

# Neuroscientific Implications for Situated and Embodied Artificial Intelligence

Keith L. Downing  
The Norwegian University of Science and Technology  
Trondheim, Norway  
keithd@idi.ntnu.no

January 1, 2007

## Abstract

While classic AI systems still struggle to properly incorporate commonsense knowledge, Situated and Embodied Artificial Intelligence (SEAI) aims to build animats that acquire a commonsense understanding of the world via interactions between simulated brains, bodies and environments. Neuroscientists believe that much of this common sense involves predictive models for physical activities, but the transfer of sensorimotor skill knowledge to cognition is non-trivial, indicating that SEAI may meet a daunting challenge of its own. This paper considers the neurological bases for implicit procedural and explicit declarative common sense, and the possibilities for its transfer from the former to the latter. This helps assess the prospects for SEAI to eventually surpass GOFAI in the quest for generally intelligent systems.

**Keywords:** Artificial Intelligence, Artificial Life, Neuroscience, Procedural and Declarative Knowledge, Skill Compilation

## 1 Introduction

The field of Artificial Intelligence (AI) began in a blaze of glory [43, 44, 48]. Within a decade or so of its inception, computers were solving geometry and physics problems at a college freshman level, playing chess like regional champions, diagnosing serious illnesses on par with expert physicians and configuring VAX machines. No problem was too complex, but, as AI researchers discovered in the 1980's, many were too *simple*.

Indeed, the capabilities that humans take for granted, our basic sensorimotor skills such as walking, climbing, and grasping for objects, turned out to be orders of magnitude more difficult to program than backgammon, bridge and biochemical analysis. By the mid 1980's, AI researchers realized that a serious shortcoming in their systems was none other than commonsense. AI systems behaved

like idiot savants, producing exceptional results on a wide range of situations, but floundering miserably on cases that demanded basic intuitions about the world, intuitions that most humans have acquired by their 2nd birthday.

Many attempts were made to force-feed this common sense into AI systems, in much the same manner and using similar knowledge-representation formats as had been successfully used to load expert rules-of-thumb into AI systems. In fact, the whole AI subfield of qualitative reasoning (QR) [17] was dedicated to this aim. Although QR produced many useful paradigms whose applications range from intelligent tutoring to plant monitoring to automobile and Mars-rover diagnosis, it did little to fortify AI systems with that broad base of general knowledge needed to transform idiot-savants into well respected (and trusted) gurus.

As AI was facing this disappointing truth in the 1980's, a related, but diametrically-opposed field began to take root: Artificial Life (ALife) [36]. Although ALife researchers primarily seek to understand the life process at a level far removed from that of neuroscience and psychology, the basic philosophy had immediate implications and inspiration for AI, although only a few AI researchers took note. The essential transferable concepts from ALife to AI are situatedness and embodiment. Although often trivial in their detail, the vast majority of ALife systems consist of simulated organisms that reside in environments (situatedness) and have a body (embodiment) whose survival depends upon a fruitful interaction with those surroundings.

Classic AI systems, often called GOFAI (Good Old-Fashioned AI) systems, assume away all environmental and bodily factors to focus on cognition in a vacuum. This works well for chess but fails consistently in robotics. As GOFAI researchers found out, general abstract-reasoning systems do not plug-and-play with any set of sensors and motors. General commonsense exists not in a platform-independent piece of software, but in a behavioral repertoire that is finely tuned to the structure and dynamics of both body and environment. And although some aspects of this repertoire are easily explained in everyday terms and rules-of-thumb, many are the unique province of biology and engineering. Logics, so common to GOFAI, are of little utility, but many of ALife's kernel concepts: emergence, competition, cooperation, etc., form the backbone of this low road to understanding intelligence [10].

For that handful of AI researchers [6, 53] who saw ALife as more than Turing-equivalent cellular-automata gliders and evolving biomorphs, but as a more fundamentally sound approach to cognition, the motivating thesis can be approximated as:

Complex intelligence is better understood and more successfully embodied in artifacts by working up from low-level sensory-motor agents than by working down from abstract cognitive mechanisms of rationality (e.g. logic, means-ends analysis, etc.).

Essentially, Situated and Embodied AI (SEAI) researchers believe that GOFAI's holy grail, common sense, comes only via the learned experiences of a body in a world. There are significant limits to how much knowledge one body (a teacher or an expert-system designer) can transfer to another (a student or an expert system), and with common sense, these limits are very stringent. Whereas "I think, therefore I am" might have been an appropriate slogan for GOFAI, it's converse more aptly summarizes SEAI. That is, by living, we acquire common sense, which then supports more

complex reasoning.

Andy Clark [8] uses the term *cognitive incrementalism* to denote this general bootstrapping of intelligence:

This is the idea that you do indeed get full-blown, human cognition by gradually adding bells and whistles to basic (embodied, embedded) strategies of relating to the present at hand.

The grand challenge to cognitive incrementalism may come from the work of Lakoff and Nunez [35], who explain mathematical reasoning, both simple and complex, as an extension of our sensorimotor understanding of the world. The neuroscientific grounding of their theory is weak, but the metaphorical ties between embedded and embodied action on the one hand and mathematical concepts on the other are striking. By linking everyday sensing and acting to one of man's most abstract cognitive endeavors, the authors implicitly motivate a Turing-type challenge for SEAI: build a sense-and-act robot that evolves the ability to do mathematics.

Ignoring the obvious difficulty of this challenge, the general SEAI philosophy and its bottom-up approach to knowledge acquisition hold some promise. After all, one can hardly deny the importance of first-hand experience in learning simple facts of life such as that wet things can be slippery, make you cold, etc. However, the cruel reality of engineering raises major obstacles, as all roboticists know.

Whereas GOFAI began with a divergent radiation of impressive applications that displayed many forms of (shallow) intelligence, SEAI seems to have converged on a menagerie of wall-following robots, all of which have very deep, functional (albeit implicit) understandings of their own body and domain: a barren floor surrounded by walls. To date, there are no biologically-inspired robots that display transferable common sense. That is, many systems behave intelligently and exhibit minimally cognitive behaviors, such as perceptual classification, attentional focusing, etc. [5], but the common sense is so tightly embedded in reactive routines (or controllers with simple notions of internal state) that it evades reuse for other tasks. So to date, sensing and acting has not produced common sense scaffolding for cognitive activities, such as the **planning** of motor sequences.

After 20 years of SEAI, one expects more. GOFAI adherents can arguably write-off SEAI as overly-optimistic biological envy in the same way that many of AI's *scruffies* criticize the *neat* logical approaches to intelligence as mere mathematical envy. Still, the fact remains that GOFAI was built on a computational foundation that was rock-solid for building useful engineering tools but flimsy and misleading as a model of intelligence. SEAI has yielded a few useful artifacts, such as floor-sweeping and lawn-mowing robots, and a host of insect-like intelligences, but the functional cornerstone (normally neural networks or other systems of distributed processors) maps much more readily to the basis of animal behavior than do GOFAI's theorem provers, frames and semantic networks. SEAI is clearly paying the price for building a proper biologically-rooted foundation, but as GOFAI's failures in cognitive science indicate: a field's initial progress is not a clear indicator of its ultimate success.

This article assesses the eventual success of SEAI by consulting neuroscience and asking, "What is

the sensorimotoric basis of commonsense, and can it be transferred to the cognitive domain?”

## 2 Evolutionary Origins of Situated and Embodied Cognition

In *i of the Vortex* [37], Rudolfo Llinas presents an enlightening account of the sensorimotor foundations of cognition, based on both developmental and evolutionary neurobiology. He argues that brains evolved to support motion - stationary organisms do not need them - and gives the convincing example of a sea squirt, which is mobile during its early life stages but becomes sessile later, whereupon it digests its own brain!

Llinas gives embryological and phylogenetic evidence for an *encephalization* process whereby motor activity patterns that were originally emergent from the direct electrical couplings between muscle cells have, through the course of evolution, become controlled by, first, spinal motor neurons, and later, neurons of higher brain regions. This higher-level control increases the potential complexity of the actions, since emergent oscillatory patterns - typically, spatial waves of muscle contractions - cannot approach the intricacy needed for walking a balance beam or playing the piano. But with a neural hierarchy, spatial activation waves at the higher levels can, via tangled top-down connections, differential propagation delays, etc., cause spatially diverse firing patterns at lower levels.

Llinas believes that muscle oscillations in primitive animals and in developing vertebrate embryos have become internalized/encephalized to the 40 Hz activity of the thalamus, thus serving as a binding signal for mental activity. In a strong sense, this 40Hz signal is the heartbeat of the brain, and it arose via an evolutionary process that gradually translated emergent muscle-activation patterns into a high-level dynamic that coordinates all activity: perceptual, motor and cognitive. In short, thought is encephalized motricity.

The nobel laureate Gerald Edelman [14] uses the concept of Neural Darwinism to support his theory of neural group selection (TNGS). In this view, neurons undergo a selective process wherein only those that a) grow connections to other neurons, and b) frequently use those connections, will survive. Essentially, neurons are involved in a *survival of the best networkers* competition. Terrence Deacon [9] uses TNGS as the basis for his Displacement Theory (DT), wherein the networking competition during early development leads to brains that essentially scale to fit the body’s sensory and motor apparatus. In short, primary sensory and motor areas of the brain are sized according to their immediate inputs or outputs, respectively. Secondary region sizes derive from those of the primary regions, and deep cortical structures grow or shrink to fit their input and output sources.

Developmental neuroscience clearly supports and employs the key tenets of DT. For example, Fuster [18] documents the earlier maturation of posterior brain regions (used in early sensory processing) and late maturation of frontal areas such as the prefrontal cortex (PFC). Striedter [54] combines DT with the work of Finlay and Darlington [16] to explain the trend of greater neocortical (and especially PFC) control of lower brain regions in higher organisms. First, Finlay and Darlington show that *late equals large* in neurogenesis: larger brain regions are those that mature later in development. Second, a key corollary of Deacon’s DT is that *large equals well-connected*: big brain regions send many axons to other regions and thereby have a significant amount of control over

them. Together, these show how small changes in developmental timing (in higher mammals) have enabled the frontal regions to mature later, hence grow larger, and hence exhibit greater control over a wide variety of cortical areas. And greater frontal control correlates well with behavioral sophistication, as illustrated by Nakajima et al.'s [42] comparisons of manual dexterity vis-a-vis frontal control of motor areas in mammals such as cats, monkeys and humans.

This links back to Llinas and even revitalizes the 19th century work of Ferrier [15] on functional encephalization or *neocorticalization*. In short, the control of many functions has *moved up* the cortical hierarchy precisely because the higher (evolutionarily younger and developmentally later) regions have grown larger and *sent down* more axons.

For situated and embodied AI, Edelman and Deacon's work has several interesting implications. First, since bodies evolve to match their environments, and brains grow to fit bodies, the physical structure of the brain is ultimately determined by the environment and body. Clearly, we cannot ignore these factors when analyzing cognitive activities, since the mental machinery itself stems directly from embodiment and situatedness. Second, Deacon goes to great lengths to show that the brains of primates and other higher organisms are largely redimensioned versions of the brains of lower organisms. For example, in humans, the big winners in Neural Darwinism and Displacement are the cerebellum and prefrontal cortex (PFC). These, Deacon believes, are the basis for our advanced speech and symbol-processing capabilities. Thus, cognition stems from a straightforward redimensioning of a primarily sensorimotor brain, not from the addition of extra cognitive modules. Considerable neurobiological evidence, summarized and integrated by Fuster [18], supports the general idea that the prefrontal cortex is the center of both motor and cognitive activity.

In summary, evolution has discovered the brain as a solution to the movement-control problem, since the emergent oscillatory dynamics of coupled muscles has severely limited complexity. It may suffice to pump blood through a multi-chambered heart, but it cannot control arm movements during tree climbing. Higher and higher layers of control evolved to realize more advanced sensing and acting, yet their ultimate coordination retains a 40 Hz motoric basis. Throughout this ascent in complexity, brains have redimensioned to fit evolving body types via an internal competition for cranial space and synaptic connections. In the transition from large-bodied apes to humans, the massive reduction of sensory and motor targets combined with a relatively constant cranial size reduced the relative demand for primary sensory and motor neurons, leaving extra space for higher-level structures, such as the prefrontal cortex, whose expansion and takeover of many control tasks arises naturally from a few genetic changes in developmental timing along with the principles of *late equals large* and *large equals well-connected*. Finally, the PFC appears to be a key prerequisite to symbolic reasoning and cognition.

However, the symbol-processing brain is merely a re-dimensioned sensorimotor brain that is partially exapted for cognitive endeavors. The complexities of sensing and acting are certainly no less problematic for us than for our distant mammalian ancestors, so the brain is still a sensorimotor controller, but with impressive reuse possibilities. The key questions for SEAI concern *what* is reused and *how*.

### 3 Prediction: The Core

Motion control - for all but trivially slow motions of simple appendages - demands an ability to predict the immediate future states of body and surroundings. Elementary control theory proves that delays in sensory feedback cause regulatory instability, and living organisms have sensory processing abilities that are too slow to properly complement their motor abilities. Hence, to rely on processed body-environment feedback after each action in order to select the next action is too inefficient, by several orders of magnitude, and would produce only very awkward motion.

To combat this mismatch, control theorists often add predictive models, such as Kalman filters, into their systems. Given a current sensorimotor context, these predict the next such context and use this estimate as the basis for further adjustments. Since the predictor runs internally, it produces an estimated future state long before the sensory system can provide the actual state. If the predictor is accurate, smooth control occurs. If not, the predictor must self-adapt based on differences between the predicted and actual states. This learning constitutes a second-order control problem.

Neuroscientists [58, 37] generally agree that the brain needs similar predictive abilities, and areas such as the cerebellum [58], basal ganglia [30] and hippocampus [20] are often cited as centers for the acquisition and use of these models. However, detailed explanations are far from consistent across the different research groups. This stems from both the general difficulty of empirical neuroscience and a theoretical mismatch between our own conceptions of situations, actions, causes and effects versus cerebral dynamics, which almost certainly uses other currencies in formulating a functional understanding of the world. Brains appear to build associations between sensorimotor contexts, i.e. states of both body and nearby environment, which we may interpret as *causal* or *common sense* knowledge.

In short, the brain must have elaborate predictive mechanisms and these may underlie our basic common sense. However, the exact associations between our explicit, verbal understanding for words such as *slap* and *shatter* and our *feel* for them are very unclear. Yet, it seems that the ultimate success of SEAI's vertical scaling to cognition critically depends upon the existence of these links, or, at least, the possibility of realizing them in artificial intelligences.

In fact, Jeff Hawkins, computer scientist and founder of the Redwood Neuroscience Institute, attributes many of AI's shortcomings to a failure to recognize prediction as the essence of intelligence [24]:

Intelligence and understanding started as a memory system that fed predictions into the sensory stream. These predictions are the essence of understanding. To know something means that you can make predictions about it. We can now see where Alan Turing went wrong. Prediction, not behavior, is the proof of intelligence.

### 3.1 Prediction and Learning

From the perspective of AI and Machine Learning (ML), prediction has an interesting appeal. In ML (and occasionally in neuroscience [11]), researchers differentiate between three types of learning: unsupervised, reinforced, and supervised. In the first, the agent learns recurring patterns without any tutoring input, while in the second case, the agent receives an occasional reward or punishment that essentially indicates that the net result of many agent activities was good or bad. Conversely, in supervised learning, the agent receives very frequent feedback that not only signals good/bad, but also indicates the action that **should** have been taken in an erroneous case. Clearly, the supervised approach, when feasible, can yield much faster learning of useful situation-action pairs than the other two methods.

Unfortunately, omniscient tutors are not available for the brunt of biological learning, and their assistance, when available, normally comes closer to reinforced than supervised. In non-human animals, the level of detailed supervision is even less. Yet, all animals are capable of learning useful behavioral information in a relatively short period. So if there is any element of supervision in this process, where is the feedback coming from? Prediction provides an answer.

Essentially, an agent can be its own teacher when doing predictive tasks. At time  $t$ , it predicts the future state of its body and immediate surroundings at  $t+1$ . Then, at time  $t+1$ , it learns the correct answer and can use that knowledge to adjust its own mapping from current to future states. This can be done almost continuously as the agent moves about the world, as researchers in evolutionary robotics have long recognized and exploited [45].

Furthermore, if the learning substrate is a network of simple, distributed processors (i.e., neurons), then the neural network literature shows that supervised learning via backpropagation works very well, although it requires numerous trials (i.e., hundreds or thousands). Also, the invention of Recirculation algorithms [27, 49] lend biological plausibility to backpropagation-like algorithms.

In sum, backpropagation is a very popular and effective method for training neural networks [25, 39], but it requires many trials and has traditionally suffered from biological implausibility. Recirculation remedies the latter problem, thus making it conceivable that real brains use a form of error-driven synaptic modification. Regarding the problem of excessive training requirements, predictive tasks clearly provide the necessary data. So if a neural-network-based agent is going to learn something useful, predictive models of world-body interaction may be the most amenable form of knowledge.

### 3.2 A Tempting Justification of SEAI

Prediction, which Llinas calls the ultimate function of the brain [37], is therefore the linchpin in a very tempting, motivating argument for continued SEAI research:

1. Complex movement requires an ability to predict the immediate future.
2. Prediction involves various functional mappings between and among states and actions.

3. These mappings constitute basic common sense.
4. Thus, the demands of movement provide the basis for cognition.

Although superficially straightforward, the neurological foundations for this line of reasoning are rather unstable, as discussed below.

## 4 Learning Models of the World

Though SEAI began with Rodney Brook's [6] direct assault on GOFAI, with the famous battle cries *intelligence without representation* and *the world is its own best representation*, two decades of experience with robotic embodiments of pure behavioralism reveal critical limitations of representation-free minds. In mammals, neuroscientists have discovered strong correlations between neural states and rich sensorimotor contexts [7], thus implying that models of the world and/or body-world coupling, are prevalent, albeit in forms that may not mirror those of textbook physics and geography.

Where are these models stored and how are they acquired? As a starting point, neuroscientists differentiate between two types of memories: declarative and nondeclarative [51].

A good deal of explicit, conscious human knowledge resides in the cortex, particularly the higher-level association regions of the parietal, temporal and frontal lobes. These declarative memories are of either a) specific objects or situations (episodic) or b) general concepts (semantic). Declarative memories are easily formed from single-exposure incidents, often those of emotional significance. However, the consolidation process is far from an observe-and-cache scenario. Rather, the hippocampus (HC) appears to store a reasonably rich snapshot of a situation and then gradually off-loads the memory back to those neocortical areas that stimulated the HC during the original experience. The transfer process generally takes several days or weeks [34, 38].

Amnesiacs, e.g., individuals with damaged hippocampi, can recall old memories (since they have been successfully logged in the neocortex), but experiences occurring only a few weeks or months prior to hippocampal damage are partially or fully corrupted, since their transfer to the cortex was interrupted. Basically, HC failure abolishes the ability to form declarative memories [52].

Procedural (or nondeclarative) memories are directly tied to sensorimotor activity. This knowledge is implicit in the sense-and-act machinery and appears inaccessible to consciousness, yet, it is believed to be the basis of our intuitive (i.e. commonsense) understanding of the world [51]. Whereas declarative memories are processed in the HC and later off-loaded to the cortex, procedural memories are learned *in place*, in areas such as the basal ganglia, cerebellum, and amygdala, along with sensory and motor cortices. Typically, procedural memories require multi-trial learning.

From the standpoint of acquisition, the two memory types are intuitive. A declarative (particularly an episodic) situation may only arise once, yet survival can be greatly enhanced by storing (at least a very abstract) representation of it. For example, a key landmark on the path to a newly-discovered temporary food source must often be remembered after a single encounter to support return trips. However, as described in [38], a new pattern cannot be forcefully added to an associative memory

(i.e. the cortex) without corrupting other memories. Hence, the continual, subconscious, *re-presentation* of the memory by the HC to the cortex allows its smooth integration with earlier patterns, just as a low learning rate is typically less corruptive in backpropagation.

Conversely, procedural skills lend themselves to frequent rehearsal; i.e., the situation arises many times and need not be cached in an HC-like organ. Here, the world truly can be its own best representation and the neural sensorimotor areas can gradually adapt to the recurring context on their own. Procedural skill learning includes operant and classical conditioning, and sequence learning. Interestingly, it also includes many forms of categorization: we can form many sensory classes without using the hippocampus. Thus, many *concepts* that our brains work with are neither hard-wired, nor explicit, but learned directly by cognitively inaccessible areas of the sensory cortex [52, 51].

At the molecular and cellular level, the formation of declarative and procedural memories are nearly identical [34, 51], but from a systems neuroscience perspective, they differ dramatically. Specifically, two key areas for procedural learning, the cerebellum and basal ganglia, perform supervised and reinforcement learning, respectively, while the hippocampus and cortex realize unsupervised associative learning of declarative information [11]. Supervised learning of predictive models is also believed to occur in the cortex [24].

Interestingly enough, back in the glory days of GOFAI, Winograd [57] developed the SHRDLU system for solving BlocksWorld problems. SHRDLU's knowledge, stored in if-then rules, was used to both a) perform block-stacking operations with a robot arm and b) plan sequences of such actions. Since these same rules could both drive motors and *think* about the world, Winograd made the very powerful generalization that procedural and declarative knowledge are interchangeable. A good many GOFAI researchers labored under this assumption, one which contemporary memory research now contradicts.

## 5 The Sensorimotor Control Hierarchy

When viewed as a sensorimotor controller, the brain exhibits a clear hierarchy of tightly interconnected modules [4, 34, 18], as shown in Figure 1. For simple reflex actions, signals travel from sensory receptors to the spinal cord and then immediately back out to the muscles. Activities involving proper timing or finesse often call on the cerebellum, which recommends actions based on learned associations between complex sensory contexts and motor responses that often enlist many muscles and body parts. Further up the hierarchy, the basal ganglia select and sequence discrete, high-level contexts that represent enduring (for seconds or minutes) sensorimotor or cognitive states. Finally, the neocortex, comprised of sensory, motor and frontal areas, provides long-term storage for the many contexts that trigger basal gangliar and cerebellar loops, and it too may engage in supervised learning of predictive models. Interestingly, the hierarchical arrangement of Figure 1 closely mirrors brain anatomy, with the cerebellum lying low along the brainstem, the basal ganglia up in the mid-brain, and the cortex covering the top [4].

In general, higher points in the hierarchy are more consciously accessible, while lower modules perform unconscious acts. However, at least the 3 highest layers are strongly involved in both

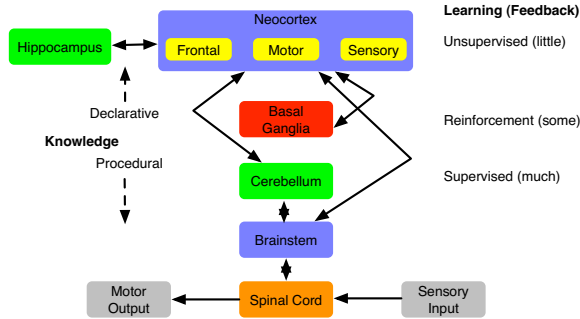


Figure 1: Mammalian sensorimotor control hierarchy

sensorimotor and cognitive activity. Hence, a lot of cognition a) utilizes the classic sensorimotor machinery, but b) is beyond conscious monitoring or control.

With respect to commonsense and predictive knowledge, their localization is highly ambiguous in the neuroscientific literature, due to both cross-experimental and cross-species differences. In general, we consider predictive knowledge as that involving associations between body-world states/contexts and other such states and/or actions. For example, knowledge that one context normally follows another is predictive, as are mappings from state-action pairs to consequent states.

We begin our investigation of predictive knowledge and its acquisition at the relatively low level of the cerebellum and work our way up to the hippocampus and neocortex.

## 5.1 Supervised Learning in the Cerebellum

Approximately half of the human brain’s neurons reside in the cerebellum[34, 4], which has long been known for its vital role in the learning and control of complex motions . As shown in Figure 2, the cerebellum receives convergent inputs from the sensory and high-level cortices. These are transferred from mossy fibers to granular cells, whose axons form parallel fibers (PFs) along the outer layer of the cerebellum. Purkinje cells (PCs) then *read* the parallel lines, with  $10^5 - 10^6$  synapsing on each PC dendritic tree. When PCs fire, they inhibit cells in the deep cerebellum, thus blocking or reducing particular muscular contractions.

Each PC receives signals from a single climbing fiber (CF); each CF contacts 1-10 nearby PCs. CFs relay signals from the inferior olive, which is stimulated by somatosensory inputs such as touch, temperature and pain. The sensory afferents of an inferior olive cell are located near the muscles controlled by the Purkinje cell’s of the olivary cell’s climbing fiber. Hence, the CF gives feedback directly related to the local action controlled by its PC. For example, if a muscular contraction causes a nearby joint to rotate excessively, the pain signal from the joint to the CF (via the inferior olive) would train the nearby PCs to reduce the strength of future contractions. Although the feedback does not provide the correct motion in erroneous situations, as with classic applications of the backpropagation algorithm, the signals are frequent and spatially focused, thus providing a reasonable facsimile of supervised learning [11].

Plasticity at the CF-PC synapse relies on post-synaptic long-term depression (LTD). When a CF forces a PC to fire strongly, those PC dendrites that were recently activated by parallel fibers undergo chemical changes that reduce their sensitivity to glutamate (the neurotransmitter used by PFs). Hence, the influence of those PFs on the PC declines [4].

Although there is good topographic correspondence between CF-PC pairs and the body areas that they serve, an individual PC does not control a single muscle, but, via its indirect connections to the motor cortex, is one of many influences to several higher-level neurons, each of which affects several muscles. In short, each PC participates in a *muscle synergy* [1]. Thus, the learning induced by a single climbing fiber involves a graded behavioral change in several functionally-related muscles, and proper motion control emerges from a multitude of these microscopic adjustments.

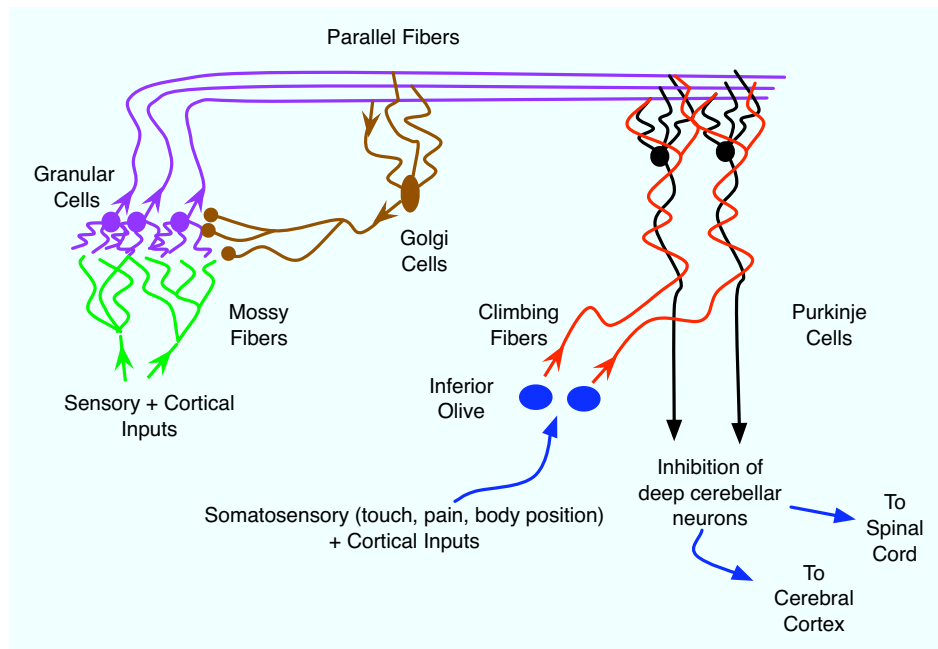


Figure 2: Basic organization of the cerebellum, an abstraction and combination of more complex diagrams in Bear et al. [4].

Several neuroscientists postulate causal models in the cerebellum, as summarized in [58]. In a nutshell, the cerebellum consists of many modular tracts or microzones, and each is believed to encode both a forward and inverse causal model related to a specific situation. The forward model computes expected future states when given the current state and action, while the inverse model computes an action when given a desired future state. As a feedforward controller, the cerebellum utilizes inverse models to provide motor-response recommendations. As a feedback controller, it uses predicted future states from the forward model to compute predicted errors, which are then used to generate the next round of motor signals. In familiar situations, the forward model circumvents the need for actual sensory feedback, whose time-consuming processing causes delays that degrade regulation. Also, [58] sees a cooperative arrangement in a model pair whereby the inverse model whose corresponding forward model provides the most accurate prediction of the future will, in turn, have higher precedence among all the inverse models when recommending the next action.

## 5.2 Reinforcement Learning in the Basal Ganglia

Animals accrue a strong selective advantage from learning associations between a) bodily and environmental indicators, and b) value-laden consequences, e.g. rewards or punishments. By predicting desirable or dangerous situations from their antecedent clues, animals can behave proactively, instead of merely reactively, to enhance survivability.

The basal ganglia (BG) are frequently viewed as the center of this *reinforcement learning* (RL) in mammalian brains [11, 29, 46]. Sketched in Figure 3, the BG are large midbrain structures that receive convergent inputs from many cortical areas onto the striatum. The striatal cells appear to function as a layer of competitive context detectors, since a) each neuron receives inputs from circa 10,000 cortical neurons, b) their electrochemical properties are such that they only fire if many of those inputs are active, and c) they have inhibitory connections to other nearby striatal neurons [28]. Strong evidence [55, 22] indicates that the BG are arranged in parallel loops wherein a striatal cell's inputs come from a region of a particular cortex, such as the motor cortex (MC). Its outputs to the subthalamic nuclei, substantia nigra and pallidum are eventually channeled back to the MC in the form of both action potentials (via the thalamus) and the neuromodulator dopamine. A great majority of these loops appear to involve the prefrontal cortex (PFC)[28, 34, 55], thus indicating BG contributions to attention, possibly as the mechanism for gating new patterns into working memory [22, 49]. In addition, the loops are not necessarily completely segregated [47], so an MC loop may include projections from and back to the PFC, for example.

Striatal modules consist of striosomal cells surrounded by matrix cells. The former project outputs to the Substantia Nigra (SN) either directly or indirectly via STN. Conversely, the matrix cells send signals to the pallidum, again directly and indirectly [30]. In Figure 3, notice that the direct paths are inhibitory, while the indirect are excitatory. Several prominent researchers [3, 29] characterize the BG as a combination of actor and critic, with the matrix cells and pallidal neurons as the actor's input and output ports, respectively, while the striosomes and substantia nigra demarcate the critic.

From an abstract perspective, the BG map contexts to actions. When a context-detecting matrix cell fires, it inhibits a few downstream pallidal neurons. In stark contrast to the striatum, the pallidum consists of low-fan-in neurons, most of which are constantly firing and thereby inhibiting their downstream counterparts in the thalamus [30]. When a striatal cell inhibits a pallidal neuron, this momentarily disinhibits the corresponding thalamic neuron, which then excites a cortical neuron, often in the PFC. The cortical excitation links back to the thalamus, creating a positive feedback loop that sustains the activity of both neurons, even though pallidal disinhibition may have ceased. Thus, the striatal-pallidal actor circuit momentarily gates in a response whose trace may reside in the working memory of the PFC for many seconds or minutes [28, 46].

Since the PFC is the highest level of motor control, its firing patterns often influence activity in the pre-motor (PMC) and motor (MC) cortices, while the MC sends signals to the muscles via the spinal cord. In addition, the sustained PFC activity provides further context for the next round(s) of striatal firing and pallidal inhibition that embody context detection and action selection, respectively. Via this recurrent looping, the basal ganglia execute high-level action sequences.

Of critical importance to the philosophical underpinnings of SEAI, the PFC is also the highest

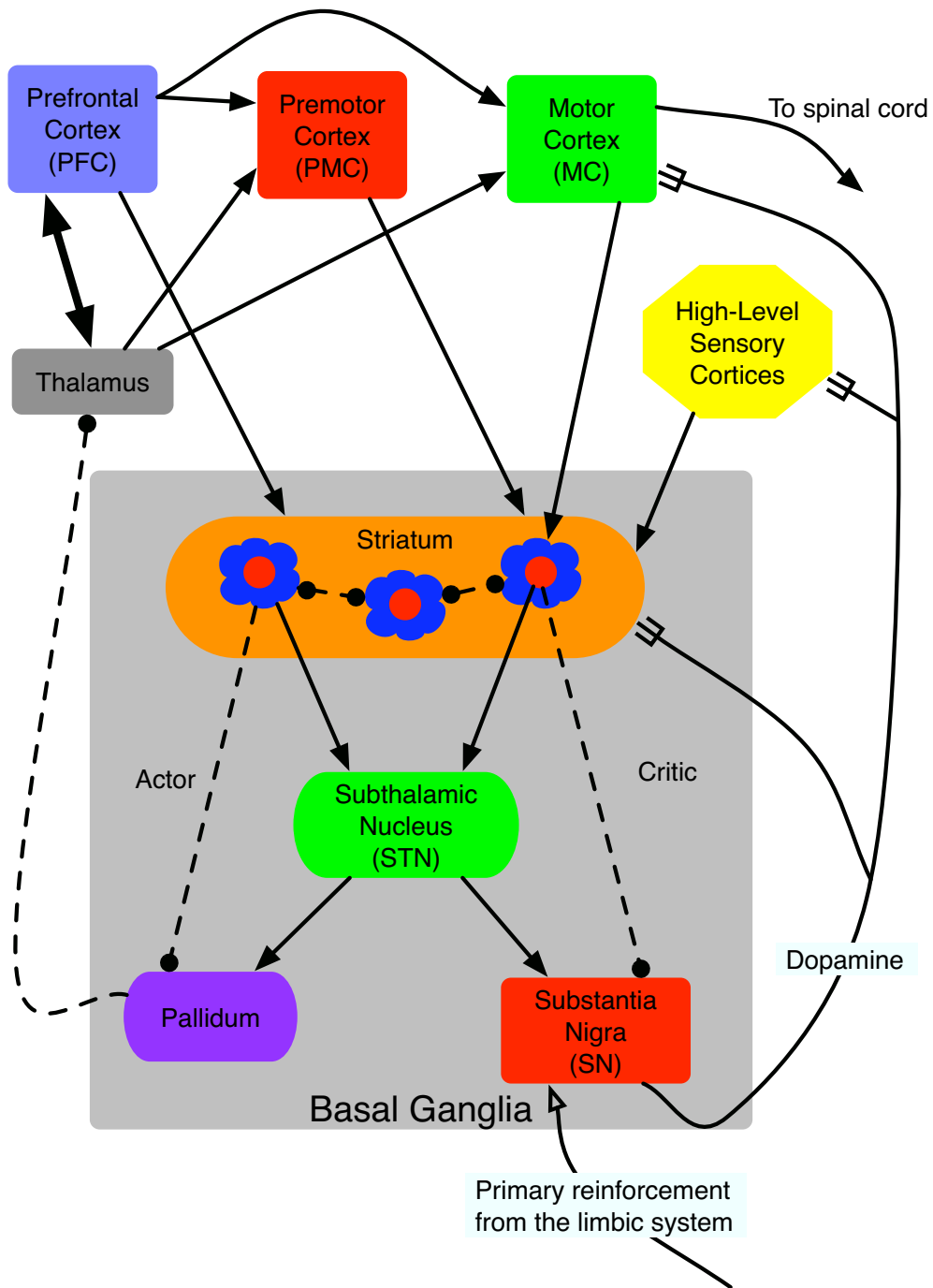


Figure 3: Basic topology of the basal ganglia and their main inputs, derived from text and diagrams in [28, 29]. Solid lines with arrows denote excitatory links, while dashed lines with circular heads are inhibitory. Forked heads denote diffuse neuromodulator (i.e. dopamine) transmission. Flowered figures in the striatum are matriosomes, whose cells project into the actor circuit, while a flower's center denotes striosomal cells, which project into the critic circuit. Some connections are omitted for ease of readability. For example, dopamine is transmitted to all cortical areas in the figure.

level of **cognitive** control [18]. Hence, PFC activation patterns can affect both motor activity and thought processes. It serves as a blackboard onto which many neural regions can log indicators of their current activity, which may then serve as inputs to other regions. Our conscious awareness of various aspects of a motor sequence may depend upon PFC activation. As motor sequences become more automatic, their control is believed to shift from a PFC-dominated BG circuit to a loop in which the thalamic outputs go directly to the motor cortex [11]. So these *compiled* sequences are no longer accessible to (PFC-mediated) conscious thought, although they still govern behavior.

The situation-action rules housed within the BG may comprise significant portions of our common sense understanding of body-environmental interactions, whether consciously or only sub-consciously accessible. Our smooth execution of both motor and cognitive tasks requires healthy BG. Major BG ailments, such as Parkinson's and Huntington's disease, cause significant cognitive impairments along with the well-documented physical deterioration [22].

However, the source of Parkinson's disease, and the key to reinforcement learning in the BG, resides not in the actor, but in the critic circuitry [3, 29] (see Figure 3). Here, dopamine (DA) signals from the Substantia nigra influence the synaptic plasticity of the regions onto which they impinge. DA acts as a second messenger that strengthens and prolongs the response elicited by the primary messenger.

For example, when a striatal neuron, S, is fired via converging inputs from the cortex, the primary messenger is the neurotransmitter from the axons of the cortical neurons (C) that recently fired. The immediate response of those S' dendrites (D) connected to the active axons is to transmit an action potential (AP) toward S's cell body. The summation of these D inputs will lead to S's production of a new AP. If dopamine enters these dendrites shortly after AP transmission, a series of chemical (and sometimes physical) changes occur which make those dendrites more likely to generate an AP (and a stronger one) the next time its upstream axon(s) produce neurotransmitters. Since the chemicals involved in this strengthening process are conserved, those dendrites that did not receive neurotransmitter may become less likely to fire an AP later on, even when neurotransmitters reach them. Thus, in the future, when the C neurons fire, the likelihood of S firing will have increased, whereas other cortical firing patterns will have less chance of stimulating S. In short, S has become a detector for the context represented by C. Without the dopamine infusion, S develops no bias toward C and may later fire on many diverse cortical patterns.

In unfamiliar situations, the SN fires upon receiving stimulation from various limbic structures, such as the amygdala (the seat of emotions), which triggers on painful or pleasurable experiences. The ensuing dopamine signal encourages the striatum to remember the context that elicited those emotions - the stronger the emotion, the greater the learning bias. Due to the biochemical temporal dynamics [29], the striatal neurons that become biased (i.e., learn a context) are those that fired approximately 100 ms **prior** to the emotional response. Hence, the BG learns a context, C, that **predicts** the reinforcing situation, R:  $C \Rightarrow R$ .

Furthermore, the topology of the critic network enables these predictions to regress backwards in time. In Figure 3, notice the indirect excitatory and direct inhibitory links from the striatum to the SN. The promoters act faster but for a shorter duration than the inhibitors. So when a striosome fires on a particular context, it will briefly stimulate the SN before blocking it for a longer period.

Consider the simplified scenario of Figure 4, in which an animal experiences a temporal series of sensorimotor contexts:  $X$ ,  $Y$  and  $Z$  before attaining the reinforcing state  $R$ . When this sequence first occurs, the attainment of  $R$  will be the first indicator of success, and the limbic reward signal will excite SN, causing dopamine-induced learning of context  $Z$  in a striosome, which sends fast excitatory and slow inhibitory signals to SN. The BG has learned the  $Z \Rightarrow R$  predictive rule.

On a later trial, the occurrence of  $Z$  will initially stimulate SN, and the ensuing dopamine will assist learning of a salient context immediately prior to  $Z$ , which is  $Y$ . Thus, another striosome is recruited to recognize a new context and notify SN upon its detection. The system has thus learned  $Y \Rightarrow Z \Rightarrow R$ . However, when  $R$  is attained, the limbic system still signals SN, but by that time,  $Z$ 's inhibitory signal has reached SN, thereby preventing further dopamine dissemination. Neuroscientists [34] have long known that dopamine signals only occur when a reinforcement is not expected, i.e., not predicted by a prior context. The temporal aspects of the biochemistry and the critic-circuit topology provide a clear explanation: when a context predicts a reward, its latent inhibition of SN blocks subsequent attempts to stimulate it.

Finally, on a still later trial, the occurrence of  $Y$  will stimulate SN, causing  $X$  to be encoded by a striosome and  $X \Rightarrow Y \Rightarrow Z \Rightarrow R$  (or  $X \Rightarrow R$ ) to be learned.

Since dopamine signaling is diffuse, the matriosomes and striosomes in a striatal module are both stimulated to learn. Hence, the critic not only learns to predict important states, but assists in the learning of proper situation-action pairs by the actor circuit. However, the neurophysiology of this potential interaction is not well-understood.

Under a slightly different interpretation of BG functionality, the striatum selects contexts consisting of situations plus proposed actions. For example, since various neurons in the prefrontal and premotor cortices are thought to represent high-level intentions to fire motor units, they may have similar connotations as striatal inputs. Gating of a context through the striatum and pallidum would indicate its approval, with portions of it (e.g. the action) held active in the PFC.

Again, with respect to the goals of SEAI, it is important to note that  $X$ ,  $Y$  and  $Z$  could be anything from a series of visual scenes processed by a maze-following rat to a sequence of motor acts during cross-country skiing (push with left foot, shift weight to right foot, glide and balance on right ski) to a string of intonations needed to pronounce a word to the steps of long division (multiply, subtract, transfer next digit). For a complicated task, such as playing a piano piece or performing a hockey slap shot, one cannot expect to learn the entire sequence from back to front. Instead, subsequences are probably learned incrementally, via intermediate reinforcements, and eventually coupled together.

### 5.3 The Hierarchy of Action Loops

In comparing the basal ganglia and cerebellum, the anatomy is somewhat ambiguous. Brain regions are densely interconnected, and contemporary *brain maps* are not yet conclusive. Functionality attributed to the basal ganglia in one paper is given to the cerebellum or hippocampus in another. Although the definitive story has yet to be written, some assumptions seem to have strong support:

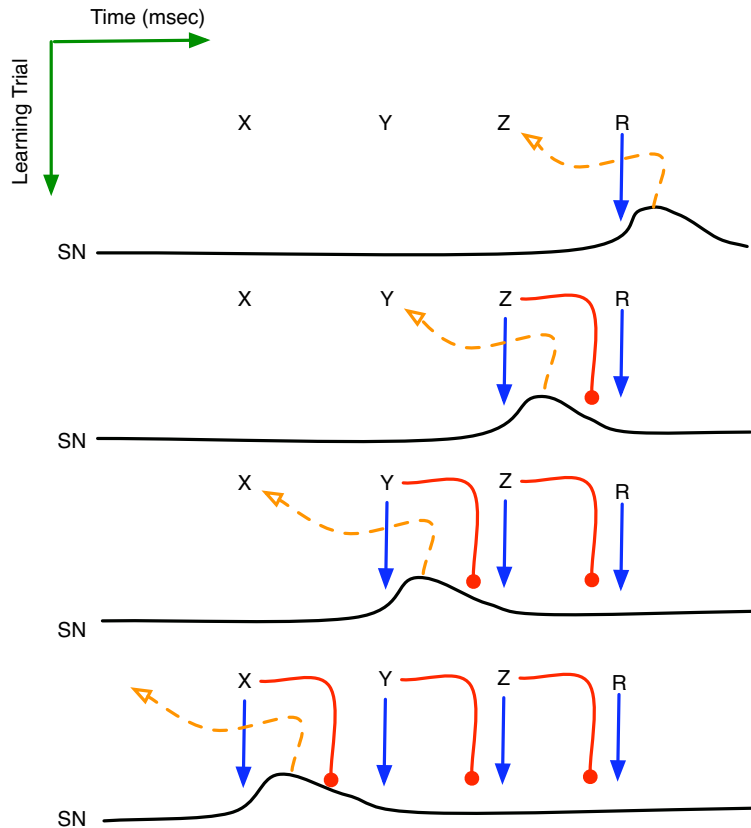


Figure 4: Reinforcement learning of sequence  $X \Rightarrow Y \Rightarrow Z \Rightarrow R$ . Horizontal plots are the time series activation levels for the Substantia nigra (SN). Solid arrows denote excitatory effects upon the SN, while round heads represent the delayed inhibition. Dashed arrows portray the learning of a new context governed by the SN's dopamine signal.

1. The prefrontal cortex is vital for working memory [21, 18, 46]. Neurons in the PFC can stay active for relatively long periods: seconds, even minutes, and are known to do so during psychological delay tasks in which subjects must remember a situation for some time before taking an action.
2. The basal ganglia have strong connections to the PFC on both input and output sides [22, 28, 4].
3. The cerebellum also connects to the PFC, but its main inputs are from the motor and sensory cortices, and its main outputs are to the spinal cord and motor cortex [4, 34].

From this, it seems safe to consider the aggregate of basal ganglia as a higher-level controller than the cerebellum. The latter is generally considered the center for timing of relatively specific actions, often based on direct sensory input, while the former selects among general contexts, allowing some to persist and influence cortical activity. Thus, the basal ganglia are driven more by complex internal contexts than by immediate sensory feedback. Some of these contexts can function as plan segments (in a very abstract sense) in that they involve **enduring** activation patterns in the PFC. By residing in the PFC, they can have strong effects upon premotor and motor areas, and by persisting, they can influence cerebellar actions for many seconds, as expected of plans, intentions, etc. The basal ganglia swap these plan segments into the cortex, with older segments gradually fading unless actively maintained.

From this perspective, the learning differences between the two regions make perfect sense. The results of cerebellar decisions are immediate and of short duration, so frequent feedback (via the climbing fibers) is appropriate for assigning credit to the most recent choices. In contrast, the basal ganglia's context choices can have broader temporal and spatial consequences such that immediate feedback is of less utility than an occasional (dopamine) reinforcement that provides a more holistic evaluation.

## 5.4 Unsupervised Episodic Learning in the Hippocampus

Often viewed as the center of long-term memory formation, the hippocampus (HC) resides in the temporal lobe and receives inputs from various cerebral regions, particularly the higher-level associational areas. As shown in Figure 5, the HC and surrounding regions perform a drastic compression (via high convergence) of information between the neocortex and area CA3, and a complementary expansion (via divergence) on the return path through CA1 and Subiculum. The topology of the HC proper is a main loop with several shortcuts from the EC to CA3 and CA1.

Only CA3 contains extensive recurrence, with each neuron connected to approximately 4% of the others [50]. This indicates that CA3 performs associative learning by standard Hebbian means: neurons that fire together wire together. The high convergence from a diverse array of neocortical areas onto CA3 hints of the holistic nature of these patterns. In rats, individual neurons in CA3 and CA1 are known as place cells [7], since they fire only when the animal is at a particular location, while in monkeys, they are called view cells, since they fire when the primate merely looks at such a location [50]. In either case, the cells represent a massive compression of sensory data.

The hippocampus' importance to long-term memory formation is well-established [51, 20], as is the fact that memories reside in the HC only until they can be off-loaded to the cortex for more permanent storage [38]. However, details of the actual learning process remain somewhat vague; the following scenario incorporates some of the more popular theories.

Since particular memorable situations may only occur once, episodic memory formation must involve a snapshot mechanism. However, McClelland et al. [38] point out that forcefully caching a new pattern in an associative network, via extensive synaptic modification (corresponding to a high learning rate in an artificial neural network), can corrupt pre-existing memories. Instead, new patterns should be repeatedly presented to the network in interleaved fashion, with only small synaptic changes occurring each time. The HC acts as the trainer for the cortex by a) temporarily caching snapshots of episodes via very fast Hebbian associative learning and then b) repeatedly re-presenting these patterns to the cortex until it too learns the associations. The basic topology and behavior of the HC indicates how this might work.

Consider a situation, S, summarized by the firing of 10,000 neurons, N1, in different high-level association cortices (many of which are in the temporal lobe along with the HC). The firings of N1 will coincide with the random background firing of certain nodes in pre-HC areas such as the perirhinal cortex, which, in turn, will coincide with random firings in the entorhinal cortex (EC), dentate gyrus (DG), CA3, CA1 and subiculum. Via Hebbian learning, the inter-layer synapses between these simultaneously-active neurons will be strengthened, as will the recurrent connections in CA3. Most HC synapses are capable of Hebbian Learning via long-term potentiation (LTP), but the recurrent collaterals within CA3 and the EC-DG, DG-CA3 and CA3-CA1 links exhibit very fast learning via rapid long-term potentiation (LTP). This learning rate far exceeds that of the cortex. So, although cortical neurons N1 will not initially form strong direct synaptic bonds with one another, their signals will gradually converge through several pre-HC and HC layers until, in CA3, a set of neurons, N2, that summarizes N1, will be co-active and immediately form strong synapses between one another.

Since a) the connections between the pre-HC layers are bi-directional, with both directions being strengthened during S's occurrence, and b) synaptic connections within the HC form a ring that begins and ends at the entorhinal cortex, a return pathway of enhanced synapses from N2 to N1 will form as well. Then, in the future, whenever N2 cells fire, the return pathway will stimulate the N1 cells. Through frequent re-stimulation, the N1 cells will gradually establish intra-cortical connections that will eventually consolidate the memory of S in the cortex. Detailed neural anatomy supports this slow-learning hypothesis, since HC afferents project onto distal (i.e. far from the cell body) portions of cortical dendrites, thus indicating a small effect for each return signal. Restimulation of the cortex may involve repeated waves of CA3 and CA1 activation during sleep.

Prior to long-term memory consolidation, recall could be achieved when portions of N1 send signals through the HC to CA3, stimulating part of N2, which the CA3 association network would then complete; the full N2 would then stimulate the rest of N1 via the return pathway. Rolls and Treves [50] point out that learning and recall are problematic to implement in the same naturally-occurring association layer, since a (complete) pattern to be learned could easily activate recurrent collaterals and be further expanded, based on the memories of other stored patterns. They believe that the two processes are segregated into two pathways, with the direct EC-CA3 connections used for recall, while the EC-DG-CA3 route achieves learning, since the DG neurons are believed to have stronger

influences on CA3 cells, thus overshadowing the effects of recurrent stimulation. However, in the same manner that dopamine signals may gate in activation patterns from the posterior cortical areas to the frontal cortex [46], neuromodulators that stimulate the EC-CA3 connections (during learning but not recall) without enhancing the CA3 recurrent collaterals could also explain the bi-modal behavior of CA3.

Evidence of place and view cells in rat and primate HCs, respectively, have motivated a cottage industry of ANN models of HC-based navigation, as summarized in [7]. Many of these involve implicit predictive knowledge in CA3 and CA1, wherein place cells fire before the animal arrives at the corresponding location. Place cells are mutually inhibitory and adapt to encode spatial contexts via competitive learning.

In one of the more popular models, Burgess et al. [41] propose a layer of goal cells (possibly in the subiculum) that receive inputs from many CA1 place cells. Goals represent very salient locations, often those involving reinforcements. As a simulated rat moves about, the goal cells fire at frequencies correlated with the rat's proximity to them, so navigation is achieved by choosing movements that increase the firing of focal goal cells.

Another interesting variant [33] posits CA3 as the site of predicted situations and CA1 as the site of real situations (via direct inputs from EC). Mismatches between the two drive learning in CA3 and thus improved predictions in the future.

From this plethora of diverse models and hypotheses, a few general, common themes emerge:

1. The pre-HC and HC compress and integrate huge amounts of sensory data into relatively sparse activation patterns in CA3 and CA1.
2. These patterns represent holistic high-level contexts which can aid the animal during repeated performance trials in the same environment, e.g. by allowing shortcut-taking during navigation.
3. There is an enduring memory for these contexts, whether in the HC or back in the cortex after off-loading.
4. Learning involves LTP in an unsupervised manner, although neuromodulators stemming from emotional experiences may be involved to help solidify memories of more salient contexts.
5. With its abundance of recurrent collaterals, most evidence points to CA3 as the center of this associative learning.

To generalize further, the HC seems instrumental in building memories, possibly quite detailed, of contexts that may later be recalled and analyzed by conscious, cognitive processes. This contrasts sharply with the basal ganglia and cerebellum, which also learn contexts, but these are often so tightly tied to sensorimotor machinery as to be unintelligible for explicit cognitive processing.

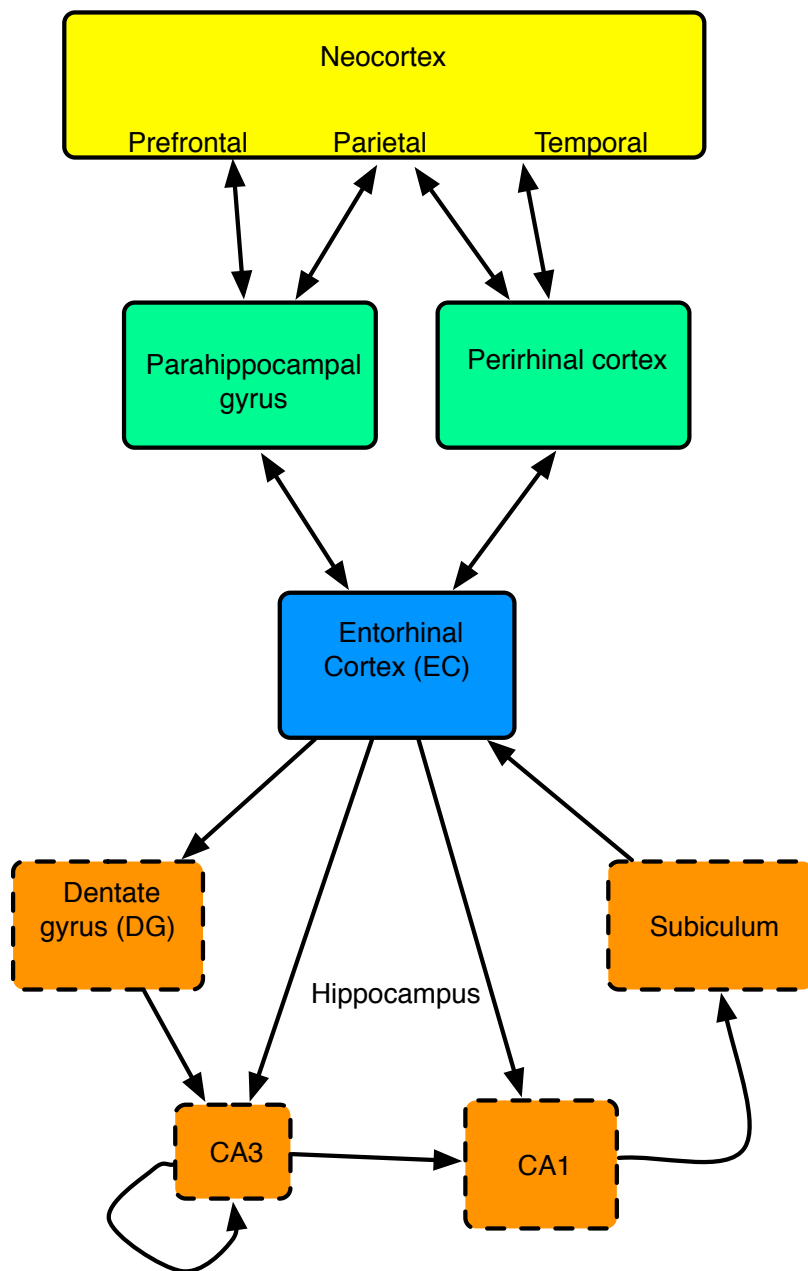


Figure 5: Basic topology of the hippocampus (boxes with dashed outlines) and surrounding areas. Box dimensions roughly illustrate relative sizes of neural populations in each area, and all connections are excitatory. Based on diagrams in [41, 50].

## 5.5 Learning Predictive Models in the Neocortex

The neocortex is the thin outer surface of the cerebrum, only a few millimeters in thickness and composed of 6 cell layers. These cells are vertically grouped into columns, with high intra-column connectivity between neurons and much weaker inter-column linkage. Hence, the columns function as semi-isolated modules. In mammals, it is convenient, and reasonably accurate, to view cortical columns near the back of the head as processors of primitive sensory information, particularly visual, while the more anterior columns process and represent higher-level concepts. Under this abstraction, bottom-up sensory-driven interpretation processes involve cascades of activation patterns moving from the back to the front of the neocortex. Conversely, top-down, memory-biased processing moves front to back, as shown in Figure 6.

Within each cortical column, the neurons capable of emitting the strongest and most influential signals to other columns are large pyramidal cells with cell bodies residing in the lower layers, 5 and 6. Axons emanating from these two layers tend to synapse on lower-level cortical columns, particularly motor neurons, and the thalamus, a subcortical structure known as a key relay station for sensory signals and believed to play a key role in integrating information. The dendrites of these large pyramidal cells extend up to layer 1, which is essentially a mat of axons coming from both higher-level cortical columns and subcortical structures such as the thalamus, hippocampus and basal ganglia. Signals from layer 1 reach the large pyramidal cells either directly, via the latter's dendrites, or indirectly via small neurons in layers 2 and 3.

Incoming axons from other columns can synapse with the large pyramidal cells at just about any point along their dendrites, from layer 4 up to layer 1. Proximal synapses (i.e., those close to the cell body, such as in layer 4) typically have a stronger effect upon the pyramidal cell's firing activity than will distal synapses at layer 1 or relay pathways through layers 2 and 3.

Of critical importance to understanding predictive-model learning in the neocortex is the fact that the axonal inputs from lower-level (i.e., posterior) cortical columns tend to enter higher-level cortical columns in layer 4, with some synapses also forming at layer 6 [24, 40, 34]. Thus, the low-level inputs form synapses near the cell bodies of the large pyramidal cells, whereas the inputs from higher-level (i.e., anterior) columns normally connect via layer 1. The immediate implication is that low-level sensory signals, which essentially represent the organism's current sensation of *reality*, have a stronger influence upon a cortical column than do the high-level thoughts (i.e. predictions) that often bias perception.

From the viewpoint of synaptic electrophysiology, the acquisition of predictive models within this hierarchical network of cortical columns has a very plausible explanation based on bi-modal thresholding. Artola et al. [2] have shown that weak stimulation of neurons (in the visual cortex) leads to long-term depression (LTD) of the synapses that were active during this stimulation, while stronger stimulation incurs long-term potentiation (LTP) of the active synapses.

Three learning cases are worth considering with respect to a particular cortical column, C, and its low-level sensory inputs, S, with synapses in layer 4 of C, and its high-level predictive inputs, P, with synapses in layer 1 of C. See Figure 6 for a graphic overview.

First, if S is active but P is not, then the effects of S on C's large pyramidal cells will produce a high

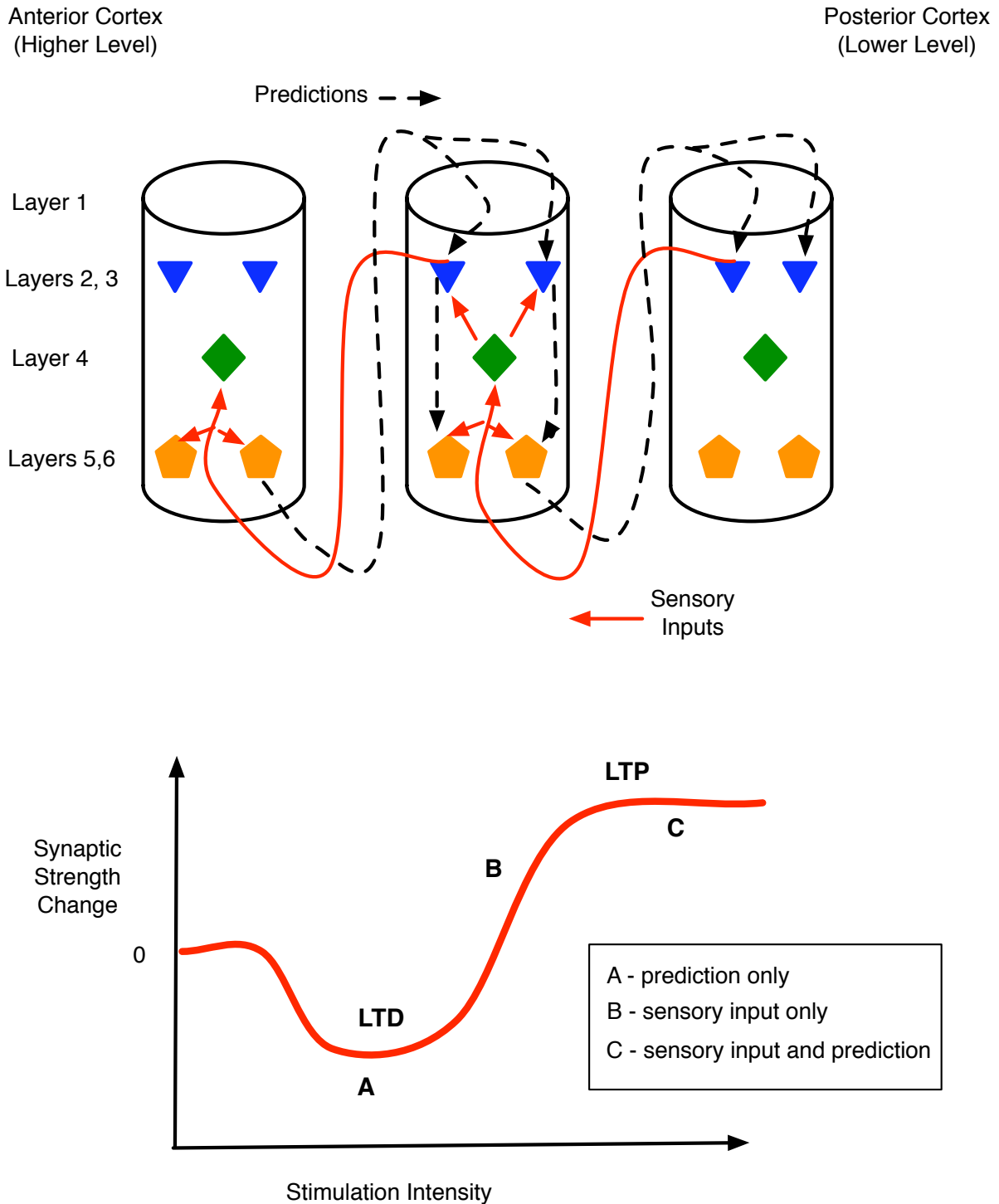


Figure 6: Top: Abstract view of cortical columns and top-down versus bottom-up information flow. Bottom-up flow (solid lines) goes from layers 2 and 3 of the sending column to layer 4 of the higher-level column, but with additional synapses onto the large pyramidal cells in layers 5 and 6 (pentagons). In the top-down pathway (dotted line), large pyramidal cells output to layer 1 of lower-level columns, with the signal eventually reaching layers 5 and 6 via either the layer 2-3 relays or directly via long dendrites from the large pyramidal cells. Bottom: Changes in synaptic strength as a function of post-synaptic stimulation intensity.

enough firing rate in those pyramidal neurons to incite LTP of the S-to-C synapses. Hence, C will learn to recognize certain low-level sensory patterns.

Second, if both P and S provide active inputs to C, then an even higher firing rate of C's layer-5 pyramidal neurons can be expected, so LTP of both the S-to-C and P-to-C synapses should ensue. In essence, the predictive and sensory patterns create a meeting point in column C by tuning the synapses there to respond to the P-and-S conjunction. In fact, after repeated co-occurrences of S and P, the synapses in C may strengthen to the point of responding to the P-or-S disjunction as well, in effect saying that it *trusts* the prediction P even in the absence of immediate sensory confirmation.

In the third case, when only P is active, the distal contacts of the P axons in layer 1 may only suffice to weakly stimulate C's large pyramidal neurons, thus leading to LTD: a weakening of the P-to-C synapses. Hence, future signals from P will not suffice to fire column C's pyramidal neurons, and thus P's predictions will not be propagated through C in the absence of verification from S. In short, the system learns that P is not a good predictor of S.

Another key aspect of this model concerns temporal relationships. Assume an initial sensory scenario,  $S_1$  at time  $t_1$ . This will propagate up the cortical hierarchy but will also evoke top-down predictive signals in layers that house expectations associated with  $S_1$ . For these expectations to propagate further down the hierarchy, they will need to match new sensory data. Due to the inherent time lag in neural signalling pathways, the predictions related to  $S_1$  will need to match up with ascending sensory data from situation  $S_2$ , which occurred at a slightly later time,  $t_2$ . Hence, the natural time delays in the system will insure the learning of predictions between states at time  $t_1$  and time  $t_2$ , thus neatly corresponding with our general conception of *causal* knowledge.

To complete this model of predictive learning, the branching factors of bottom-up versus top-down pathways are important to consider [24]. In general, the number of cortical columns decreases in moving up the processing hierarchy. Hence, bottom-up pathways appear convergent in that many primitive sensory neurons feed into the same higher-level neurons. Likewise, top-down pathways appear divergent, with one associative neuron signalling many lower-level columns. Hence, a *predictive* high-level neuron, P, may initially supply many lower-level neurons, in effect encoding an expectation that many different sensory patterns will be active when P fires. Through experience, many of these divergent connections will be pruned as their synapses weaken due to unfulfilled expectations and the resulting LTD.

Over time, a network of cortical columns with the topology and learning mechanisms described above will adapt its synaptic strengths to form a system that can both interpret sensory data and use top-down expectations to a) complete partial sensory states, and b) predict future states. In fact, under the general view of a context as an amalgamation of related information with some degree of temporal scope, the completion of a partial context could essentially involve the recall of future states associated with states closer to the present. Hence, most acts of pattern completion embody prediction as well, and the cortex seems particularly adept at this task.

## 6 Knowledge Compilation within the Hierarchical Prediction Machine

By tying the different brain regions and different learning methods together, we can begin to understand skill compilation, from declarative to procedural knowledge. However, in so doing, we also recognize how unnatural, and in many cases unnecessary, the inverse process seems.

Consider some physical task,  $T$ , being learned for the first time, such as shooting a basketball in a large stadium with fiberglass backboards. The player has presumably shot baskets with opaque wooden backboards in small gyms, above garage doors, etc., but never in a professional stadium. The key differences will lie in the visual contexts that both trigger and accompany various shooting motions. Training time is needed to learn these new contexts and seamlessly link them to the player's existing motor patterns for shooting.

Using Hawkins' model [24] of cortical functioning, the initial exposure to the new basketball court will primarily illicit surprise from the brain. Sensorimotor contexts at time  $t$  will not be easily predicted by contexts at time  $t-1$ , and thus the unique/surprising firing patterns will propagate all the way up the cortical hierarchy to invoke the hippocampus, which will take snapshots of these new contexts and gradually off-load them to the cortex in the course of days or weeks.

Imagine one such context,  $C$ , now safely encoded in synapses relatively high up along the cortical hierarchy, presumably spanning the prefrontal cortex and portions of the parietal and temporal lobes. At this early stage, the player is probably consciously aware (i.e. attending to) several aspects of the situation, hence the involvement of the PFC. If  $C$  is frequently activated by continued shooting practice, the connections between the thousands of neurons comprising  $C$  will become stronger such that many subsets of  $C$  will easily excite the remaining neurons to complete the pattern. In effect,  $C$  becomes a stable attractor in activation space.

As the attractor becomes more stable, it is more easily completed and less easily disrupted. This means that  $C$  requires less input from the prefrontal cortex to a) hold  $C$  stable, and b) perform the necessary motor movements, such as eye saccades, to *gather* information needed to complete  $C$ . Basically, the cortical predictive mechanisms would learn to complete  $C$  by assumption, rather than by additional sensory input.

As  $C$  becomes more frequently activated due to its attractiveness in activation space, it will often co-occur with dopamine-signalling reward situations, due to the emotive value of actually making baskets. Thus, a context detector for  $C$ ,  $D_C$  will emerge in the striatum, and the presence of  $C$  in the cortex will trigger  $D_C$ , thereby initiating a particular type of shot, such as a jump shot from approximately 5 meters in front of the basket (i.e., near the foul line).

At this point, the action will have become fairly automatic: a few important sensory cues will move up the cortical hierarchy, such as a rough judgement of the distance to the basket, an assessment of obstacles (i.e. opposing players), the position of the feet and hands and their recent movements, etc. This partial pattern will trigger assumptions/predictions in cortical columns that will complete  $C$  and activate  $D_C$ , without engaging much of the frontal cortex.

Of course, patterns at lower levels in the cortical hierarchy will also be active during shooting, but they will not form stable attractors at such an early stage. Since it resides at a relatively high level,  $C$  has the advantage of hippocampal involvement, and thus, accelerated learning via consolidation. The majority of hippocampal afferents and efferents involve high-level sensory and motor associational areas. Note that the hippocampus has been shown to be non-essential for basic skill learning in various tightly-controlled laboratory scenarios [51], but it would seem essential in learning a complex task in a new environment full of important sensory clues, just as it strongly affects rodent performance in maze-navigation tests [51, 7].

With continued practice, stable attractors would arise at lower levels of the cortical hierarchy, as sensory cortices began to learn new invariants of the basketball-stadium scenario. These would be patterns that were not previously discovered in other everyday situations, such as particular combinations of visual edges that might indicate direction and angle to the basket. Many such patterns would be inaccessible to consciousness, so the player could not necessarily explain why he was able to toss in 3-point shots even without a clear view of the hoop itself.

Furthermore, some of these patterns might engage the cerebellum. Consider the case of learning to release a jump shot at the apex of the jump. Initially, the player would need to consciously attend to the jump, estimate the apex, etc. With practice, the cerebellum would learn the timing constraints such that the initial upward acceleration of the body would engage a particular set of parallel fibers that a) connected to purkinje fibers controlling hand motion, and b) had an implicit time lag that delayed the hand motions until the optimal body height.

Consider one such low-level context,  $C^*$ , whose activity often coincides with that of  $C$ , i.e., much of the detailed information embedded in  $C^*$  is subsumed by the abstraction represented by  $C$ . Although it will take longer for  $C^*$  than  $C$  to be learned and thus to become a stable attractor, once formed,  $C^*$  will activate faster than  $C$ , due to its lower position in the cortical hierarchy.

Returning to the earlier discussion of reinforcement learning in the basal ganglia, recall that when a context,  $Z$ , that predicts reward,  $R$ , is learned, the dopamine signal is evoked by  $Z$  but no longer by  $R$ ; and precursors to  $Z$ , such as  $Y$ , can back up dopamine secretion even further in time. In the basketball example, context  $C$  would become associated with the emotional reward of making baskets, and dopamine signalling could be expected to commence upon activation of  $D_C$ . In turn, this would lead to the formation of striatal detectors for contexts that preceded  $C$ , contexts such as  $C^*$ .

Thus,  $C$  and  $C^*$  would both become detectors of the shoot-from-5-meters scenario, but  $C^*$  would activate first, thus engaging the basal ganglia in the 5-meter-jump-shot activity. Context  $C$  would become active slightly later, but the striatal detector  $D_{C^*}$  would already be inhibiting  $D_C$  as a part of the normal intra-striatal competition, so  $D_C$  would remain silent.

Over time, the synapses between  $C$  and  $D_C$  could easily atrophy, leaving  $C^*$  as the only initiator of the 5-meter jump shot. Thus, the pattern-recognition skills necessary to invoke the 5-meter shot would have become very fast, but very low level. Although  $C$  might remain a stable attractor, it would normally activate *after the fact* and thus have relatively low salience for basketball shooting. Moreover, in Hawkins' cortical model,  $C$  would not activate at all, since all predictions would be confirmed at the level of  $C^*$ , thus inhibiting propagation up the hierarchy. Disuse would lead to

the deterioration of the inter-C synapses, and the absence of C as a stable attractor would make it difficult to explain or reason about the activity.

A return to C-level processing might only be necessary in unusual situations where C\* no longer applies, such as games against exceptionally tall opponents, where the jump shot might have less chance of being blocked if the ball is released sometime during the ascent, prior to the peak height, in order to surprise the defender. The recalibration of cerebellar circuitry for this strategy might initially require some conscious C-level effort.

Along with illustrating possible interactions between the hippocampus, cerebellum, basal ganglia and neocortex during learning, this example has important implications for the declarative-procedural discussion: the dynamics of striatal context detection and competition, along with the temporal aspects of dopamine secretion and reception, give a clear competitive advantage to fast, early, low-level context detectors for skill acquisition. And, not incidentally, a faster response is conducive to better performance in most activities. So in the absence of extreme environmental variability, the tendency to compile decisions into subconscious circuitry will dominate the need to maintain declarative control of the process, and the *conscious hooks* into the activity may disappear. Only errors in normal behavior will necessitate cognitive involvement in an activity, and with experience, these erroneous cases will become a rarity.

In short, although our brains may use commonsense knowledge at all levels, there are many advantages to compiling it away into the subconscious and few reasons for retaining the *declarative source code*. Only by maintaining a need for the high-level representations is there a reasonable chance of its persistence. Hence, a basketball player may have little ability to explain proper shooting technique, while a coach could give a step-by-step account due to his background studying basketball fundamentals, but importantly, not necessarily due to his previous playing experience.

Interestingly, many explanations elude even seasoned experts. Gladwell [19] cites numerous examples of gurus in sports, psychiatry, and other endeavors who have nearly perfect predictive abilities based on very limited sensory data - such as being able to determine whether a tennis player will hit a service fault or whether a couple will remain married - but cannot explain their choices. The decision-making rules clearly exist but are not consciously accessible.

## 7 The Barrier to Reuse

One striking difference between the cerebellum and basal ganglia versus the cortex (and CA3 of the hippocampus) is the high density of excitatory recurrent intra-layer connections in the latter. These support associative memories in which a) partial patterns can be completed via spreading activation, and b) stable attractor firing patterns can emerge and persist. Stability seems prerequisite to the focus of attention that underlies conscious cognition, while spreading activation seems critical for recalling associations between concepts in memory, that is, for reasoning.

Conversely, the basal ganglia and cerebellum consist of thousands of parallel tracts through a series of layers which have inhibitory intra- and inter-layer recurrent collaterals. These areas appear designed to map cortical contexts onto other contexts, which may involve motor acts, intentions to

perform such acts, etc. Although the basal ganglia may be involved in pattern stabilization via their interactions with the frontal cortex, neither region seems wired to support spreading activation.

Using these two regions, one can sing a song perfectly but cannot, without the aid of external media (e.g. paper and pencil), recall and compare the 7th and 23rd lines. Each word or phrase of the song may be stored in the cortex, but extraction is mediated by the combination of preceding cortical context (declarative) and basal gangliar wiring (procedural). There are similar restrictions on commonsense skill knowledge for writing, riding, throwing, etc. Essentially, the glue that holds pieces of complex cortical patterns together (across time) resides in a procedural bottleneck - only a few striatal neurons fire simultaneously - such that the whole pattern cannot be activated and analyzed as one stable cortical image. Hence, our causal knowledge may involve disjoint abstract cortical snapshots whose combinations generally evade conscious contemplation.

Convincing evidence of the strong barrier between declarative and procedural knowledge comes from patients with hippocampal damage (i.e., amnesiacs), who can learn a wide variety of complex procedural tasks as well as normal patients, but who have no awareness of what they have learned [51]. For example, they may learn to classify a set of training cases, but afterwards, they cannot *recognize* any of the cases.

Classification often occurs in the striatum of the basal ganglia, where striatal cells learn invariants among similar contexts via a competitive learning process. However, categorization can occur very quickly, using only lower levels of the cortical hierarchy. These levels are less accessible to PFC-mediated attention and have more continuously-changing patterns (i.e., less stability) than higher levels. In short, a pattern could easily stimulate the striatum and trigger a correct yes/no response, thereby exhibiting a cognitive competence, but without engaging conscious declarative thought.

In general, declarative and procedural memories are not the sole provinces of cognition and sensorimotor behavior, respectively. For example, one of the procedural tasks on which amnesiacs and controls perform equally well involves dynamic staffing of a simulated sugar factory to achieve an optimal production level . This clearly qualifies as cognitive, since it involves non-trivial arithmetic reasoning. But, somewhat surprisingly, it lacks a strong declarative component and thereby remains fully accessible to amnesiacs [51]. In addition, brain imaging reveals high activity in procedural areas such as the basal ganglia and cerebellum during pure thinking tasks.

Similarly, sensorimotor tasks surely require the long-term storage of certain explicit contexts, particularly of detailed sensory cues that help initiate tasks, after which the brain may run in a purely procedural auto-pilot mode. For example, in navigation, an intelligent animal would not want to launch into a particular movement plan until it had a fairly definitive recognition of a key landmark, as opposed to performing an impulsive reaction to a few low-level stimuli. This careful recognition requires the invocation of declarative memories.

## 8 Insights from Alternative Computer Metaphors

Although GOFAI could not successfully cross the chasm between declarative and procedural, human brains frequent this route. As discussed above, skills begin as conscious, declarative activities

involving the frontal cortices and often become *compiled* into purely procedural routines, which run much more efficiently but are less adaptive to novel situations. Is it unreasonable to imagine traffic in the other direction: from procedural to declarative, should the need (e.g., coaching or teaching) arise?

Under GOFAI's legacy to cognitive science, the computational brain metaphor, the answer appears to be "No", or at least, "Not easily". Although a high-level computer language is compiled into machine language, the process is non-invertible. Of course, machine language is easily converted into Fortran or Pascal, but the mapping is one-to-many, so recovering the original source code is nearly impossible.

So one could argue that the mapping between neuronal patterns in subconscious sensory and motor areas and corresponding high-level patterns in the frontal cortices is one-to-many and thus hard to reverse compile, but the computational analogy still misses the mark. To reverse the analogy, if computers worked like the brain, then program compilation might work as follows:

- Begin with source code S.
- Run S via an interpreter. When errors occur, modify S to avoid them.
- When a particular module, M, of S, seems to be working well, compile it into machine language.
- Test S where M is compiled but the remainder, S-M is not and thus must be interpreted on each run.
- Gradually compile and test more modules until all of S is in machine code.
- If, at any time, errors are found in a compiled module, M, create a hybrid module, H, that includes both M and error-corrected source for M,  $S^*(M)$ . Whenever H runs, it must stochastically choose between M and  $S^*(M)$  in generating an output. Over time, the probability of consulting M should go to zero and  $S^*(M)$  can be compiled into  $M^*$ , the updated, machine-code version of M.

The final step reflects the fact that subconscious routines are difficult to unlearn. This approach, except for the last step, is not totally foreign to software engineering, but only a few high-level languages, such as Lisp, can run in interpretive or mixed interpretive-compiled mode.

Perhaps a more fitting computational metaphor comes from large-scale software design, where one or a few high-level managers organize and instruct a group of lower-level programmers. The managers determine what needs to be built and the tests to perform on the successive code versions, but during testing, it is the programmers who monitor the step-by-step behavior of the program modules and perform the minute-by-minute source-code debugging changes. At day's end, the manager may get a report consisting of a) overall test performance and b) very high-level descriptions of program changes (e.g., the spell-checker now handles hyphenated words). This information will then influence managerial choices of future tests.

A similar metaphor comes from mentoring. In the same way that a teacher or coach educates by giving basic conceptual instruction and then exposing students to practice drills, the forebrain can help structure the environment for the procedural regions by putting the body in drill situations and steering a few high-level actions.

To consider this metaphor at two levels, a coach may instruct a poor foul shooter to take a few hundred shots from the foul line. In turn, the player uses her forebrain to decide to stand at the foul line and shoot. However, neither coach nor forebrain has a deep understanding of the internal changes realized by their subordinates (player and subconscious motor-control brain regions, respectively) during learning. The coach only sees the players taking shots and, hopefully, over time, some improvements. The player feels herself focusing on various movements (e.g. keeping the elbows in, flexing the wrists, bending the knees, etc.) but cannot understand exactly what the basal ganglia and cerebellum are doing to enable or improve these movements. Only by observing overt behavioral improvements can the coach and player's forebrain confirm the success of their training regimes. Very little else of the distributed change gets back to the higher level.

Similarly, in animal skill learning, conscious activity determines the general context to which a procedural module is exposed, but the exact location and detailed nature of the changes is unconscious and intrinsic to the procedural circuits. LTP and LTD are pinpoint synaptic changes based on specific local information plus very vague global broadcasts indicating only that *something* significant has just occurred. Thus, the forebrain cannot monitor procedural learning at any detailed level.

Basically, evolution designed brains this way. Higher-level cortical areas arose to *enhance* procedural activity, not to understand or explain it; and this provides a disconcerting precedent for SEAI.

## 9 Conclusion: SEAI and the divide

Given its direction of approach, SEAI may have an even harder time crossing the barrier between procedural and declarative knowledge than GOFAI did. Although natural evolution shows that brains evolved to control sophisticated sensorimotor activity, and that cognition arose by borrowing that same machinery, there is little evidence of direct reuse of procedural knowledge for declarative purposes.

Not only does the gap appear real in our own brains, but it has a conceptual depth that may be hard to circumvent in even artificial systems. The very essence of episodic/semantic knowledge and its acquisition contrasts so strongly with that of physical skills that any single mechanism to accomplish both might be unrealistic for tasks more complex than SHRDLU's block-stacking.

As a general assessment of cognitive incrementalism, Clark writes [8]:

It may indeed be that the neural mechanisms of higher thought and reason are fully continuous with mechanisms of on-line action control. But it may be quite otherwise. Most likely, what we confront is a subtle and complex mixture of strategies, in which new kinds of information-processing routine peaceably coexist with, and at times exploit

and coopt, more primitive systems...In sum, we must treat the doctrine of cognitive incrementalism with great caution. It is a doctrine that is both insufficiently precise (concerning what is to count as continuity, incremental change, etc.) and empirically insecure. Attention to the shape of nature's solution to basic problems of real-time response and sensorimotor coordination will surely teach us a lot. Whether it will teach us enough to understand mindfulness itself is still unknown.

In the very least, by understanding the neural processes of sensorimotor control, we should have a good biologically-based start in building cognitive controllers by exaptive means, since considerable neurological evidence now points to the use of traditionally motor areas for cognition as well [18, 30, 22, 34].

Still, this is a disappointing consolation given our original goal of the cognitive reuse of procedural commonsense knowledge. Those memories now seem so tightly intertwined with the body and its control that they cannot be consciously accessed, abstracted or otherwise manipulated. In effect, procedural memories *make no sense* with respect to abstract reasoning about the world, or even with respect to relatively specific situation-action pairs. They only have meaning relative to the specific sensory and motor *apparatus* of the organism. In their basic essence, they are abstraction free.

Since, by definition, SEAI systems must first crawl before they can contemplate, the designed or evolved architectures will naturally be optimized for sensing and acting, and as nature reveals, this entails a daunting cognitive impenetrability of the acquired common sense. Conceivably, transfer may occur via the environment: the agent could perform actions, observe results, and inductively form declarative causal representations. Observing other agents would also work, but clearly, the somatosensory feedback involved in one's own activity has powerful effects on both declarative and procedural memory formation. For example, when a particular movement causes pain, we often do consciously attend to the situation and learn explicit heuristics.

Societal approaches, as explored in [53], seem promising, since if animats must evolve to both perform tasks and to transfer behavioral tips, then their brains will not necessarily become optimized for solipsism. The demands of communication may force an early coupling between declarative and procedural knowledge such that common sense could in fact transfer directly from procedural to declarative realms.

Also, Squire and Zola [52] observe that subjects with functioning hippocampi form a parallel, auxiliary declarative representation (of no immediate performance benefit) while doing a purely procedural task. SEAI systems might utilize a similar mechanism, wherein both competitive classifiers in the striatum-like module and associative memories in a simulated CA3 are tuned during sensorimotor adaptation.

These parallel pathways are the key contribution of Sagita and Tani's [56] recent neural network models. One ANN maps situations to behaviors, while the other maps them to simple sentences. The two networks are correlated using a vector of reals (known as a *PB vector*) that is a) a common input to both networks, and b) modified along with the connection weights during gradient-descent training of each network. As a result, the ANNs synchronize and generalize enough to both a) produce sentences for activities that it has never previously performed, and b) carry out actions

described by sentences that it has never before seen. To paraphrase the authors, *the link between symbols and behaviors is implicit, dynamic and self-organizing, not hard-wired via an ad-hoc symbol-grounding theory.*

Sagita and Tani’s work is clearly seminal in addressing (albeit implicitly) the coordinated interaction of procedural and (somewhat) declarative modules. The PB vector weakly resembles the prefrontal cortex (PFC), since it caches a set of activation patterns that influence and are strongly influenced by activity in other areas. The next step might be to implement the PB vector as a third ANN with reciprocal connections to the other networks but relatively stringent gating preconditions for the entry of outside information into the working memory that it embodies.

All of this work [56, 52, 53] indicates that the best approaches to spanning the procedural-declarative divide are indirect ones. Systems essentially do work on both sides of the gap, building a series of loose connections between various pivotal procedural and declarative representations. Then, a pattern, A, emerging on one side drives the formation of its counterpart, B, on the other side. From B, many other patterns are spawned, only a few of which may map back to the other side.

The gap becomes less ominous in light of these indirect approaches, but only slightly. The canyon metaphor is used to imply that we cannot simply march from one side to the other, *hitting the ground running*. Careful and coordinated preparatory work on both sides is necessary, since each involves different machinery and constraints. The bridges that are built cannot carry complete representations, only coordinating information. In short, the indirect approaches should improve our optimism for eventually creating impressive artificial intelligences, but they should not diminish our awareness of nor respect for the extreme differences between procedural and declarative information.

There are interesting parallels between our research and the simulation and emulation theories of motor and perceptual imagery. In simulation theories [26, 32, 59, 31], the same internal (i.e., non-motor) neuronal patterns used to generate motor activity are employed when imagining those actions. Imagery thus involves running the internal neurons with deactivated sensory- input and motor-output systems. Emulation theories [23] postulate a separate predictor mechanism that runs parallel to the action-controlling circuits. This extra module serves as a fast, feedforward, causal predictor of the next sensory state of the world, given the current state and motor requests. This enables the action-generating module to begin preparing the next motor commands without having to wait for the actual sensory inputs, which, as discussed earlier, are observed and interpreted at rates that are much too slow to support effective motion control. Comparisons of actual and predicted future states are then used to modify the predictor module, as in a Kalman filter.

A superficial analysis would indicate that our ideas support emulation over simulation theories, since we claim that the procedural knowledge cannot simply be reused when reasoning about sensorimotor activity: a separate declarative representation seems essential. However, at the higher levels of motion control, our views are completely compatible with simulation theory, which assumes that activity patterns in preparatory cortical areas such as the prefrontal cortex and premotor cortex can be reused for sensorimotor imagery. However, if these preparatory patterns involve lower level structures such as the motor cortex and cerebellum, thereby making it more likely that they evade conscious attention, then in our view, they would need complementary declarative structures in order to support reasoning about action.

With respect to emulation theory, our work is somewhat orthogonal. For motion control, an emulator needs to run fast, faster than the combined motor and sensory tasks. In contrast, conscious declarative reasoning runs relatively slowly. Hence, the declarative models that we propose could not serve as stand-alone emulators for real-time motion control. Of course, they may indirectly affect motor activity by providing general constraints that persist during many rounds of motor-command choice, such as the beginner’s conscious efforts to keep a) the head down while swinging a golf club, or b) the elbow in while shooting a basketball.

In our view, fast emulators must work with activity patterns in a *procedural language* that is *understood* by motor neurons, while conscious reasoning about sensorimotor activity involves predictive machinery that runs on more abstract patterns. No single emulator is enough. The entire control hierarchy, with its many levels of abstraction (see Figure 1) seems necessary. We assume that its lower levels house basic sensorimotor controllers as well as emulators. As the tasks become more complex and thus action selection requires more detailed contextual information and possibly conscious reasoning, higher levels of the hierarchy are needed. These too include predictive machinery, but coded in a more declarative form. For example, knowledge of the form, ”Everytime I throw the ball to Frank, he drops it,” might be used to convince a player to pass to someone besides Frank. We assume that these declarative models are learned by declarative processes as well. For example, one might need to think about the game and recall the many passes that Frank dropped before formulating the behavior-influencing generalization that Frank is a risky passing outlet.

In general, our work has no direct conflicts with either simulation or emulation theories, but it centers on the question of how the predictive structures originate. To the emulation supporters, we might ask how the feedforward model is generated. Does it grow out of the basic action controller or does it arise independently? The simulation promoters would receive similar questions regarding the posited relationships between lower- and higher-level representations and how each may influence the emergence of the other.

With respect to these origins, we share Dreyfus and Dreyfus’s interest in skill acquisition and the multiple representational levels it seems to encompass [13]. However, as discussed earlier, our focus is less on the conversion of declarative to procedural knowledge and more on the inverse process and its neurological barriers, since these have direct relevance for cognitive incrementalism. In more recent work, Hubert Dreyfus [12] acknowledges that the subtleties of physical expertise cannot be articulated, as they are too tightly tied to the body and the intricate details of its control.

Beer’s work on minimally-cognitive tasks [5] might serve as a useful guiding philosophy when considering tasks such as the aforementioned sugar-factory controller, which are cognitive but largely procedural. The question then becomes how far you can stretch procedural mechanisms to perform cognitive activities in the absence of declarative representations. An analysis of these minimally *declarative* cognitive tasks might yield some indication of when purely procedural common sense breaks down and demands declarative encodings.

Finally, the cortical model described by Hawkins [24] has yet to be implemented. Although he admits to *cortical chauvinism*, a reasonably competent artificial intelligence would probably require the general predictive mechanisms of a cortical hierarchy along with a hippocampal structure, for consolidation without destruction of older memories, and a context/action-selecting region similar to the basal ganglia.

At present, contemporary SEAI animats perform action selection and little more. To successfully reason about a situation, the artificial brain should be capable of holding a concept in memory and allowing spreading activation to make behavior-influencing predictions of expected inputs, related concepts, etc.

Of course, in the end, SEAI may simply reach the opposite bank of a wide chasm that GOFAI discovered 20 years ago: knowing how and knowing about are non-interchangeable. Currently, however, SEAI appears motivated by an implicit belief that bottom-up sensorimotor agents will eventually scale to general intelligences, since worldly experience is the only way to achieve GOFAI's Achilles heel, common sense.

However, nature has stumbled upon a formidable barrier between the evolutionarily older procedural faculties and relatively new declarative abilities: organisms can *possess* common sense without being able to analyze or explain it. SEAI, with its focus on predominantly procedural controllers, appears headed down a similar path. Fortunately, neuroscience and psychology post several warnings of the procedural-declarative chasm along this route. Careful attention by SEAI researchers to this divide will hopefully lead to the design of appropriate bridges. However, we must avoid one of GOFAI's cardinal sins of assuming that separate cognitive faculties can be built in isolation over several years and then tied together in a weekend.

## References

- [1] N. S. A.G. BARTO, A. H. FAGG AND J. HOUK, *A cerebellar model of timing and prediction in the control of reaching*, *Neural Computation*, 11 (1999), pp. 565–594.
- [2] A. ARTOLA, S. BROCHER, AND W. SINGER, *Different voltage-dependent thresholds for inducing long-term depression and long-term potentiation in slices of rat visual cortex*, *Nature*, 347 (1990), pp. 69–72.
- [3] A. BARTO, *Adaptive critics and the basal ganglia*, in *Models of Information Processing in the Basal Ganglia*, J. Houk, J. Davis, and D. Beiser, eds., Cambridge, MA, 1995, The MIT Press, pp. 215–232.
- [4] M. BEAR, B. CONNERS, AND M. PARADISO, *Neuroscience: Exploring the Brain*, Lippincott Williams and Wilkins, Baltimore, MD, 2 ed., 2001.
- [5] R. BEER, *The dynamics of active categorical perception in an evolved model agent*, *Adaptive Behavior*, 11 (2003), pp. 209–243.
- [6] R. BROOKS, *Cambrian Intelligence: The Early History of the New AI*, The MIT Press, Cambridge, MA, 1999.
- [7] N. BURGESS AND J. O'KEEFE, *Hippocampus: Spatial models*, in *The Handbook of Brain Theory and Neural Networks*, M. Arbib, ed., The MIT Press, Cambridge, MA, 2003, pp. 539–543.
- [8] A. CLARK, *Mindware: An Introduction to the Philosophy of Cognitive Science*, The MIT Press, Cambridge, MA, 2001.

- [9] T. DEACON, *The Symbolic Species: The Co-evolution of Language and the Brain*, W.W. Norton and Company, New York, 1998.
- [10] K. DOWNING, *Artificial life and natural intelligence*, in Proceedings of the 6th Genetic and Evolutionary Computation Conference, Seattle, WA, 2004, The MIT Press, pp. 81–92.
- [11] K. DOYA, *What are the computations of the cerebellum, the basal ganglia, and the cerebral cortex?*, Neural Networks, 12 (1999), pp. 961–974.
- [12] H. DREYFUS, *Intelligence without representation - merleau ponty's critique of mental representation*, Phenomenology and the Cognitive Sciences, 1 (2002).
- [13] H. DREYFUS AND S. DREYFUS, *Mind Over Machine*, Free Press, 1982.
- [14] G. EDELMAN AND G. TONONI, *A Universe of Consciousness*, Basic Books, New York, NY, 2000.
- [15] D. FERRIER, *The Functions of the Brain*, G.P. Putnam's Sons, New York, 1886.
- [16] B. FINLAY AND R. DARLINGTON, *Linked regularities in the development and evolution of mammalian brains*, Science, 268 (1995), pp. 1578–1584.
- [17] K. D. FORBUS, *Qualitative reasoning*, in The Computer Science and Engineering Handbook, 1997, pp. 715–733.
- [18] J. FUSTER, *Cortex and Mind: Unifying Cognition*, Oxford University Press, Oxford, 2003.
- [19] M. GLADWELL, *Blink*, Little, Brown and Company, New York, 2005.
- [20] M. GLUCK AND C. MYERS, *Gateway to Learning: An Introduction to Neural Network Modeling of the Hippocampus and Learning*, Kluwer Academic Publishers, Norwell, Massachusetts, 1989.
- [21] P. GOLDMAN-RAKIC, *Working memory and the mind*, Scientific American, 267 (1992), pp. 110–117.
- [22] A. M. GRAYBIEL AND E. SAKA, *The basal ganglia and the control of action*, in The Cognitive Neurosciences III, M. S. Gazzaniga, ed., The MIT Press, Cambridge, MA, 2004, pp. 495–510.
- [23] R. GRUSH, *The emulation theory of representation: motor control, imagery, and perception.*, Behavioral and Brain Sciences, 27 (2004), pp. 377–442.
- [24] J. HAWKINS, *On Intelligence*, Henry Holt and Company, New York, 2004.
- [25] S. HAYKIN, *Neural Networks: A Comprehensive Foundation*, Prentice Hall, Inc., Upper Saddle River, N.J., 1999.
- [26] G. HESSLOW, *Conscious thought as simulation of behavior*, Trends in Cognitive Science, 6 (2002), pp. 242–247.
- [27] G. HINTON AND J. MCCLELLAND, *Learning representations by recirculation*, in Neural Information Processing Systems, D. Anderson, ed., American Institute of Physics, New York, 1988, pp. 358–366.

- [28] J. HOUK, *Information processing in modular circuits linking basal ganglia and cerebral cortex*, in Models of Information Processing in the Basal Ganglia, J. Houk, J. Davis, and D. Beiser, eds., Cambridge, MA, 1995, The MIT Press, pp. 3–9.
- [29] J. HOUK, J. ADAMS, AND A. BARTO, *A model of how the basal ganglia generate and use neural signals that predict reinforcement*, in Models of Information Processing in the Basal Ganglia, J. Houk, J. Davis, and D. Beiser, eds., Cambridge, MA, 1995, The MIT Press, pp. 249–270.
- [30] J. HOUK, J. DAVIS, AND D. BEISER, *Models of Information Processing in the Basal Ganglia*, The MIT Press, Cambridge, MA, 1995.
- [31] M. JEANNEROD, *Neural simulation of action: A unifying mechanism for motor cognition*, Neuroimage, 14 (2001), pp. 103–109.
- [32] D.-A. JIRENHED, G. HESSLOW, AND T. ZIEMKE, *Exploring internal simulation of perception in mobile robots*, in Proceedings of the 4th European Workshop on Advanced Mobile Robots, Arras, Baerveldt, Balkenius, Burgard, and Siegart, eds., 2001, pp. 107–113.
- [33] S. KALI AND P. DAYAN, *The involvement of recurrent connections in area ca3 in establishing the properties of place fields: a model*, Journal of Neuroscience, 20 (2000), pp. 7463–7477.
- [34] E. KANDEL, J. SCHWARTZ, AND T. JESSELL, *Principles of Neural Science*, McGraw-Hill, New York, NY, 2000.
- [35] G. LAKOFF AND R. NUNEZ, *Where Mathematics Comes From*, Basic Books, New York, 2000.
- [36] C. LANGTON, *Artificial life*, in Artificial Life: Proceedings of an Interdisciplinary Workshop on the Synthesis and Simulation of Living Systems, C. Langton, ed., Addison-Wesley, Reading, Massachusetts, 1989, pp. 1–49.
- [37] R. R. LLINAS, *i of the vortex*, The MIT Press, Cambridge, MA, 2001.
- [38] J. MCCLELLAND, B. MCNAUGHTON, AND R. O’REILLY, *Why there are complementary learning systems in the hippocampus and neocortex: Insights from the successes and failures of connectionist models of learning and memory*, Tech. Rep. PDP.CNS.94.1, Carnegie Mellon University, Mar. 1994.
- [39] T. MITCHELL, *Machine Learning*, WCB/McGraw-Hill, Boston, MA, 1997.
- [40] V. MOUNTCASTLE, *Perceptual Neuroscience: The Cerebral Cortex*, Harvard University Press, Cambridge, Massachusetts, 1998.
- [41] M. R. N. BURGESS AND J. O’KEEFE, *A model of hippocampal function*, Neural Networks, 7 (1994), pp. 1065–1083.
- [42] K. NAKAJIMA, M. A. MAIER, P. KIRKWOOD, AND R. LEMON, *Striking differences in transmission of corticospinal excitation to upper limb motoneurons in two primate species*, Journal of Neurophysiology, 84 (2000), pp. 698–709.
- [43] A. NEWELL AND H. SIMON, *Human Problem Solving*, Prentice Hall, Englewood Cliffs, N.J., 1972.
- [44] N. NILSSON, *Principles of Artificial Intelligence*, Tioga Publishers, Palo Alto, CA, 1980.

- [45] S. NOLFI AND D. FLOREANO, *Evolutionary Robotics: The Biology, Intelligence, and Technology of Self-Organizing Machines*, The MIT Press, Cambridge, MA, 2000.
- [46] R. C. O'REILLY AND Y. MUNAKATA, *Computational Explorations in Cognitive Neuroscience*, The MIT Press, Cambridge, Massachusetts, 2000.
- [47] H. PARTHASARATHY, J. SCHALL, AND A. GRAYBIEL, *Distributed but convergent ordering of corticostriatal projections: Analysis of the frontal eyefield and the supplementary eyefield in the macaque monkey.*, *Journal of Neuroscience*, 11 (1992), pp. 4468–4488.
- [48] E. RICH, *Artificial Intelligence*, McGraw-Hill Book Company, New York, NY, 1983.
- [49] O. R.O, *Biologically plausible error-driven learning using local activation differences: the generalized recirculation algorithm*, *Neural Computation*, 8 (1996), pp. 895–938.
- [50] E. ROLLS AND A. TREVES, *Neural Networks and Brain Function*, Oxford University Press, New York, 1998.
- [51] L. SQUIRE AND E. KANDEL, *Memory: From Mind to Molecules*, Henry Holt and Company, New York, 1999.
- [52] L. SQUIRE AND S. ZOLA, *Structure and function of declarative and nondeclarative memory systems*, *Genetic Programming and Evolvable Machines*, 93 (1996), pp. 13515 – 13522.
- [53] L. STEELS, *Intelligence with representation*, *Philosophical Transactions: Mathematical, Physical and Engineering Sciences*, 361 (2003), pp. 2381–2395.
- [54] G. F. STREIDTER, *Principles of Brain Evolution*, Sinauer Associates, Sunderland, Massachusetts, 2005.
- [55] P. L. STRICK, *Basal ganglia and cerebellar circuits with the cerebral cortex*, in *The Cognitive Neurosciences III*, M. S. Gazzaniga, ed., The MIT Press, Cambridge, MA, 2004, pp. 453–461.
- [56] SUGITA AND J. TANI, *Learning semantic combinatoriality from the interaction between linguistic and behavioral processes*, *Adaptive Behavior*, 13 (2005).
- [57] T. WINOGRAD, *Frame representations and the procedural-declarative controversy*, in *Representation and Understanding: Studies in Cognitive Science*, D. Bobrow and A. Collins, eds., Academic Press, 1975, pp. 185–210.
- [58] D. WOLPERT, R. C. MIALL, AND M. KAWATO, *Internal models in the cerebellum*, *Trends in Cognitive Sciences*, 2 (1998), pp. 338–347.
- [59] T. ZIEMKE, D.-A. JIRENHED, AND G. HESSLOW, *Internal simulation of perception. a minimal neurobotic model*, *Neurocomputing*, 68 (2005), pp. 85–104.