACHER, AND K. WEHRBERGER, *The smile of the lion man. Recent excavations in Stadel cave (Baden-Württemberg, southwestern Germany) and the restoration of the famous upper palaeolithic figurine*, Quartär, 61 (2014), pp. 129–145.
- [63] P. J. KWANTES, *Using context to build semantics*, Psychonomic Bulletin & Review, 12 (2005), pp. 703–710, <https://doi.org/10.3758/BF03196761>.
- [64] P. LIEBERMAN, *Language did not spring forth 100,000 years ago*, PLoS Biology, 13 (2015), pp. e1002064, [doi.org/10.1371/journal.pbio.1002064](https://doi.org/10.1371/journal.pbio.1002064).
- [65] P. M., *Wired for culture: The natural history of human cooperation*, Penguin, London, UK, 2012.

- [66] S. MCBREARTY AND A. BROOKS, *The revolution that wasn't: A new interpretation of the origin of modern human behavior*, *Journal of Human Evolution*, 39 (2000), pp. 453–563.
- [67] J. M. MCCLELLAND, *Memory as a constructive process: The parallel distributed processing approach*, in *The memory process: Neuroscientific and humanistic perspectives*, S. Nalbantian, P. M. Matthews, and J. L. McClelland, eds., MIT Press, Cambridge, MA, 2011, pp. 129–155.
- [68] J. D. MEDAGLIA, M. E. LYNALL, AND D. S. BASSETT, *Cognitive network neuroscience*, *Journal of Cognitive Neuroscience*, 27 (2015), pp. 1471–1491.
- [69] S. A. MEDNICK, *The associative basis of the creative process*, *Psychological Review*, 69 (1962), pp. 220–232.
- [70] P. MELLARS, *Going east: New genetic and archaeological perspectives on the modern human colonization of Eurasia*, *Science*, 313 (2006), pp. 796–800.
- [71] A. MESOUDI, A. WHITEN, AND K. N. LALAND, *Towards a unified science of cultural evolution*, *Behavioral and Brain Science*, 29 (2006), pp. 329–347, <https://doi.org/10.1017/S0140525X06009083>.
- [72] S. MITHEN, *The prehistory of the mind: A search for the origins of art, religion, and science*, Thames and Hudson, London, UK, 1996.
- [73] S. MITHEN, *Creativity in human evolution and prehistory*, Routledge, London, UK, 1998.
- [74] S. MITHEN, *Ethnobiology and the evolution of the human mind*, *Journal of the Royal Anthropological Institute*, 12 (2006), pp. S45–S61.
- [75] E. MOSSEL AND M. STEEL, *Random biochemical networks and the probability of self-sustaining autocatalysis*, *Journal of Theoretical Biology*, 233 (2005), pp. 327–336, <https://doi.org/10.1016/j.jtbi.2004.10.011>.
- [76] J. MULVANEY AND J. KAMMINGA, *Prehistory of Australia*, Smithsonian Institution Scholarly Press, Washington, 1999.
- [77] M. MUTHUKRISHNA, M. DOEBELI, M. CHUDEK, AND J. HENRICH, *The cultural brain hypothesis: How culture drives brain expansion, sociality, and life history*, *PLoS Computational Biology*, 14 (2018), p. e1006504, <https://doi.org/10.1371/journal.pcbi.1006504>.
- [78] S. NELSON, *Diversity of the Upper Palaeolithic Venus figurines and archaeological mythology*, *Archeological Papers of the American Anthropological Association*, 2 (2008), pp. 11–22.
- [79] B. A. NOSEK, *Implicit-explicit relations*, *Current Directions in Psychological Science*, 16 (2007), pp. 65–69.

- [80] D. N. OSHERSON AND E. E. SMITH, *On the adequacy of prototype theory as a theory of concepts*, *Cognition*, 9 (1981), pp. 35–58, [https://doi.org/10.1016/0010-0277\(81\)90013-5](https://doi.org/10.1016/0010-0277(81)90013-5).
- [81] M. OTTE, *The management of space during the Paleolithic*, *Quaternary International*, 247 (2012), pp. 212–229.
- [82] A. PIKE, D. HOFFMANN, M. GARCÍA-DIEZ, P. PETTITT, J. ALCOLEA, R. DE BALBIN, AND J. ZILHÃO, *U-series dating of Paleolithic art in 11 caves in Spain*, *Science*, 336 (2012), pp. 1409–1413.
- [83] M. PORR, *Palaeolithic art as cultural memory: A case study of the Aurignacian art of Southwest Germany*, *Cambridge Archaeological Journal*, 20 (2010), pp. 87–108.
- [84] A. POWELL, S. SHENNAN, AND M. G. THOMAS, *Late Pleistocene demography and the appearance of modern human behavior*, *Science*, 324 (2009), pp. 1298–1301.
- [85] R. RAPPAPORT, *Ritual and religion in the making of humanity*, Cambridge University Press, Cambridge, 1999.
- [86] E. ROSCH, C. B. MERVIS, W. D. GRAY, D. M. JOHNSON, AND P. BOYES-BRAEM, *Basic objects in natural categories*, *Cognitive Psychology*, 8 (1976), pp. 382–439.
- [87] J. M. SMITH AND E. SZATHMARY, *The major transitions in evolution*, Oxford University Press, Oxford, UK, 1997.
- [88] P. SOWDEN, A. PRINGLE, AND L. GABORA, *The shifting sands of creative thinking: Connections to dual process theory*, *Thinking & Reasoning*, 21 (2015), pp. 40–60.
- [89] M. STEEL, W. HORDIJK, AND J. C. XAVIER, *Autocatalytic networks in biology: structural theory and algorithms*, *Journal of the Royal Society Interface*, 16 (2019), p. rsif.2018.0808, <https://doi.org/10.1098/rsif.2018.0808>.
- [90] D. STOUT, N. TOTH, K. SCHICK, AND T. CHAMINADE, *Neural correlates of early stone age toolmaking: technology, language and cognition in human evolution*, *Philosophical Transactions of the Royal Society B: Biological Sciences*, 363 (2008), pp. 1939–1949, <https://doi.org/10.1098/rstb.2008.0001>.
- [91] T. SUDDENDORF, D. R. ADDIS, AND M. C. CORBALLIS, *Mental time travel and the shaping of the human mind*, *Philosophical Transactions of the Royal Society B: Biological Sciences*, 364 (2009), p. 1317, <https://doi.org/10.1098/rstb.2008.0301>.
- [92] I. TATTERSALL, *The origin of the human capacity*, American Museum of Natural History, 1998.
- [93] J. J. TEHRANI AND F. RIEDE, *Towards an archaeology of pedagogy: Learning, teaching and the generation of material culture traditions*, *World Archaeology*, 40 (2008), pp. 316–331.

- [94] M. TOMASELLO, *A natural history of human thinking*, Harvard University Press, Cambridge MA, 2014.
- [95] V. VASAS, C. FERNANDO, M. SANTOS, S. KAUFFMAN, AND E. SZATHMÁRY, *Evolution before genes*, *Biology Direct*, 7 (2012).
- [96] T. VELOZ, L. GABORA, M. EYJOLFSON, AND D. AERTS, *Toward a formal model of the shifting relationship between concepts and contexts during associative thought*, in *Proceedings of the Fifth International Symposium on Quantum Interaction*, D. Song, M. Melucci, I. Frommholz, P. Zhang, L. Wang, and S. Arafat, eds., Springer, Cognitive Science Society, 2011, pp. 25–34, [https://doi.org/10.1007/978-3-642-24971-6\\_4](https://doi.org/10.1007/978-3-642-24971-6_4).
- [97] T. VELOZ, I. TEMPKIN, AND L. GABORA, *A conceptual network-based approach to inferring cultural phylogenies*, in *Proceedings of the 34th annual meeting of the Cognitive Science Society*, N. Miyake, D. Peebles, and R. P. Cooper, eds., Cognitive Science Society, Austin TX, 2012, pp. 2487–2492.
- [98] A. WHITEN, *The scope of culture in chimpanzees, humans and ancestral apes*, *Philosophical Transactions of the Royal Society, Series B.*, 366 (2011), pp. 997–1007.
- [99] A. WHITEN, *Cultural evolution in animals*, *Annual Review of Ecology, Evolution, and Systematics*, 50 (2019), pp. 1–22.
- [100] J. C. XAVIER, W. HORDIJK, S. KAUFFMAN, S. M., AND W. F. MARTIN, *Autocatalytic chemical networks at the origin of metabolism*, *Proceedings of the Royal Society of London. Series B: Biological Sciences*, 287 (2020), p. 20192377.
- [101] X. ZHANG, D. WANG, AND T. WANG, *Inspiration or preparation? Explaining creativity in scientific enterprise*, in *Proceedings of the 25th ACM International on Conference on Information and Knowledge Management*, ACM, 2016, pp. 741–750.
- [102] J. ZILHÃO, *Modernity, behavioral*, in *The International Encyclopedia of Anthropology*, H. Callan, ed., American Cancer Society, 2018, pp. 1–9, <https://onlinelibrary.wiley.com/doi/abs/10.1002/9781118924396.wbiea1787>.

9. APPENDIX: MATHEMATICAL DETAILS AND JUSTIFICATION OF PREDICTIONS BASED ON EQUATIONS (1) AND (2).

We begin with some preliminary remarks. In the following arguments, we treat  $\lambda$  as a constant over the short time-frame considered in the dynamics of CCPs, since the dependence of  $\lambda$  on  $M$  applies over considerably longer time-scales. Moreover, in treating  $\lambda$  as a constant and setting  $S = 0$ , Eqn (1) is not of the form  $\Phi(W, dW/dt) = S(t)$ , since  $\mathbb{E}[f(\mathcal{W})]$  is not, in general, a function of  $W$ . For example, for  $f(x) = x(1 - x/K)$ , Eqn (1) becomes:

$$\frac{dW}{dt} = -\mu W + \lambda \left( W \left( 1 - \frac{W}{K} \right) - \frac{\mathbb{V}(\mathcal{W})}{K} \right) + S,$$

where  $\mathbb{V}(\mathcal{W})$  is the variance of  $\mathcal{W}$  at time  $t$ . Note also that the dynamics of  $\mathcal{W}(t)$  is not determined by the behaviour of  $W(t)$ ; the latter just represents the expected (average) value of the former.

Returning to the first prediction of this model, observe that:

$$(3) \quad \mathbb{E}[f(\mathcal{W})] \leq f(\mathbb{E}[\mathcal{W}]) = f(W) \leq f'(0) \cdot W.$$

The first inequality in (3) is by Jensen's inequality for the concave function  $f$  (see e.g. [44]). The second inequality also uses the concavity of  $f$  together with the condition  $f(0) = 0$ . Thus if  $\lambda < \mu/f'(0)$ , we have:

$$\frac{dW}{dt} \leq -cW + S$$

for  $c = (\mu - \lambda f'(0)) > 0$ . Consequently, once  $S$  declines to zero, so too does  $W$ , and by the Markov inequality (see e.g. [44]), we have:

$$\mathbb{P}(\mathcal{W}(t) > \epsilon) \leq \frac{\mathbb{E}[\mathcal{W}(t)]}{\epsilon} = \frac{W(t)}{\epsilon} \rightarrow 0$$

as  $t$  increases, which establishes the first prediction.

Now suppose that  $\lambda > \mu/f'(0)$ . Let  $\beta = \mu(1 + \eta)/\lambda < f'(0)$  for a sufficiently small  $\eta > 0$ . By the concavity of  $f$ , it follows that, for some  $\gamma \geq 0$ , we have:

$$(4) \quad f(x) \geq \beta x, \text{ for all } x \in [0, \gamma],$$

since the line  $y = \beta x$  and the function  $y = f(x)$  both pass through the origin; however, the latter function has a strictly greater slope at the origin.

Now suppose that  $\mathcal{W}(t_1) = w$ , where  $w \in (0, \gamma)$ . Considering the process moving forward from time  $t_1$ , with the initial condition  $\mathcal{W}(t_1) = w$  (and thus  $W(t_1) = w$ ) at  $t = t_1$ , we then have:

$$\frac{dW}{dt} = -\mu w + \lambda f(w) + S \geq -\mu w + \lambda \beta w = \eta w > 0,$$

where the first inequality is from (4) together with  $S(t_1) \geq 0$ . In summary, when  $\lambda$  passes above the threshold  $\mu/f'(0)$  and  $\mathcal{W}(t)$  is small but non-zero (and even with  $S = 0$ ), the expected value of  $\mathcal{W}(t)$  begins to increase, due to CCPs in working memory.