# A Cognitive Model Fleshes out Kahneman's Fast and Slow Systems

Usef Faghihi[1], Clayton Estey[2], Ryan McCall[2], Stan Franklin[2]

[1]*Cameron University, OK, USA.*
[2]*University of Memphis, TN, USA.*
*ufaghihi@cameron.edu,* {cestey, rmccall, stan.franklin }*@memphis.edu*

**Abstract**

Daniel Kahneman (2011) posits two main processes that characterize thinking: "System 1" is a fast decision making system responsible for intuitive decision making based on emotions, vivid imagery, and associative memory. "System 2" is a slow system that observes System 1's outputs, and intervenes when "intuition" is insufficient. Such an intervention occurs "when an event is detected that violates the model of the world that System 1 maintains" (Kahneman, 2011, p. 24). Here, we propose specific underlying mechanisms for Kahneman's Systems 1 and 2, in terms of the LIDA model, a broad, systems-level, cognitive architecture (Stan Franklin et al., 2014). LIDA postulates that human cognition consists of a continuing, overlapping iteration of cognitive cycles, each a cognitive "atom," out of which higher-order processes are built. In LIDA terms, System 1 employs consciously mediated action selection in which a stimulus is acted upon within one or two cognitive cycles. In contrast, System 2, which LIDA posits to operate according to James' ideomotor theory (William James, 1950) , requires more cognitive cycles in its deliberative decision making. Thus, we suggest that System 2 employs multiple occurrences of System 1 in its operation. To test the proposed mechanisms, we perform an in silico experiment using a LIDA-based software agent.

*Keywords:* Learning Intelligent Distribution Agent (LIDA,) Kahneman's fast and slow systems, cognitive architecture, consciously mediated action selection, deliberative decision making

## 1  Introduction

As human beings, we interact with our environment and integrate implicit and explicit knowledge for decision making. Many researchers in psychology and neuroscience have suggested models for decision making (Pfeiffer, Whelan, & Martin, 2000; Scott & Bruce, 1995; Shiv & Fedorikhin, 1999). However, the main target of these models is the functional level process of decision making, rather than their underlying mechanisms. Some other questions that have been discussed in the literature are: Would our brains use implicit or explicit knowledge or a combination of both to make a decision (Peters &

Levin, 2008)? How can we explain the underlying processes that affect errors of judgment? Or choices under pressure? What are the prerequisites, such as training, required to have a clear enough understanding of a given situation to make appropriate decisions? How can professional training improve an expert's decision making?

Kahneman's Thinking, Fast and Slow (TFS) (2011) discusses all the aforementioned questions. The book postulates that two systems are responsible for decision making; namely, "System 1," a fast system, which is responsible for intuitive decisions based on emotions, vivid images and associative memory, and "System 2," a slow system, which observes System 1's output, and intervenes when the output is considered to infringe on decision making rules. In order to pay attention to the output of System 1, System 2 requires a great deal of extra energy, and also can sometimes be "lazy" (see below). In order to understand the decisions we make, it is necessary to understand such mental processes underlying decision making.

The Learning Intelligent Distribution Agent (LIDA) model is a broad, systems-level, cognitive architecture that attempts to model how minds work (Stan Franklin et al., 2014). LIDA conceptually and computationally implements the Global Workspace Theory of Baars (1988). Global Workspace Theory (GWT) is perhaps the most widely accepted psychological and neurobiological theory of the role of consciousness in cognition (Bernard J. Baars, 2002). LIDA postulates that human cognition consists of a continuing, overlapping iteration of cognitive cycles, each acting as a cognitive "atom," out of which higher-level cognitive processes are built. We suggest that Kahneman's fast System 1 is characterized by decision making in one to three cognitive cycles, while the slower System 2, implementing James' ideomotor theory of decision making (William James, 1950), employs more than 3 cognitive cycles in its deliberative decision making.

In the next section of this paper, we will give a brief review of Thinking, Fast and Slow. We then present what we take to be the conceptual core of the book, and describe how the LIDA model and its cognitive mechanisms implement that core. In particular, we will describe how Kahneman's two systems are implemented by two forms action selection in LIDA. To begin validating LIDA as a model of the underlying processes of Systems 1 and 2, using a LIDA agent, we replicate some of the experiments described by Kahneman in TFS.


## 2  Kahneman's Fast and Slow Systems

The focus of Thinking, Fast and Slow (TFS) is on two systems of decision making. "System 1" is the fast system, which is responsible for intuitive decisions based on emotions, vivid imagery and associative memory. "System 2" is the slow system, which observes System 1's output, and intervenes when the output is considered to either infringe on more rational decision-making rules, or when an agent's intuition[*] is insufficient in handling a situation. For example, a person feels thirsty and immediately reaches for a glass of water on the table (System 1 in action). Or, he considers having a beer instead, but thinks that it's too early in the morning for that, and decides to drink orange juice instead (System 2 in action). In the former example, the person had an immediate need, which could be satisfied by an action based on intuition alone. No alternative possibility needed to be accounted for. In the latter, it was initially intuitive for the person to consider having a nice, refreshing beer. However, the time of day was sufficient context to provoke an inconsistency involving this initial impression and the long-term consequences of consuming beer so early. In order to address inconsistencies between one's intuitive impressions about how to decide and evaluate, and aspects of our situation conflicting with those impressions, System 2 becomes vigilant and resolves the issue.

Relevant to the above, humans have to adapt to both fast-paced, chaotic environments and to slower-paced, more stable environments. Evolutionarily, we needed quick, heuristic decision making when

---

[*] The act or faculty of knowing or sensing without the use of rational processes; immediate cognition (*Freedictionary.com*)

there was no time for long-term planning in such fast-paced situations. Dealing with present opportunities for short-term gain is an example, much like the beer-versus-water dilemma above. Another is when there was no time to deliberate in the rational sense (System 1) (e.g., with an immediate threat). However, we also need a slower system for long-term planning, which would best fit "slower," more stable situations. Examples would include planning for the winter, when there is less food available, or choosing what to order in a restaurant. We would have to make rational (System 2) predictions, and reason about multiple consequences, all over extended periods of time. System 1 is incapable of such long term processing. Likewise, it would be disastrous if System 2 interfered when time was of the essence. It is imperative that these two systems collaborate to deal with these conflicting demands. Kahneman's account of Systems 1 and 2 involves many circumstances in which such collaboration occurs, whether culturally enforced or built-in by evolution.

TFS, discusses attention and effort, cognitive ease and strain, norms, surprise and causes, causal versus statistical reasoning, expert intuition, intuition verses formulas, associative coherence, attribution substitution, the availability heuristic, availability cascades, the affect heuristic, the halo effect, the representativeness heuristic, and the anchoring effect. It also mentions some of Slovic's work (2000) regarding human judgment of risk. In this paper, from the above list, we will give a brief description of attention and effort, cognitive ease, expert intuition, associative coherence, the availability heuristic, the affect heuristic, the representativeness heuristic, and the anchoring effect and then in the Kahneman's Systems 1 and 2 à la LIDA section we will discuss them in the subsection below.

**Attention and effort.** The effort described by Kahneman refers to what a subject is doing, instead of what is happening to him (Kahneman, 1973). According to Kahneman, attention is made up of various components such as selective and focal attention. Selective attention is the purposeful direction of effort to a particular mental task. Kahneman assumed that focal attention to an object increases sensitivity to matters related with that object (Kahneman, 1973). To explain attention and effort Kahneman postulates that, System 1 continuously monitors what is going on outside and inside the mind. Then, it continuously generates assessments of various aspects of the situation with little or no effort. System 2 receives questions from System 1 or it generates them, and directs attention and searches memory to find the answers. Kahneman and others (Bijleveld, Custers, & Aarts, 2009; Marshall, 2002; Peavler, 1974) postulate that the pupils of the eyes are sensitive indicators and good measures of mental effort. For example, they dilate substantially when people multiply two-digit numbers in their heads. In other experiments, when subjects were exposed to more digits than they could remember, their pupils stopped dilating or actually shrank. To describe the adoption and termination of task sets, Kahneman uses "executive control." Kahneman describes executive control using the following example: You are asked to count all occurrences of the letter x in this page and at the end of the page you are asked to count all "." in the next page. The second task is harder because you need to overcome the recently developed tendency to focus attention on the letter x.

**Cognitive ease†.** To illustrate cognitive ease, Kahneman suggests that a message can be made more convincing by first maximizing legibility. For instance, a study on intuitive judgment presented its participants with the following two false statements:
**"Adolf Hitler was born in 1892**.
Adolf Hitler was born in 1887 (page 64)."
Participants were more likely to choose the statement in bold type. In the state of "cognitive ease" our rational, "thinking" (System 2) brain is put on hold due to being distracted by our intuition (System 1). In the above example with two false statements, System 1 produces the feeling of familiarity, and System 2 accepts that impression. In the case of System 2 being lazy or ignorant, if you cannot recall

---

† "A sign that things are going well—no threats, no major news, no need to redirect attention or mobilize effort (Page 61)." (Kahneman, 1973)

the origin of a statement, and have no clue with which to associate it to other things, you have no choice but to go with the sense of cognitive ease. Kahneman suggests additional rhetorical strategies exploiting cognitive ease with statements, such as simplifying the language and making them memorable.

**Expert Intuition: When Can We Trust It?** A fundamental question Kahneman asks is: When can we trust an expert's intuition? He suggests that intuition plays an important role in the decision making process of experts. Kahneman defines expertise in a domain as a large collection of "miniskills," be it for firefighters, chess masters, physicians, nurses, athletes, etc. For instance, chess masters can assess a complicated position intuitively, but years of practice are needed to achieve that level of skill. Kahneman mentions that to become a chess master you need to practice for at least 10,000 hours. Accordingly, he argues that we can trust an expert's intuition when the environment is regular, predictable, and the expert has enough opportunity to practice and get immediate feedback. For example, a military commander will identify a narrowing ravine as a threat in an intuitive way due to his training, which taught him to associate them with ambushes. However, we cannot trust an expert's intuition in the absence of stable regularities in the environment, for example, an economist's tips on the stock market (see below).

Kahneman describes how the influence of emotions can bring an expert's decision to erroneous conclusions. In such situations human decision making is made by the emotionally influenced System 1, which passes the careless and permissive review of System 2. Furthermore, System 1 often creates correlation where there is none. For instance, forecasting the stock market in which Google's shares jump above $1000 because of the recent government shutdown. The shares' prices also can change sharply if the US government, all of a sudden, decides to make war with other countries. Thus, no one can claim to have sufficient information with which to predict stock markets.

**Associative Coherence.** Associative coherence occurs when an agent interprets the salient content of its current situation as being consistent with its beliefs, irrespective of that content's reliability or validity. In associative coherence only evidence not conflicting with the agent's beliefs is considered. Evidence conflicting with its current beliefs requires deliberation to deal with, and such content requires reasoning about multiple possibilities "at once" (deliberation). Deliberation requires System 2, and is outside the capacities of System 1. Yet, System 1's actions dominate an agent's processing, with associative coherence being the core of intuition à la System 1. Kahneman gives an example of associative coherence in Chapter 4, "The Associative Machine." If one is presented with the two words "Banana" and "Vomit" juxtaposed, one may experience all sorts of memories depending on your past. The unconscious aspects of our mind immediately assume causality/categorical relatedness is at work, and you may be perceptually primed for a coherent thought like "The banana made me vomit." The words evoke memories of an event, which in turn evokes emotions, which in turn evokes a gut reaction. Within a short time, your mind made sense of two known words presented in a very novel way, and it was probably effortless to you. This seemingly effortless attempt of the mind to fill in the gaps of its understanding by creating a reality more coherent than what is presented before it; that is the nature of associative coherence, and it was entirely System 1.

**Affect Heuristic.** The affect heuristic occurs when an agent evaluates a stimulus based on its "likes and dislikes." Or for more complicated situations, it occurs when an agent either weighs benefits more than costs (positive), or weighs costs more than benefits (negative) in its decision making. An example of the former is when you really like how a car looks, so you ignore its poor gas mileage. An example of the latter is disliking a competing cognitive architecture since it doesn't account for consciousness, so much so that its other virtues are ignored. These examples are extreme versions of the affect heuristic just to illustrate the point. In reality, this heuristic accounts for any non-neutral weighting of benefit vs. cost. One of Kahneman's examples was the Hsee et al. study (Hsee, Weber, & Welch, 2001). In one of the three contexts which the study explored, participants rated the quality of a 24-piece dinner set as better than a set with the same pieces, plus 16 additional pieces including 9 broken ones. This is an

example of the "less-is-better effect," which is based on the affect heuristic (Hsee et al., 2001). In this example, detailed summation of each piece would have required too much effort, so System 1 engaged in a lump, emotive perception/ reaction of the material. This resulted in ignoring the fact that the second set had the same good pieces plus more, because of the inclusion of a few flaws.

**Availability Heuristic:** The availability heuristic is used when an agent judges the plausibility/potency/importance of a category from a retrieved memory based on either the ease/fluency of its retrieval (System 1 usage) or on the frequency of certain content that is retrieved (System 2 usage). For System 1, Kahneman gives an example in which a group of psychologists in the early 1990's, led by Norbert Schwartz, investigated how fluency of retrieval related to self-assessments of assertiveness. What they found was that the more instances of assertiveness they were asked to list, the less assertive they viewed themselves to be (Schwarz et al., 1991). This is because as the number of instances to be listed increased, the fluency of producing examples decreased. Because the fluency was less than they expected, they associated the negativity of that with their own self-assessment. The participants were thus victims of the availability heuristic a la System 1. For System 2, Kahneman gives an example of the same researchers recruiting two groups of students whose task were to recall instances of their routines influencing their cardiac health. The two groups were those who had a family history of heart disease and those who did not. The latter group produced the same effects as from the assertiveness study. The former group, however, had the opposite effect. Because it was about them, and due to their family history, they were put into a higher state of vigilance than would be expected from the other group. As they recalled more instances of safe behavior, they felt safer. When they recalled more instances of dangerous behavior, they felt more at risk. In their vigilance, they did not "go with the flow" as what would happen if they were in the other group. Instead, they engaged in deliberation to evaluate patterns in the content, and assessed long term consequences based on such. This is the difference between the two systems and how they implement the availability heuristic. An earlier analysis of the mechanics of the availability according to the LIDA model appears elsewhere (Stan Franklin, B. J. Baars, Uma Ramamurthy, & M Ventura, 2005).

**Representativeness Heuristic:** When an agent judges the likelihood of an event, the judgment is affected by how much the event represents/resembles a parent category or process. To put it another way: an agent will often appeal to categorical/causal stereotypes instead of base rates and statistical regularities. One of the examples Kahneman gives is in Chapter 14, "Tom W's Specialty." Tom W is a graduate student with certain personality traits (e.g. hard-working, detail oriented, little sympathy for others, etc.), and the task of a participant is to judge the likelihood, from 1 to 9 that Tom W is a student in certain disciplines. The judged likelihood changes depending on the personality description given because the descriptions represent a stereotype of how people in certain fields behave. When participants listed the disciplines in descending order of likelihood, they didn't consider the number of students in each field as a base rate. The disciplines with the greatest number of people (e.g. humanities and education, social science and social work) were rated as least likely because of the way Tom W was described. In the absence of knowledge of base rates, we often substitute the statistical term "likelihood" for the causal/stereotyping term "plausibility", and evaluate that instead. It is often the case that people with a lazy (possibly in this case, ignorant as well) System 2 make causal/categorical membership judgments based on insufficient evidence. This is also seen when people make sweeping generalizations about a person's abilities from very few observations, where there is little chance to control for lucky or unlucky performance. These biases occur because while System 1 is sufficient for causal/categorical inference, reasoning about statistical patterns over time requires the deliberative System 2.

**Anchoring Effect:** The anchoring effect is being used when an agent makes judgments using available information, such that the judgment does not deviate very much from the magnitude (e.g. a number) or relatedness (e.g. qualitative content) of the available information. This available information

is called the "anchor." Both System 1 (through subliminal priming of a reference point to be used in a later, unrelated task) and System 2 (through deliberate adjusting from a reference point, due to uncertainty) use the anchoring effect. For the former, Kahneman gives an example of two German psychologists, Mussweiler and Strack, asking participants:

> "Is the annual mean temperature in Germany higher or lower than 20 degrees Celsius (68 F)?" or "Is the annual mean temperature in Germany higher or lower than 5 degrees Celsius (41 F)?"

Then they were briefly shown words they were asked to identify. They found that 68 F made identifying summer words such as "sun" and "beach" easier, while 41 F made identifying winter words such as "frost" and "ski" easier (Mussweiler & Strack, 2000). In the case of a deliberative anchoring effect, he gives an easily accomplished example. These were his instructions:

> "…Take a sheet of paper and draw a 2 1/2 – inch line going up, starting at the bottom of the page—without a ruler. Now take another sheet, and start at the top and draw a line going down until it is 2 1/2 inches from the bottom. Compare the lines." (Kahneman, 2011)

He noted it was likely that the first line drawn would be shorter than the second. Because we do not know exactly how long the first line is, there is a range of uncertainty as we approximate a past threshold while drawing the second line, the threshold being the length we think the first line is. We are likely to stop when we are no longer sure we should go further, and this usually results in stopping just short of the actual first line's end. Although there was imperfect adjustment, there was still deliberate adjustment from a reference point we acquired in the recent past. This adjustment was more effortful than what we would expect from System 1, as it was based on an attempt to find reasons to move away from the anchor, which requires engagement of System 2.

# 3  LIDA architecture and its cognitive cycle

The LIDA architecture is grounded in the LIDA cognitive cycle. Every autonomous agent (Stan Franklin & Graesser, 1997), be it human, animal, or artificial, must frequently sample (sense) its environment and select an appropriate response (action). Humans process (make sense of) the input from such sampling in order to facilitate their action selection. The agent's "life" can be viewed as consisting of a continual, overlapping sequence of these cognitive cycles. Each cycle consists of three phases, a perception/understanding phase, an attending phase, and an action selection/learning phase. It is commonly referred to as the action-perception cycle (Dijkstra, Schöner, & Gielen, 1994; Freeman, 2002; Neisser, 1976). A cognitive cycle can be thought of as a cognitive "moment." Higher-level cognitive processes are composed of many of these cognitive cycles, each a cognitive "atom."

Just as atoms have inner structure, the LIDA architecture includes a rich inner structure for its cognitive cycles (Bernard J. Baars & Franklin, 2003; S Franklin, B J Baars, U Ramamurthy, & Matthew Ventura, 2005). During each cognitive cycle a LIDA agent first makes sense of its current situation as best as it can by updating its representation of both external (coming through the external senses) and internally generated features of its world. This is the perception/understanding phase of the cycle. By a competitive process to be described below, it then decides what portion of the represented situation is most in need of attention, that is, most salient. This portion is broadcast to the rest of the system, making it the current contents of consciousness. This competition and the subsequent broadcast constitute the attending phase. The contents of a broadcast facilitate the recruitment of internal resources, that is, potential actions, from which the action selection mechanism chooses. This is the action selection phase.

Figure 1 shows this highly parallelized asynchronous processing in more detail. It starts in the upper left corner and proceeds roughly clockwise.
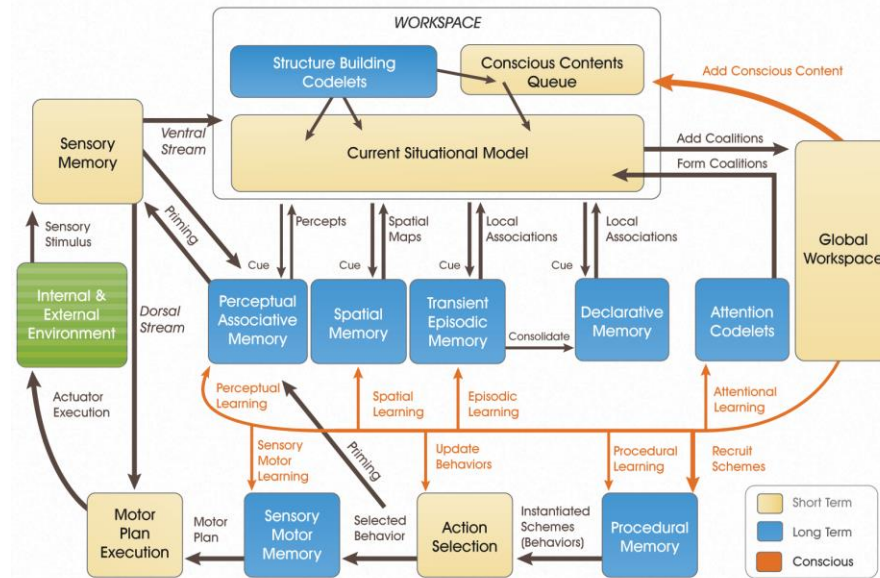


Figure 1 The LIDA cognitive cycle

The cycle begins with sensory stimuli from sources in the agent's external and internal environment being interpreted in Sensory Memory (e.g., the iconic buffer). Low-level feature detectors in Sensory Memory begin the process of making sense of the incoming stimuli. These low-level features are next processed by Perceptual Associative (recognition) Memory (PAM) where higher-level features, such as objects, actions, feelings, categories, relations, events, situations, etc., are recognized. These entities, when recognized preconsciously, ‡ comprise the current percept, which enters the Workspace asynchronously. Here a preconscious model of the agent's current situation, termed the Current Situational Model, is assembled (updated). This percept, and items from the Current Situational Model, serve to cue Spatial Memory and the two forms of episodic memory (the memory for events): transient (Conway, 2005) and declarative (autobiographical and semantic). Responses to the cues consist of a currently appropriate cognitive map, as well as local associations, which are remembered events from these two memory systems that were associated with the various elements of the cue. In addition to the current percept and the Current Situational Model, the Workspace contains the Conscious Contents Queue, a queue of the recent contents of consciousness (see below).

A new Current Situational Model (CSM) is assembled from the percepts, the cognitive map, the associations, and the undecayed parts of the previous model. This assembling process will typically require structure-building codelets. These are small, special-purpose processors, each of which has some particular type of structure it is designed to build. These codelets are continually monitoring the Workspace for opportunities to fulfill their particularly specified task. They may draw upon recognition memory and even Sensory Memory to enable the recognition of relations and situations. The newly assembled model constitutes the agent's best understanding of its current situation within its world. This completes the perception/understanding phase for this cognitive cycle.

For an agent operating within a complex, dynamically changing environment, the contents of the Current Situational Model (CSM) may be much too large for the agent to consider all at once (in humans

---

‡ In the LIDA architecture a preconscious representation is not conscious at that moment, but has the potential to become conscious (Franklin & Baars, 2010).

within ~100 ms (T Madl, Baars, & Franklin.Stan, 2011)). It needs to selectively attend to a salient portion of the CSM. Portions of the CSM compete, based on their activation, for attention. These portions take the form of coalitions of structures from the CSM. Such coalitions are formed by special-purpose attention codelets, whose function is to bring certain structures of concern to the particular attention codelet into the Global Workspace (hence the name Global Workspace Theory (Bernard J. Baars, 1988)). The most salient (important, urgent, insistent, novel, arousing, unexpected, bright, loud, moving, etc.) coalition wins the competition. In effect, the agent pre-consciously inferred what to attend to at the moment. The contents of the winning coalition are then broadcast globally, bringing its contents to consciousness and, thereby, completing the attending phase of the cycle. Thus the architecture's functional consciousness mechanism acts as an attention filter.

The purpose of all this attentional processing is twofold. First, to help the agent choose what to do next. Second, to select what to learn in updating each of its several memories. Learning and action selection take place concurrently in LIDA. We will describe action selection first.

Though the contents of this conscious broadcast are available globally, a primary recipient is Procedural Memory, which stores schemes of possible actions including their contexts and possible results. It also stores an activation value for each such scheme that attempts to measure the likelihood that an action taken within its context will produce the expected result. Schemes whose contexts sufficiently intersect the contents of the conscious broadcast instantiate specified copies of themselves. These instantiated copies are bound to the specifics of the current situation as specified in the broadcast conscious contents. As yet undecayed instantiated schemes remaining from previous cycles may also continue to be available. These instantiations are next processed by the action selection mechanism, which chooses a single action from one of these instantiations. The selected action then goes to Sensory-Motor Memory where it is executed by an appropriate algorithm (motor plan). The action taken affects the environment and thus completes that cognitive cycle.

Although we addressed the mechanisms behind how LIDA implements the perception-attention-action cycle, learning is also an important conceptual component of the model, and is especially important if we want to address evaluation and decision making in agents. Learning in LIDA is the encoding of knowledge about the past for use in the present. GWT supports the Conscious Learning Hypothesis: significant learning takes place via the interaction of consciousness with the various memory systems (B. J. Baars, 2003; S Franklin et al., 2005). That is, rely on conscious cognition for their updating, either in the course of a single cycle or over multiple cycles. LIDA's multiple modes of learning all occur continually, simultaneously, and online using each global broadcast of the contents of consciousness (S Franklin et al., 2005; Stan Franklin & Patterson, 2006). Perceptual learning is learning to recognize objects, feelings, categorizations, relationships, events, etc. As new objects, categories, and the relationships among them and between them and other elements of the agent's ontology are learned, nodes (objects and categories) and links (relationships) are added to Perceptual Associative Memory (Figure 1). Spatial learning refers to the building and updating of cognitive maps which serve to locate objects in the environment (Tamas Madl, Franklin, Chen, & Trappl). Episodic learning is the encoding of information into episodic memory, the associative, content-addressable, memory for events -- the what, the where, and the when (A. Baddeley, Conway, & Aggleton, 2001; S Franklin et al., 2005). Relatively little studied by memory theorists, attentional learning refers to the learning of what to pay attention to. In the LIDA architecture attentional learning is the learning of new attention codelets, and the updating and reinforcing of the existing ones (Faghihi, McCall, & Franklin 2012). Procedural learning is the encoding of procedures for executing behaviors into Procedural Memory (Figure 1). It is the learning of new actions and action sequences with which to accomplish new tasks (Stan Franklin et al., 2014). Here we must distinguish between action selection and action execution. LIDA's Procedural Memory is composed of schemes concerned with the selection of actions. Algorithms (motor plans) for the execution of actions are founfd in Sensory-Motor Memory where sensory-motor learning takes place.

In the LIDA architecture, all cognitive processing takes place via a continuing iteration of such cognitive cycles. A general design principle of the architecture demands that these cycles, and practically all of their processes, occur asynchronously. These cycles also cascade; that is, several cycles overlap, having different processes running simultaneously. This cascading must, however, respect the serial nature of conscious processing that is necessary to maintain the stable, coherent image of the world (Stan Franklin et al., 2005).

# 4 Action Selection/Decision Making in LIDA

Higher-level cognitive processing in humans might include imagination, deliberation, volitional decision-making, metacognition, reasoning, planning, scheduling, problem solving, etc. In the LIDA Model such higher-level processes require multiple cognitive cycles. Every higher-level cognitive process can be implemented by one or more behavior streams[§], that is, streams of instantiated schemes and links from procedural memory.

Cognitive processes have differing levels of control. Sloman (1999) distinguishes three levels that can be implemented by the architecture of an autonomous agent – the reactive, the deliberative, and the metacognitive[**]. The reactive, is the level we would typically expect of many insects, that is, a relatively direct connection between incoming sensory data and the outgoing actions of effectors. The key point is the relatively direct triggering of an action once the appropriate environmental situation occurs. Though direct, such a connection can be almost arbitrarily intricate, requiring quite complex algorithms to implement in an artificial agent. In the LIDA model, the action selected at the end of any single cognitive cycle is done so at the reactive level, though other reactive selection of actions may require several cognitive cycles.

The reactive level is perhaps best defined by what it is not. "What a purely reactive system cannot do is explicitly construct representations of alternative possible actions, evaluate them and choose between them, all in advance of performing them." (Sloman, 1999) Reactive control alone is particularly suitable for agents occupying relatively simple niches in reasonably stable environments, that is, for agents requiring relatively little flexibility in their action selection. Such purely reactive agents typically require relatively few higher-level, multi-cyclic cognitive processes.

Consciously planning a driving route from a current location to the airport is an example of deliberative, volitional decision making. Choosing to turn left at an appropriate intersection along the route requires consciously acquired information about the identity of the cross street, but the choice itself is most likely made unconsciously—the choice was consciously mediated even though it was unconsciously made—it did not require deliberation. Though heavily influenced by the conscious broadcast (i.e., the contents of consciousness), action selection during a single cognitive cycle in the LIDA model is not done consciously. A cognitive cycle is a mostly unconscious process. When speaking, for example, a person usually does not consciously think in advance about the structure and content of the next sentence, and is sometimes even surprised at what comes out. When approaching the intersection in the example above, no conscious thought need be given to the choice to turn left. Consciousness serves to provide information on which such action selection is based, but the selection itself is done unconsciously after the conscious broadcast. We refer to this very typical single cycle process as *consciously mediated action selection*. Each such is reactive in the sense of Sloman above.

---

[§]A stream is a sequence with its order only partially specified. Some actions in a stream may be taken in either order.

[**]Sloman speaks of meta-management rather than metacognition. We prefer the more common psychological term; they are synonyms in this context.

On the other hand, deliberative control typically employs such higher-level cognitive processes as planning, scheduling and problem solving. Such deliberative processes in humans, and in some other animals[††] (Mulcahy & Call, 2006; Tarsitano, 2006; Werdenich & Huber, 2006; Wilcox & Jackson, 2002), are typically performed in an internally constructed virtual reality. Such deliberative information processing and decision-making allows an agent to function more flexibly within a complicated niche in a complex, dynamic environment. An internal virtual reality for deliberation requires a short-term memory in which temporary structures can be constructed with which to mentally "try out" possible actions without actually executing them. In the LIDA Model the Workspace (4) serves just such a function. It is essentially a preconscious working memory in the sense of Baddeley (1992). The action selected during several cognitive cycles may consist of building, or adding to, some representational structures in the Workspace during the process of some sort of deliberation. Structure-building codelets, the sub-processes that create such structures, modify or compare them, etc., are typically implemented as internal reactive processes. Deliberation builds on reaction. In the LIDA Model, deliberation may be implemented as a collection of behavior streams, each selected action of which is an internal reactive process (S. Franklin, 2000).

Sloman's levels of control can be seen as being built on top of one another. However, Kahneman frames the issue of "control" as communication between two *separate* systems, as have dual-process theories of reasoning in general (see Osman, 2004, for a review). We offer a different conception of the two systems in LIDA, more analogous to Sloman, while still maintaining the essence of what these theories aim to predict.

## 4.1   Feelings and Emotions in LIDA

Every autonomous agent must be equipped with primitive motivators, drives that motivate its selection of actions (Stan Franklin & Graesser, 1997). In humans, in animals, and in the LIDA model, these drives are implemented by feelings and emotions, where emotions are taken to be feelings with cognitive content (S. Franklin & Ramamurthy, 2006). Such feelings implicitly give rise to values that serve to motivate action selection. This section is devoted to an explication of how feelings are represented in the LIDA model, the role they play in attention, and how they act as motivators, implicitly implementing values. Reference to Figure 1 may prove helpful to the reader.

Feelings are represented in the LIDA Model as nodes in its Perceptual Associative Memory. Each node constitutes its own identity, for example, distress at not enough oxygen is represented by one node, relief at taking a breath by another. Each feeling node has its own valence, always positive or always negative. The current activation of the node measures the arousal. Those feeling nodes with sufficient activations, along with their incoming links and object nodes, become part of the current percept and, are passed to the Workspace.

Like other Workspace structures, feeling nodes help to cue the two episodic memories. The resulting local associations may also contain feeling nodes associated with memories of past events. These feeling nodes play a major role in the assigning of activation to coalitions of information to which they belong, helping them to compete for attention within the Global Workspace. Any feeling nodes that belong to the winning coalition become part of the conscious broadcast, the contents of consciousness. Thus the LIDA agent becomes conscious of those feelings.

Any feeling node in the conscious broadcast that also occurs in the context of a scheme in Procedural Memory adds to the current activation of that scheme, increasing the likelihood of it instantiating a copy of itself into the Action Selection mechanism. It is here that feelings play their first role as implementation of motivation by adding to the likelihood of a particular action being selected. That

---

[††]Deliberation has been demonstrated in apes, birds, and even in arachnids (Wilcox and Jackson, 2002, Tarsitano, 2006).

feeling in the context of the scheme implicitly increases the value of the result of taking that scheme's action.

In the Action Selection mechanism, the activation of a particular scheme, and thus its ability to compete for selection and execution, depends upon several factors. These factors include how well the context specified by the scheme agrees with the current and very recently past contents of consciousness, that is, with the current situation. As mentioned earlier, the activation of this newly arriving behavior also depends on the presence of feeling nodes in its context, and their activation as part of the conscious broadcasts. Thus feelings contribute motivation for taking action to the activation of newly arriving behavior schemes.

# 5  Kahneman's Systems 1 and 2 à la LIDA

Consciously mediated action selection in LIDA likely happens in one or two cognitive cycles, whereas deliberative action selection typically requires a number of cognitive cycles. We maintain that LIDA's consciously mediated action selection (CMAS) provides the underlying mechanism for decisions taken by Kahneman's System 1. System 1 decisions are quick. The LIDA model postulates that CMAS occurs in humans in a fraction of a second (T Madl et al., 2011). System 1's intuitive decisions are based on emotions, as are CMAS as was described in the previous section. They are also based on vivid images. The LIDA model proposes that such images are the result of cued local associations in the Workspace, or possibly the output of structure building codelets, that such images can come to consciousness via the attention mechanism, and that they thereby play a role in CMAS. Finally, System 1 decisions are based on associative memory, which in the LIDA model includes both Perceptual Associative Memory (PAM) and both of the Episodic Memories. Local associations from these memory systems are cued into the Workspace, can be attended to so as to become part of the conscious contents, and so can influence CMAS.

Kahneman's System 2 is the slow, effortful system, which observes System 1's output, and intervenes when the output is considered to either infringe on more rational decision-making rules (System 2), or when an agent's intuition is insufficient in handling a situation. We suggest that System 2 is realized in humans by multi-cyclic, deliberative decision making (DDM) as modeled by LIDA. Being multi-cyclic, DDM is certainly slow relative to CMAS. The feeling of being effortful is part of fringe consciousness (Mangan, 2001). As such, it is represented by a feeling node in PAM that can come to the Workspace as part of a percept, and so possibly to consciousness, from which it can participate in deliberation. The actions taken as a result of CMAS are noted in the Workspace, and become grist for the deliberative mill. Thus we can say that System 2 (DDM) observes System 1's (CMAS's) output. In the LIDA model, rational decision-making rules (System 2) are denizens of Semantic Memory, which is a part of LIDA's Declarative Memory. Thus such rules can be cued from Semantic Memory as local associations into LIDA's Workspace, and thus play a role in System 2's (DDM's) decision making. Decision making which occurs in LIDA's Procedural Memory implements James' ideomotor theory of volition (1890). That is, among others, once a proposal, an objection, a support gets into the PM, some schemes from the Scheme net get instantiated. For instance, if PM receives a proposal, one of the schemes (time keeper) whose action is internal, starts a timer for writing the agent's current goal in the current situational model. The scheme whose action is starting the timer competes among other schemes in the behavior net and if it wins, the Ideomotor theory process starts. While the timer is running, if an objection, or a feeling of not knowing what to do, comes to consciousness a scheme whose internal action is to turn off the timer for the agent's current goal gets instantiated from the Scheme net.  If it's selected, its internal action will turn off the timer.

That is, the agent's intuition is insufficient to handle a situation when a feeling of not knowing what to do comes to consciousness. Such a feeling is part of fringe consciousness (Mangan, 2001).

However, a later supporter would likely result in the instantiation of a behavior to turn the timer for the agent's current goal back on, this time with a shorter time-lapse. Putting less time on the timer for the agent' goal makes it more likely to be chosen as winner.

Though Kahneman views System 1 and System 2 as separate mental processes, taking their underlying mechanisms as specified by the LIDA model to be CMAS and DDM respectively, implies that the System 2 actually is implemented by multiple instances of System 1. This follows since every deliberation in LIDA is carried out by a behavior stream producing a sequence of consciously mediated action selections.

***Attention and ease a la LIDA:*** Attention is the process of bringing content to consciousness (B. J. Baars, 1988, p. 299), a definition that is adopted by the LIDA model (Faghihi et al., 2012). In the LIDA model, passive or bottom-up attention is a hard-wired reflex. Active or top-down attention is achieved by a selected action, which has the agent focus on content in the Workspace (e.g., a person). Such an action creates an attention codelet that focuses on the person. In LIDA, an expectation codelet is a kind of attention codelet that tries to bring the result (or non-result) of a previously executed action to consciousness. An expectation codelet is created when a behavior is selected.

Attention is comprised of three concepts: 1) Alerting: "maintaining an alert state" which can be occurred deliberately and non-deliberately in LIDA; 2) Orienting: "focusing our senses on the information we want" (e.g., your focus on reading this document); 3) Executive attention: "the ability to manage attention towards goals and planning." (B. J. Baars, 1988, p. 229) In the LIDA model, attention codelets whose contents are sensitive to movement, location, or spatial information help bring content to consciousness—Kahneman's System 1. If such content wins the competition for consciousness an orienting action may ensue- Kahneman's System 1. Executive attention in LIDA, Kahneman's System 2, is illustrated using Miller and Cohen's example: consider a US citizen in the UK wanting to cross a road: In LIDA, the road stimulus is recognized by feature detectors in PAM which instantiates a "road" node. The context (UK) was also recognized by PAM, and would be part of the content in the Current Situational Model in the Workspace. Supposing that both the road node and the UK node entered the Global Workspace, and won the competition for consciousness, these nodes would be broadcast throughout the system including to Procedural Memory. Appropriate schemes in Procedural Memory could be instantiated, e.g. a scheme with context "road" and "UK", and action of "look- right".

In the LIDA model, the expectation codelets play an important role in resolving conflicts. An expectation codelet's content of concern comes from the result of the selected behavior that created it. This helps the system detect errors, and consequently search for solutions to fix the errors. A LIDA agent selects an action look left with the expectations of the car coming from the left. However, in UK, no car comes from left. Then the scheme that contains the action (look left before crossing the road) has its base-level activation decreased (Faghihi et al., 2012).

***Intuition & expert:*** In a LIDA agent we hypothesize that most decision-making, even by an expert, occurs in a consciously mediated manner. In the LIDA model, expert knowledge may be learned from frequent temporal patterns among conscious events. Consider Kahneman's chess expert example. In a LIDA agent every move's result is stored in Procedural Memory as schemes of possible actions including their contexts and possible results. Each scheme also stores two activation values, a base-level activation and a current-level activation. The base-level activation (used for learning) is a measure of the scheme's overall reliability in the past. It estimates the likelihood of the result of the scheme occurring by taking the action given its context. The current activation is a measure of the relevance of the scheme to the current situation (i.e., conditions, goals, etc,). Schemes whose contexts sufficiently intersect the contents of the current conscious broadcast instantiate specified copies of themselves. Given a specific situation in chess, the higher the base-level activation value of a scheme, the more this

specific scheme was chosen before. Thus, such a specific scheme would likely to be fired, once a similar situation occurs. That is, for that specific situation, the expectation codelets brought positive information to consciousness lots of times. Thus, the schemes' base-level activation gets saturated, which means it is likely to be fired in similar situations. However, for this specific situation, if the scheme receives negative reinforcement it would less likely to be fired in the future. That is, the expectation codelets brought negative reinforcement to consciousness. Thus, the schemes' base-level activation may not be enough to get the scheme fired on similar situations.

*Associative Coherence:* Let's assume a LIDA agent replicating the experiment mentioned above in connection with associative coherence. That agent would have a learned node in PAM for the word Vomit. This node would be instantiated while PAM's nodes detect the word Banana, and would likely come to consciousness. Later, the Bananas node would likely be similarly recognized with, due to the short time passage, Vomit still active in the Conscious Contents Queue of the Workspace. A temporal structure-building codelet then builds a new structure in the Current Situational Model of the Workspace, based on the Bananas node and its temporal predecessor, Vomit. Then, suppose this structure of Vomit, link, and Bananas is included in a coalition and wins the competition for consciousness. The subsequent broadcast instantiates a previously learned scheme with context "Bananas" and action "avoid eating." The ensuing behavior is selected for execution. Eating Bananas causes the agent to sense itself avoiding eating, another event, which includes a negative feeling, "Bananas displeasure." If the avoiding eating event comes to consciousness, two kinds of learning are performed: 1) a temporal link from "Bananas" to "eat" representing an affordance is added to PAM, and 2) due to the negative affective valence from the feeling node, a negative update is made to the "eat banana" node's base-level incentive salience‡‡ (Koedinger, Anderson, Hadley, & Mark, 1997). Given the negative emotional valance attached to the eat banana events, the event would likely make it to consciousness and get broadcasted. Since there is a temporal link between two events, temporal difference learning occurs and updates the earlier event's (Bananas) base-level incentive salience based on the difference between its current value and the base-level incentive salience of the latter event (avoiding eating). Consequently, the antecedent event "avoiding eating" loses base-level incentive salience, at first just for the Bananas node (immediate predecessor to avoiding eating), but later for the Bananas node as well. The result of all of this learning is that the node for the Bananas has a strong negative base-level activation, which helps it to strongly activate the "Bananas" node via the temporal links.

*Affect Heuristic:* Let's imagine a LIDA agent participating in the Hsee (1998) study mentioned earlier. To recall, participants on average rated the quality of a 24 piece dinner-set as better than a set with the same pieces plus extras, some of those extras being broken. This was meant to be an example of an intuitive evaluation instead of a rational one, because the existence of a few flaws while ignoring the total number of objects was enough to sway the participants to avoid the set with the flawed pieces, and instead choose the set with fewer pieces overall. He called this the "less is better" effect, a term which reflects intuitive decision making based on features easier to evaluate despite being less relevant.

The LIDA agent in this study would have representations like those of the human counterparts. In PAM, there is an association between a negative affective value, the form of the object, and its functional purpose (or perhaps the lack thereof). For this agent, there would be an association between a negative affective value, the broken nature of the dish that was sensed, and either the affordance of "unable to eat properly given this form," or no affordance at all. In either case, the salience of such a representation

---

‡‡ Incentive salience is a motivational "wanting" attribute given by the brain to stimuli transforming the representation of a stimulus into an object of attraction. This "wanting" is unlike "liking" in that liking is produced by a pleasure immediately gained from consumption or other contact with stimuli, while the "wanting" of incentive salience is a motivational magnet quality of a stimulus that makes it a desirable and attractive goal, transforming it from a mere sensory experience into something that commands attention, induces approach, and increases the likelihood of its being sought out.

is evaluated in the context of the agent's CSM, which calls for purchasing dinnerware with the proper affordances. In PAM, there would also be similar representations for intact pieces of dinnerware, this time with positive affective value and eating affordances as what's being associated. Because this is a set of both good and bad dinnerware the agent is sensing at the moment, both representations co-activate and become a part of the agent's CSM.

In a rational agent, the number of good plates relative to the total would be most salient in its CSM of the world at that moment. In this more intuitive agent, however, the total number of utensils is less salient than the affective weights associated with certain ones. This is because the affective weight is not associated with the total number of items, as we would expect from a rational agent. Such weights are instead associated with the individual items themselves. Therefore, the percept which makes it to consciousness and provides context for scheme-based decision making are the broken pieces of dinnerware, and the fact they are included in a purchasable set. The possible schemes triggered by such context would be "seek a better set," "avoid this set," etc.

*Availability Heuristic:* The availability heuristic is used when an agent judges the plausibility/potency/importance of a category from a retrieved memory based on either the ease/fluency of its retrieval (System 1 usage) or the frequency of certain content that is retrieved (System 2 usage).

In LIDA, what comes from the percept and items from the Current Situational Model serve to cue the two forms of Episodic Memory (the memory for events): transient and declarative (autobiographical and semantic). Responses to the cue consist of local associations, which are remembered events from these two memory systems that were associated with the various elements of the cue. In addition to the current percept included in the Current Situational Model, the Workspace contains the Conscious Contents Queue, a queue of the recent contents of consciousness. Attention codelets observe the Current Situational Model in the Workspace, trying to bring its most energetic, which is the most salient content, to the Global Workspace.

For System 1, Kahneman gives an example in which a group of psychologists in the early 1990's, led by Norbert Schwartz, investigated how fluency of retrieval related to self-assessments of assertiveness. What they found was that the more instances of assertiveness they were asked to list, the less assertive they viewed themselves to be (Schwarz et al., 1991). This is because as the number of instances to be listed increased, the fluency of producing examples decreased. Because the fluency was less than they expected, they associated the negativity of that with their own self-assessment. The participants were thus victims of the availability heuristic a la System 1.

For a LIDA agent that is supposed to simulate the above experience, there would be negative affective valences attached to assertiveness nodes in Semantic Memory. After broadcasting the coalition, given the goal is to find as many assertiveness behaviors, some schemes of possible actions including their contexts and possible results will be instantiated from PM to compete for action selection. These instantiations are next processed by the action selection mechanism, which chooses a single action from one of these instantiations. The result of the action selection will be broadcasted. Given the LIDA agents' task is to find schemes with instances of assertiveness, the numbers of Schemes that are instantiated in PM will decrease causing the instantiation of expected codelets to bring negative results to consciousness.

Furthermore, in the LIDA agent's Semantic Memory nodes, there would be a negative affective valences attached to the association between the number of assertiveness nodes retrieved from semantic memory, and the affordance of the highly activated coalition in the CSM for listing the instances of assertiveness. This is in part due to the expectation codelets bringing negative feedback when it came to the agent's self-assessment. That is, after a while, the number of the instantiated schemes that compete

in Action Selection decreases. Thus, the salience of such a representation is evaluated in the context of the agent's CSM, which calls for the proper affordances.

For System 2, Kahneman gives an example of the same researchers recruiting two groups of students whose task was to recall instances of their routines influencing their cardiac health. The two groups were those who had a family history of heart disease and those who did not. The latter group produced the same effects as from the assertiveness study. The former group, however, had the opposite effect. Because it was about them, and due to their family history, they were put into a higher state of vigilance than would be expected from the other group. As they recalled more instances of safe behavior, they felt safer. When they recalled more instances of dangerous behavior, they felt more at risk. In their vigilance, they did not "go with the flow" as what would happen if they were in the other group. Instead, they engaged in deliberation to evaluate patterns in the content, and assessed long term consequences based on such. This is the difference between the two systems and how they implement the availability heuristic.

For system 2, consider a LIDA agent that was tasked with recalling instances of her routines influencing her cardiac health. In this case, the agents' semantic memory has structures containing negative feeling nodes about family cardiac history. Being retrieved to the CSM, and having negative feeling nodes, the coalition containing family history information is more likely to be selected for the conscious competition. This could be due to the negative feelings nodes that are attached to that coalition. So, the more behaviors with positive feeling attached are retrieved to the Global Workspace the less likely coalitions with negative feeling attached to them get selected. Thus, the agent felt more secure. A more detailed description of how the LIDA model might account for the availability heuristic has appeared earlier (Stan Franklin et al., 2005).

# 6  Experiments/simulation

In this section, using a LIDA-based agent, we replicated a psychological experiment of the reinforcer devaluation paradigm through which we explain how it simulates some of Kahneman's System 1 and System 2 concepts. Our experiments illustrate the difference between the fast acting, slow adapting model-free control (consciously mediated action selection in LIDA, and Kahneman's System 1) and the slow acting, fast adapting model-based control (volitional decision making in LIDA, and Kahneman's System 2). In particular, we adapted an experiment testing the effects of orbitofrontal lesions on the representation of incentive value in associative learning (Gallagher, McMahan, & Schoenbaum, 1999). The orbitofrontal cortex (OFC) is thought to contain representations of the motivational significance of cues (conditioned stimuli) and the incentive value of expected outcomes. The significance of the reinforcer devaluation task is that normal performance depends on the ability of a conditioned stimulus (CS) to gain access to the motivational properties of an upcoming unconditioned stimulus (US).

In the original study, the experimenters divided rats into two groups: those in the first had their OFC lesioned, while those in the second maintained an intact OFC. All rats were first trained in a *conditioning phase* in accordance with standard Pavlovian conditioning: Over a series of 40 trials, rats were presented with a 10-second light CS, which was paired with (immediately followed by) a food delivery, itself followed by a ten-minute period in which the rat was allowed to eat freely. After a series of conditioning trials, a conditioned response (food cup behavior) to the CS was established. The measure during these trials was the rat's appetitive behavior towards the food cup during the last 5 seconds of the 10-second cue.

After the conditioning phase, each rat underwent three trials in a different *US devaluation phase*. Here, in each trial, there was no light cue, rather food was delivered first, and was followed by an

aversive event, the injection of LiCl, producing temporary sickness. The US devaluation phase introduces two more experimental conditions: In the *paired injection condition* the experimenters injected the rats immediately after the eating period. In the *unpaired injection condition* the rats were injected six hours after eating. Combining these conditions with the earlier OFC lesion manipulation, there were a total of four experimental groups: lesioned-paired, lesioned-unpaired, intact-paired, and intact-unpaired. The measure during the US devaluation phase was the amount of food consumed during the eating stage of each trial.

After the US devaluation phase, the experimenters performed a *devaluation test phase* that revisited the rats' conditioned responses (CRs) to the light CS. In this phase, each rat was presented with the light CS only, i.e., with no further experimental manipulations. As in the first phase, the measure for these trials was the rat's appetitive behavior towards the food cup during the last 5 seconds of the 10-second cue.

Although the light CS was absent during the devaluation phase, its previous association with the food US provides a basis for anticipating the US. The experimenters found that lesions of OFC did not affect either 1) the initial acquisition of a conditioned response to the light CS in the initial conditioning phase or 2) the learning of food aversion in the US devaluation phase. However, in the devaluation test phase, OFC lesioned rats exhibited no change in their conditioned responding to the light CS, i.e., they continued to exhibit appetitive food cup behavior. This outcome contrasts with the behavior of control rats: after the devaluation of the US, a significant decrease in the food cup approaches (appetitive behavior) occurred in the devaluation test phase. The experimenters hypothesized that, after OFC damage, the cue was unable to access the representational information about the incentive value of the associated US (Gallagher et al., 1999).

**A LIDA-based agent Account of Experimental Behavior.** Recall that, in the first phase of the experiment, a light cue is paired with food delivery, that is, food is delivered immediately after the light signal terminates. The agent's appetitive behavior towards the food cup is recorded during the last 5 seconds of the 10-second cue. The results of phase 1 of the original experiment are shown in the first graph in Figure 2. The measure of learning in phase 1 was food cup behavior recorded as a percentage of total behavior recorded during the last 5-second observation interval of the 10-second CS presentation (the light). This was achieved by recording a single behavior for each 1.25-second interval, and then computing the percentage of behavior that was food cup behavior, i.e., the frequency of food cup behavior in an observation interval was divided by the total number of observations made in that interval. Now we describe how a LIDA agent models the events of this phase, and how it learns a conditioned response.

Let us first assume that a LIDA agent replicating this experiment would early on learn a memory for the light cue in the form of a "light-cue" node in its Perceptual Associative Memory (PAM). Then, during the conditioning phase, the light-cue node would be instantiated while the light is on, and this node would typically come to consciousness. Later, a "food" node is similarly recognized with, due to the short time passage, the light node still being active in the Conscious Contents Queue of the Workspace. A temporal structure-building codelet then builds a new structure in the Current Situational Model of the Workspace, based on the food node and its temporal predecessor, the light node. If this structure of light, link, and food is formed into a coalition by an attention codelet, and wins the competition for consciousness, then a new temporal link from light-cue to food would be learned in PAM.

The conscious broadcast of this structure would also serve to recruit resources to deal with the situation; in this case, we assume the broadcast instantiates one or more previously learned schemes from Procedural Memory having "food" in their context and some appetitive food cup action, e.g. approach cup or eat food, and that a resulting behavior is selected for execution. When the *eating* event occurs, and is recognized, it would be accompanied by the feeling node, "food pleasure," having positive affective valence. If the eating event comes to consciousness, two kinds of learning are performed: 1) a

temporal link from "approach" to "eat" is positively reinforced in PAM, and 2) due to the positive affective valence from the "food pleasure" feeling node, a positive update is made to the "eat" node's base-level incentive salience.

Several repetitions of this first phase would lead to repeated conscious broadcasts of the light-food-approach-eat event sequence (System 1). Temporal difference (TD) learning would occur each time a structure with a temporal link is present in the broadcast, and would update the base-level incentive salience of the link's source event. Working backwards in order of occurrence, TD learning would first update the base-level incentive salience of the "approach" event based on the difference between the approach event's current base-level incentive salience and that of the following event, eat. Later, the two antecedent events of food delivery and light-cue would, incrementally over multiple cognitive cycles, gain base-level incentive salience. At first, this would only affect the approach node (immediate predecessor to eat), but later the other nodes would receive some "credit" in predicting the food "reward." The upshot of all of this learning is that the node for the light cue gains a high base-level incentive salience, which later helps it to strongly activate the "food" node via learned temporal links. Once "food" is strongly activated, the selection of an appetitive food cup behavior is likely to occur, even in the absence of actual food.

For the second phase of the experiment, there were two conditions: 1) In the paired injection condition, rats were given food immediately followed by (paired with) an illness inducing LiCl injection, 2) the unpaired injection condition first provided food, but the LiCl injection occurred six hours later. The results from the original experiment are shown in the second graph in Figure 2. The experimental groups that received food paired with the injection are shown in white. These rats learned to greatly reduce their food consumption (System 2). The unpaired groups are shown in black. These groups attenuated their food consumption by significantly less. It was not shown that these unpaired groups significantly decreased their food consumption across sessions. One explanation for the apparent decrease for the rats in the unpaired groups is that they might have performed some deliberative temporal learning to actually associate the food with the injection (System 2 in action). For the paired injection experimental group, the mental events occurring in a LIDA-based agent are similar to those of the first experimental phase. The agent would, via conscious learning, add a temporal link from the food node to the injection event. Additionally, after the injection, the agent would recognize a "sickness" event and, via conscious learning, add another temporal link from the injection event to the sickness event. The sickness event would come with an accompanying "sickness" feeling node with negative affective valence. Conscious learning would then lead to the assignment of a negative base-level incentive salience to the sickness event. Repeated conscious exposures of this food-approach-eat-injection-sickness sequence would then, via temporal difference learning, "devalue" or decrease the base-level incentive salience, first of the injection event, then of the earlier events as well. For an agent in the unpaired group, the processing and learning is the same as for the paired group, except that since the injection occurs six hours after the food presentation, the node for food in the Workspace has long since decayed away during the aversive events. Thus, for such a simple agent, no temporal links are ever learned between the eating event and the injection or sickness. (In this second group, the injection and sickness events would still be learned with a temporal link between them and both would be given low base-level incentive salience.) For the unpaired groups, the apparent decrease in food consumption across sessions may have been a result of deliberative association between the injection event and an earlier event (e.g. food consumption) recalled from Episodic Memory (System2).
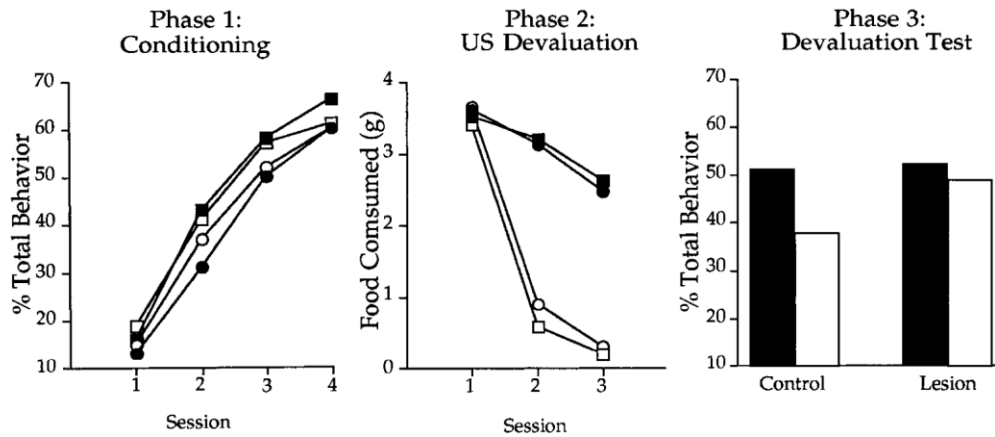
Figure 2. The results of the original reinforcer devaluation experiment with rats. *White* represents the paired injection groups and *black* the unpaired injection groups. *Squares* represent control groups, and *circles* represent the lesioned groups. The Phase 1 graph shows that all groups acquired conditioned responses to the light cue, as evidenced by their increased food cup behavior as a percentage of total behavior during the latter half of the light cue presentation. The Phase 2 graph shows that the rats receiving paired LiCl injections (white) significantly reduced their food consumption as compared to those rats receiving unpaired injections (black). There was neither a significant difference in food consumption (due to lesion) among the paired groups nor among the unpaired groups. Finally, the Phase 3 graph shows that only intact rats receiving paired injections (left white bar) significantly reduced their food cup behavior during the devaluation test.

In the final phase of the experiment the agent receives several presentations of the light cue, which are not followed by any additional experimental manipulations. As in phase 1, the agent's appetitive behavior towards the food cup is recorded during the last 5 seconds of the 10-second cue. The experimental findings for phase 3 of the experiment are shown in the final graph in Figure 2. Only the paired control (left white bar) group significantly decreased its rate of conditioned responses from the other groups, whose CR rate were statistically equivalent. The white bars represent groups receiving paired injections, and the black bars represent unpaired groups.

To describe a LIDA-based explanation of phase 3, we first note the LIDA agent can form new long-term memories based on temporal links, which affords the agent the ability to, based on the currently active nodes in the CSM of the Workspace, instantiate expected future event(s) into the CSM, each further ahead in time than the last. Observing CSM, structure building (SB) codelets try to find local associations from declarative, episodic or perceptual memory, and create future anticipation. Then, if the future anticipation built by the SB has enough activation, an attention codelet will make a coalition of that, and will send it for conscious competition. Each new future anticipation can be a proposal or objection that can be accepted or rejected by the LIDA agent's decision making system (Kahneman's expert and intuition). It is also hypothesized that expected events, their temporal links, and the original event can all form into a single coalition. This coalition competes based on the total activation and total incentive salience of each of these events including the expected one(s). In our implementation, we compute coalition activation based on 1) the average total activation of coalition nodes and 2) the average of the absolute value of the total incentive salience of coalition nodes. These two averages are combined and multiplied by the attention codelet's base-level activation, which correspond to the bottom up attentional learning (for more explanation the reader referred to the Vomit and Bananas example above).

Keeping this in mind, if we are to construct a LIDA agent replicating the results of the experiment, we must hypothesize a functional role for the OFC and relate this role to a capacity of a typical LIDA agent that must be removed to simulate an OFC lesion. The OFC has been suggested as critical for "associative learning," and the representation of "associative information, particularly information about the value of expected outcomes" (Schoenbaum, Takahashi, Liu, & McDannald, 2011). The neural activity in the OFC "increases to cues and after responses that predict rewards." Finally, the authors suggest viewing OFC function as "constructing or implementing a model-based representation" (Schoenbaum et al., 2011). Based on these ideas, we define a *lesioned OFC LIDA agent* as one that cannot use the total incentive salience of *expected* event(s) in determining the total incentive salience of an option temporally preceding those event(s). An *intact OFC LIDA agent* is one that can access the total incentive salience of expected event(s) and uses this in computing the total incentive salience of an option temporally preceding the event(s). We note that for both agent types, the base-level incentive salience of event nodes is assumed to be intact, the lesion does not affect existing memory in PAM or Procedural Memory, and, both agent types, lesioned or intact, can perform temporal difference learning.

The results for the unpaired injection groups in phase 3 can be explained simply: Since the injections were unpaired, the aversive injection event would not be active in LIDA's Workspace contemporaneously with the food event, and thus it cannot be associated with the food event by structure-building codelets in the Workspace. This fact is independent of whether the agent is lesioned or not. Thus a potential temporal link never comes to consciousness and no TD learning can occur which might devalue the base-level incentive salience of the temporally earlier events—food delivery and light-cue. As a result, the light event retains its high base-level incentive salience, originally learned from phase 1, motivating the agent to approach the food cup when the light cue is later shown in phase 3.

Now, let's consider the two paired injection groups. For the paired-lesioned OFC group, the agent is only able to evaluate a stimulus' (light cue) value based on its base-level incentive salience. This would prevent the agent from integrating any expectation of future aversive events into a coalition with an instantiated light event node. Since the light cue occurred in the initial conditioning phase of the experiment, its base-level incentive salience was positively updated by TD learning. However, since the light cue did not occur in phase 2, it could not have been altered by TD learning. Consequently, as an option, the light cue would have an overall positive incentive salience, and, via a conscious broadcast, might instantiate schemes leading to appetitive food cup behavior. The expected result is that a lesioned-paired LIDA agent would exhibit a similar percentage of food cup behavior as both unpaired groups.

Finally, why might an intact-paired injection agent reduce its food cup behavior? We suggest that this type of agent, given the instantiation of the light cue node, is able to instantiate the subsequent events it has learned to expect. It does this by repeatedly cueing with its PAM based instantiated expected events in the Workspace. The initial light-cue event cues PAM instantiating the food delivery event into the Workspace. Next, the eating event is instantiated, then injection, etc. An attention codelet can form a coalition from this integrated sequence of events and bring it to the Global Workspace. While the earlier events in this sequence may have a fairly high base-level incentive salience, the later ones would surely have negative base-level incentive salience due to the devaluation trials. Such an option would then have less overall incentive salience and, consequently, less of a chance to win the competition for consciousness. Even if it does win, it has less of a chance to induce an appetitive action selection.

# 7 Conclusion

According to Kahneman, human decision making process consists of two main processes that characterize thinking (Kahneman, 2011Kahneman, 2011): "System 1," the fast system, is responsible for intuitive decisions based on emotions, vivid imagery and associative memory. "System 2," the slow system, observes System 1's outputs, and intervenes when an agent believes its intuition is insufficient. The LIDA model postulates that human cognition consists of a continuing iteration of cognitive cycles, each a cognitive "atom," out of which higher-order processes are built. In LIDA System 1 is consciously mediated action selection that occurs in one or two cognitive cycles, while System 2 employs multiple cognitive cycles in its deliberative decision making. Furthermore, the LIDA model suggests that System 2 employs multiple occurrences of System 1 in its operation.

Throughout this paper we explained what LIDA's conceptual model suggests as the underlying mechanism for Kahneman's System 1 and System 2. We also replicated an experiment *in silico* that briefly explains what is discussed in this paper regarding Kahneman's book, using LIDA based software agents as simulated subjects.

# References

Baars, Bernard J. (1988). *A Cognitive Theory of Consciousness*. Cambridge: Cambridge University Press.

Baars, Bernard J. (2002). The conscious access hypothesis: origins and recent evidence. *Trends in Cognitive Science, 6*, 47–52.

Baars, B. J. (2003). The global brainweb: An update on global workspace theory. *Science and Consciousness Review*(2).

Baars, Bernard J., & Franklin, S. (2003). How conscious experience and working memory interact. *Trends in Cognitive Science, 7*, 166–172.

Baddeley, A., Conway, M., & Aggleton, J. (2001). *Episodic Memory*. Oxford: Oxford University Press.

Baddeley, A. D. (1992). Consciousness and working memory. *Consciousness and Cognition, 1*(1), 3-6.

Bijleveld, E., Custers, R., & Aarts, H. (2009). The Unconscious Eye Opener Pupil Dilation Reveals Strategic Recruitment of Resources Upon Presentation of Subliminal Reward Cues. *Psychological Science, 20*(11), 1313-1315.

Conway, M. A. (2005). Memory and the self. *Journal of Memory and Language*, 594-628.

D'Mello, S. K., & Franklin, S. P. (2004). A cognitive architecture capable of human like learning. *Connection Science*.

Dijkstra, T. M. H., Schöner, G., & Gielen, C. C. A. M. (1994). Temporal stability of the action-perception cycle for postural control in a moving visual environment. *Experimental Brain Research, 97*(3), 477-486.

Faghihi, U., McCall, R., & Franklin, S. (2012). A Computational Model of Attentional Learning in a Cognitive Agent. *Biologically Inspired Cognitive Architectures, 2*, 25-36.

Franklin, S. (2000). Deliberation and Voluntary Action in 'Conscious' Software Agents. *Neural Network World, 10*, 505-521.

Franklin, S., Baars, B. J., Ramamurthy, U., & Ventura, M. (2005). The Role of Consciousness in Memory. *Brains, Minds and Media, Vol.1, bmm150 (urn:nbn:de:0009-3-1505)*.

Franklin, S., & Graesser, A. (1997). Is it an Agent, or just a Program?: A Taxonomy for Autonomous Agents. *Proceedings of the Third International Workshop on Agent Theories, Architectures, and Languages, published as Intelligent Agents III, Springer-Verlag*, 21-35.

Franklin, S., Madl, T., D'Mello, S., & Snaider, J. (2014). LIDA: LIDA: A systems-level architecture for cognition, emotion, and learning. *IEEE Transactions on Autonomous Mental Development., PP(99), 1 doi: 10.1109/TAMD.2013.2277589*

Franklin, S., & Patterson, F. G. J. (2006). The LIDA Architecture: Adding New Modes of Learning to an Intelligent, Autonomous, Software Agent *IDPT-2006 Proceedings (Integrated Design and Process Technology)*: Society for Design and Process Science.

Franklin, S., & Ramamurthy, U. (2006). *Motivations, Values and Emotions: 3 sides of the same coin.* Paper presented at the Sixth International Workshop on Epigenetic Robotics, Paris, France.

Freeman, W. J. (2002). The limbic action-perception cycle controlling goal-directed animal behavior. *Neural Networks, 3*, 2249-2254.

Gallagher, M., McMahan, R. W., & Schoenbaum, G. (1999). Orbitofrontal cortex and representation of incentive value in associative learning. *The Journal of Neuroscience, 19*(15), 6610-6614.

Hsee, C. K., Weber, E. U., & Welch, N. (2001). Risk as Feelings (Vol. 127, pp. 267-286): American Psychological Association.

James, W. (1890). The Principles of Psychology. *Cambridge, MA: Harvard University Press*.

James, W. (1950). The Principles of Psychology Volume. 2. *Journal of Language and Social Psychology, 25*(4).

Kahneman, D. (1973). *Attention and effort*. Englewood Cliffs, N.J.: Prentice-Hall.

Kahneman, D. (2011). *Thinking, Fast and Slow*. New York: Farrar, Straus and Giroux.

Koedinger, K. R., Anderson, J. R., Hadley, W. H., & Mark, M. A. (1997). Intelligent tutoring goes to school in the big city. *International Journal of Artificial Intelligence in Education (IJAIED), 8*, 30-43.

Madl, T., Baars, B. J., & Franklin.Stan. (2011). The Timing of the Cognitive Cycle. *PLoS ONE, 6*(4).

Madl, T., Franklin, S., Chen, K., & Trappl, R. *Spatial Working Memory in the LIDA Cognitive Architecture.* Paper presented at the Proc. International Conference on Cognitive Modelling.

Mangan, B. (2001). Sensation's Ghost: The Non-Sensory "Fringe" of Consciousness. *Psyche, 7*, http://psyche.cs.monash.edu.au/v7/psyche−7−18−mangan.html.

Marshall, S. P. (2002). *The index of cognitive activity: Measuring cognitive workload.* Paper presented at the Human factors and power plants, 2002. proceedings of the 2002 IEEE 7th conference on.

Mulcahy, N. J., & Call, J. (2006). Apes save tools for future use. *Science, 312*(5776), 1038-1040.

Mussweiler, T., & Strack, F. (2000). The use of category and exemplar knowledge in the solution of anchoring tasks. *Journal of Personality and Social Psychology, 78*(6), 1038-1052.

Neisser, U. (1976). *Cognition and Reality: Principles and Implications of Cognitive Psychology* San Francisco: W. H. Freeman.

Osman, M. (2004). An evaluation of dual-process theories of reasoning. *Psychonomic Bulletin & Review, 11*(6), 988-1010.

Peavler, W. S. (1974). Pupil size, information overload, and performance differences. *Psychophysiology, 11*(5), 559-566.

Peters, E., & Levin, I. P. (2008). Dissecting the risky-choice framing effect: Numeracy as an individual-difference factor in weighting risky and riskless options. *Judgment and Decision Making, 3*(6), 435-448.

Pfeiffer, A. M., Whelan, J. P., & Martin, J. M. (2000). Decision-making bias in psychotherapy: Effects of hypothesis source and accountability. *Journal of Counseling Psychology, 47*, 429-436.

Schoenbaum, G., Takahashi, Y., Liu, T. L., & McDannald, M. A. (2011). Does the orbitofrontal cortex signal value? *Annals of the New York Academy of sciences, 1239*(1), 87-99.

Schwarz, N., Bless, H., Strack, F., Klumpp, G., Rittenauer-Schatka, H., & Simons, A. (1991). Ease of retrieval as information: Another look at the availability heuristic. *Journal of Personality and Social Psychology, 61*(2), 195-202.

Scott, S. G., & Bruce, R. A. (1995). Decision making style: The development and assessment of a new measure. *Educational and Psychological Measurement, 55*, 818-831.

Shiv, B., & Fedorikhin, A. (1999). Heart and Mind in Conflict: The Interplay of Affect and Cognition in Consumer Decision Making. *Journal of Consumer Research, 26*, 278–292.

Sloman, A. (1999). What Sort of Architecture is Required for a Human-like Agent? In M. Wooldridge & A. S. Rao (Eds.), *Foundations of Rational Agency* (pp. 35–52). Dordrecht, Netherlands: Kluwer Academic Publishers.

Slovic, P. E. (2000). *The perception of risk*: Earthscan Publications.

Tarsitano, M. (2006). Route selection by a jumping spider (Portia labiata) during the locomotory phase of a detour. *Animal Behaviour, 72*(6), 1437-1442.

Werdenich, D., & Huber, L. (2006). A case of quick problem solving in birds: string pulling in keas, Nestor notabilis. *Animal Behaviour, 71*(4), 855-863.

Wilcox, S., & Jackson, R. (2002). Jumping spider tricksters: deceit, predation, and cognition. *The cognitive animal*, 27-33.