

Final Workshop Report

Future Challenges for the Science and Engineering of Learning
July 23-25, 2007

National Science Foundation

Organizers

Rodney Douglas - E. T. H. and University of Zurich

Terry Sejnowski - Salk Institute and University of California at San Diego

Executive Summary

On July 23, 2007, 20 researchers gathered at the National Science Foundation Headquarters in Arlington Virginia to participate in the NSF-sponsored workshop on "Future Challenges for the Science and Engineering of Learning" that ended on July 25, 2007. The goal of the workshop was to explore technological developments that are relevant to the NSF Science of Learning Centers (SLCs). Half of the researchers at the workshops were from SLCs and the other half were experts in neuromorphic engineering and machine learning.

This report summarizes the discussions that took place at the workshop but it does not capture the learning process that took place between the participants over its course. Issues were raised, strong opinions were voiced, and work took place among smaller groups that were charged with exploring focused topics. The most important outcome was a set of open questions that concern biological learning and machine learning; in particular, a common set of issues emerged that point toward a convergence that will benefit future research into learning in man and machine and their increasing interaction.

Acknowledgment

The organizers and workshop participants are grateful to the National Science Foundation for the extraordinary opportunity to explore learning across disciplines, and especially to Soo-Siang Lim, Program Director for the Science of Learning Centers, who suggested the topic.

Table of Contents

EXECUTIVE SUMMARY	1
PARTICIPANTS	3
GENERAL QUESTIONS OF LEARNING	5
What are the General Characteristics of Learning in Biology and Machines?	5
OPEN QUESTIONS IN LEARNING BY BIOLOGICAL SYSTEMS	7
What are the different temporal scales of learning and how are they implemented?	7
What are the practical implications of the Spacing Effect?	9
How does Wake-Sleep Cycle impact learning?	9
What are the factors underlying metaplasticity of synapses?	10
How can neuromodulation be used to improve memory?	10
How are different time scales represented in network space?	11
What are neuronal network architectures that support learning?	10
What are the factors underlying metaplasticity of synapses?	10
How can neuromodulation be used to improve memory?	10
How are different time scales represented in network space?	11
How do the local interactions of neurons support the extraction of more general statistical properties of signals?	12
What is the role of feedback in autonomous adaptive learning?	13
What is the role of the areal interactions, coherence, and synchrony in cortical learning?	13
How does the failure of prediction influence the search for solutions?	13
What properties are required for "autonomous" learning in a changing world?	13
How are the known different learning mechanisms combined in autonomous agents?	14
How do the interactions of social agents promote their learning?	14
OPEN QUESTIONS IN LEARNING BY ARTIFICIAL SYSTEMS	15
What are the major challenges in 'Machine' Learning?	16
What are effective shallow learning algorithms?	16
What are effective deep learning algorithms?	16
What is real-time autonomous learning?	16
Is it possible to develop neuromorphic electronic learning systems?	17
What are electronic signals and information representations suited for learning?	18
What are the elementary units of biological learning?	18
What is the relationship between physical self-assembly and learning?	19
How can the human interface to robots and other machines be enhanced?	19
How can efficient communication across human/machine interface be learned?	20
How can synergistic learning between humans and machines be enhanced?	20
How should portable human-machine learning systems be implemented physically?	21
APPENDIX	21
Schedule	21
Participants and Groups	22
Initial set of questions and starting points for discussions	23
The 'Willow Wish List'	25
Personal views offered by the participants	26

PARTICIPANTS

Andreas Andreou, Johns Hopkins University

*Tony Bell, University of California at Berkeley

Kwabena Boahen, Stanford University

Josh Bongard, University of Vermont

Rodney Douglas, E.T. H. /University of Zurich

Stefano Fusi, Columbia University

*Stephen Grossberg, Boston University

Giacomo Indiveri, E.T. H. /University of Zurich

*Ranu Jung, Arizona State University

*Pat Kuhl, University of Washington

Yann LeCun, New York University

Wolfgang Maass, Technical University Graz

*Andrew Meltzoff, University of Washington

*Javier Movellan, University of California at San Diego

Andrew Ng, Stanford University

*Howard Nusbaum, University of Chicago

*Terry Regier, University of Chicago

*Terry Sejnowski, Salk Institute/ University of California at San Diego

*Barbara Shinn-Cunningham, Boston University

Vladimir Vapnik, NEC Research

*Members of NSF Science of Learning Centers

Description of the Workshop

The NSF-sponsored workshop on "Future Challenges for the Science and Engineering of Learning" was held at the NSF Headquarters in Arlington VA from the evening of July 23 to the afternoon of July 25, 2007.

The goal of the workshop was to explore research opportunities in the broad domain of the Science and Engineering of Learning, and also to provide NSF with a Report identifying important open questions. It is anticipated that this Report will be used to encourage new research directions, particularly in the context of the NSF Science of Learning Centers (SLCs), and also to spur new technological developments.

The format of the meeting was designed to encourage open discussion. There were only relatively brief formal presentations. This format was very successful: So much so that the participants finally focused on a rather different set of questions than those that were originally proposed by the organizers.

The organizers circulated seed questions to the Participants prior to the Workshop. Initially these questions clustered into roughly 5 general areas: Science of Learning, Learning Theory, Learning Machines, Language Learning, and Teaching Robots. The Participants were assigned in pairs to one of 5 general areas, and each asked to give a 10 minute position statement in that domain. We hoped that their independent presentations would provide two separate attempts to identify, and motivate, some initial issues for discussion. These statements were not expected to offer solutions to the open problems, but merely to identify some of them, and to provoke and steer discussion about them. These position statements were planned for the mornings, and breakout discussions on the same topic were planned for the afternoons (see Schedule in Appendix). In practice though, discussions soon began to focus on the slightly different topics that are now described in this Report.

GENERAL QUESTIONS OF LEARNING

What are the General Characteristics of Learning in Biology and Machines?

Biological learners have the ability to learn autonomously, in an ever changing and uncertain world. This property includes the ability to generate their own supervision, select the most informative training samples, produce their own loss function, and evaluate their own performance. More importantly, it appears that biological learners can effectively produce appropriate internal representations for composable percepts -- a kind of organizational scaffold -- as part of the learning process.

By contrast, virtually all current approaches to machine learning (ML) typically require a human supervisor to design the learning architecture, select the training examples, design the form of the representation of the training examples, choose the learning algorithm, set the learning parameters, decide when to stop learning, and choose the way in which the performance of the learning algorithm is evaluated. This heavy-handed dependence on human supervision has greatly retarded the development and ubiquitous deployment autonomous artificial learning systems. This deficit also means that we have not yet understood the essence of learning, and so research on a variety of topics outlined below should be encouraged and supported by NSF.

One key theme is that of understanding the importance of the computational architecture of learning systems. It is widely believed that, unlike the shallow architectures typically employed in ML (machine learning), the learning in neural systems uses a deep, or heterarchical architecture, in which increasingly higher-level, or more abstract, concepts (or features) are composed of simpler ones. Re-representation across different neural mechanisms is an important feature of neural architecture. This deep architecture first extracts simple fundamental features that can be used to typify objects in a natural scene. These features then feed to later stages of processing, where more complex features are composed from these low-level building blocks. This hierarchical architecture is ubiquitous, and critical for enabling robust learning and characteristic of Hebb's notion of the nesting of levels of abstraction within neural structures. It is this proposed functional organization that would provide "scaffolding" by finding the elemental features that are most universal (across target learning domains) and fundamental within natural signals in the earliest sensory processes in the system. These simple features promote the ability to separate and segment a source of interest from the competing sources in the everyday environment. Once these features are tuned by experience or learned outright, higher-order regularities and structure can then be extracted with more refined processing.

Unlike supervised approaches to machine learning, most biological learning appears to be based on largely "unlabeled" data. Large amounts of natural data are available to the learning system without external supervision or labels, although the structure of the environment may provide causal or correlated signals that have distributional properties that may function akin to labeling. Generally, there is at most very weak external supervision, such as can be obtained by real-valued environmental

“reward” signals. Although a behaving organism rarely learns in a supervised manner (in a strict mathematical sense), the organism often appears to act with some implicit or explicit intent that presumably reflects an internal goal. The success or failure in achieving such an internally derived goal can provide a strong reward signal that drives learning. Even in situations where no specific outcome is desired, the organism may nevertheless be generating internal predictions relevant to possible actions and their consequences, and so can continually evaluate the success or failure of these predictions to improve learning.

How does a learner select relevant data? In most real-world situations there is no explicit teacher, and so a major problem is determining which of the deluge of impinging information are relevant for the solution of a particular task. In order to solve this problem, the learner must first be able to segment and segregate the sources of information in the environment, and then focus attention on the information that is relevant for a given task. The heterarchy of sensory processing probably supports these processes of segmenting and directing attention to construct the scaffold. However, the exact nature of this process and its biological implementation are essentially unknown.

There has been some progress in understanding the development of this scaffold in some important areas. A key example showing the evolution from general to specific learning is in language acquisition. For example, the learning *in utero* of simplest speech attributes (pitch, prosody, coarse phonotactics). Or the use of “Motherese” to accelerate segmentation and segregation by providing clear, over-articulated prosodically modulated speech that gives efficient information about what attributes and features matter in a particular language. Learning then progresses more quickly to enable understanding of words, semantics, and grammar. Understanding this principle would be beneficial for education in general.

During the past decade there has been an explosion of interest and success in ML. However, those successes are largely in the classification domain, operating on electronic data sets which are not typical of the ever-changing sensorium of the real world. The success of ML within its domain has had the negative effect of drawing that field away from the more challenging problems of biological / autonomous learners. NSF could encourage a return to the synergistic middle ground by promoting research at the interface between machine learning and biology, in addition to more basic research on the biological bases of learning that can guide future ML developments.

It is certainly realistic to start now a concentrated effort towards the design of autonomous learners. Bits and pieces of principles and methods have already been discovered that automate some of the tasks which a human supervisor formerly had to do. Examples include learning algorithms that are self-regulating in the sense that they automatically adapt the complexity of their hypotheses to the complexity of a learning task (e.g. SVMs, MLPs and ART systems); or automatically produce suitable internal representations of training examples (e.g. deep learning). A few examples of autonomous learning have been demonstrated in robotics. For example the winning robot in the DARPA Grand Challenge used a form of self-supervised learning called “near-to-far” learning, in which the general principle is to use the output of a reliable but specialized module (such as a stereo-based short range obstacle detection system, or a bumper) to provide supervisory signals (labels) for a trainable module with wider applicability (such as a long-range vision-based obstacle detector). Recent promising advances in deep learning have relied on unsupervised learning to create hierarchical representations.

OPEN QUESTIONS IN LEARNING BY BIOLOGICAL SYSTEMS

What are the different temporal scales of learning and how are they implemented?

Learning is not a monolithic process that can be described by any single mechanism. Moreover, learning may depend on a variety of mechanisms that are not themselves experience-dependent, but rather provide prior constraints or structure that is necessary for different aspects of learning to occur. All these mechanisms operate on extremely different time scales ranging from evolutionary time - during which the genetic foundations of learning mechanisms were established; to developmental time - during which the various functional periods in an organism's life play out; to task-processing time - during which information from a particular experience is encoded, processed, stored, and ultimately consolidated. The relationship among these time scales requires further research.

What is the functional relationship between brain evolution and behavioral success? This is a key question underlying the emergence of autonomous intelligence, for superior neurons will not survive Darwinian evolution if they cannot work together in circuits and systems to generate successful adaptive behaviors. In order to understand brain autonomy, one therefore needs to discover the computational level on which the brain determines indices of behavioral success. Decades of modeling research support the hypothesis that this success is measured at the network and system level rather than the single neuron level.

This evolutionary principle does not imply that individual neurons are irrelevant, but rather that the relevance of individual neurons lies in their ability to configure themselves in relation to their neighbors and thereby establish the network conditions that generate emergent network properties that map, in turn, onto successful behaviors. Thus, in order to understand brain autonomy, we need to explore principles and mechanisms that can unify the description of multiple organizational levels, and time scales of organization.

At the much shorter time scales, it is clear that some types of learning require repeated experiences. The timing of those experiences and the relative timing of information about the importance or meaning of those experiences (reinforcement or feedback) is critical to the process of learning and causal inference. However the mechanisms that mediate these effects are not completely understood. By contrast, other types of learning (such as the Garcia Effect in which food avoidance is rapidly learned) occur with a single experience that may be temporally remote from the ultimate outcome of the experience, and yet somehow an association between the antecedent and the consequent is formed.

Understanding how different brain mechanisms are able to interact across these widely different time scales, and how such time scales are linked to the nature of the differences among mechanisms, pose basic challenges for understanding biological learning and similar problems obtain for artificial learning: Are there universal principles of learning that transcend different time scales, and are there mechanisms that can account for the interactions that must hold across these different time scales? When mechanisms operate on different time scales, what kinds of representational differences emerge, and

how can these diverse mechanisms be coordinated? The operation of mechanisms at different time scales raises the problem of appropriate credit assignment bridging asynchronous processes.

A rather separate problem from the mechanisms of learning is the temporal structure of relevant events in the world from which we learn. The events that occur in a learning environment have their own set of time scales, from responding in real time to speech, to observing and understanding the unfolding of a set of physical behaviors in relation to an interactive tutor, or to the relative timing of feedback or reinforcement. How do the internal time scales of learning processes relate to this variability in timing of events in the world? In the following sections we consider a set of issues that highlight some of these questions in the context of biological mechanisms and psychological processes.

What are the practical implications of the Spacing Effect?

Over 100 years ago, Ebbinghaus reported that the timing of learning is critical to the strength of learning. The most effective learning occurs when practice is distributed over time such that learning experiences are separated in time. This *spacing effect* is remarkably robust in establishing long-term retention of learning. Perhaps most remarkable is that the spacing effect itself holds over a wide range of time scales, from spacing out learning within the course of a single day, to spacing out learning over days and even months (the maximum study interval is limited by forgetting and has been shown to be effective out to 14 months, with an 8 year test interval). The fact that this principle holds over such a wide range of time scales has argued against specific psychological theories of the spacing effect and presents a basic challenge to understanding the neural substrates of this effect. The fact that the spacing effect holds as well for memorizing a list of arbitrary words as for improving your tennis game, suggests some general feature of learning may be involved rather than a specific mechanism. What are the practical applications of this principle for the classroom and for the development of a skilled workforce?

How does Wake-Sleep Cycle impact learning?

Although learning is usually considered to be a process that takes place in an awake animal, there is a growing body of research indicating that sleep is an important part of the learning process. Animals spend an enormous amount of their lives asleep, but the biological and psychological processes of sleep are poorly understood. Sleep has been traditionally viewed as important for learning mainly because of the absence of experiences that could interfere with prior learning. It now appears that sleep is an active process that serves to consolidate learning that took place prior to sleeping. Theories of sleep have considered how the relative duration of sleep stages such as slow-wave sleep (SWS), or rapid-eye movement (REM) sleep are important to the consolidation process; however, understanding how the timing of these sleep stages occurs and changes over the course of a sleep period remains an important question.

One well established fact about sleep in mammals is that normal sleep duration is highest in infants and decreases with age, consistent with a parallel decrease in learning abilities. It is, however, clear that some aspects of consolidation following learning occur over different periods of waking time, and the relative roles of time-dependent consolidation and sleep-dependent consolidation need to be further investigated. Different types of learning (e.g., rote vs. generalization or simple vs. complex) may be consolidated through different types of mechanisms with different time courses. Although there has been speculation about possible mechanisms (synaptic downscaling, protein synthesis changes, etc.),

this is an important aspect of learning research that has been neglected. Understanding these mechanisms could lead to biological interventions (pharmacological) or behavioral interventions that could improve learning in parts of the work force (e.g., students at different ages, shift workers) and lead to improved learning algorithms that mimic aspects of the biological processes of learning.

What are the factors underlying metaplasticity of synapses?

There is accumulating evidence that synapses may be “metaplastic”; that is, the plasticity of a synapse may depend on the history of previous synaptic modifications. This means that the learning rule changes in time according to a meta-rule that reflects the interaction of mechanisms working on multiple timescales. What is the role of metaplasticity in learning, especially in non-stationary environment? Can metaplasticity explain the variable effectiveness of experiments that attempt to induce long-term synaptic modifications in vitro?

How can neuromodulation be used to improve memory?

Although computational models of learning typically focus on synaptic transmission and neural firing, there are significant modulatory influences in the nervous system that play important roles in learning. These influences are difficult to understand because the time scale of some of these modulatory mechanisms is much longer than the processing that occurs synaptically. For example, consideration of the time scale over which LTP/LTD develops suggests that this kind of mechanism may differ from the mechanism by which consolidation takes place. Modulatory neurotransmitters such as serotonin (5HT) can operate on a short time scale through ligand-gating (5HT₃ receptors) or a longer time scale through G-protein binding, and these time-scale differences may have different effects on neural processing. Hormones such as testosterone can have effects on memory consolidation that take as long as 24 hours to develop whereas other pharmacological effects on memory occur relatively quickly. These mechanisms are understood less well than mechanisms such as LTP, and the interaction of these modulatory systems with LTP needs to be investigated.

Increased understanding of the principles of neuromodulatory systems, and understanding how mechanisms interact across different time courses, could lead to the development of pharmacological aids to learning and memory that could speed up the process of initial learning or memory consolidation to improve retention. We learn faster and remember better when we are motivated. The state of the brain changes depending on the level of arousal and, in particular, with the expected reward. Neuromodulators such as dopamine and acetylcholine are involved in regulating what is attended, what is learned and how long it is remembered. How do these neuromodulatory systems regulate the time scales of learning? It has recently shown that slow learning in monkeys of the meaning of cues can be greatly speeded when the cues signal the amount and quality of reward. Moreover, rapid reversal of meaning and category learning occur much more slowly when the amount of reward is fixed.

How are different time scales represented in network space?

The interaction between different brain areas has been shown to be fundamental to understanding complex cognitive functions, flexible behavior and learning. For example the cortical-basal ganglia loop is known to be involved in reward prediction at different time scales. In particular the neurons predicting immediate rewards have been shown to be segregated from those predicting future rewards. Such a representation of time in neuronal network space might also help to perform a complex cognitive task in a changing environment where the rewarded visuomotor associations change in time.

What are the detailed neuronal network mechanisms that would implement such a reward prediction system operating on multiple timescales? In particular, how can neurons, which process signals on a millisecond range, participate in networks that can adaptively time behaviors with delays of many seconds? What is the role of these neurons in reinforcement learning? Can these neurons be used to generate a representation of the context, and even more generally, to abstract general rules? Such an issue has a fundamental importance because most reinforcement learning algorithms operate under the assumption that the states of the learning system are given, whereas in a real world the animal has to create autonomously the neural representations of these states.

What are neuronal network architectures that support learning?

For many decades the results of cortical neuroanatomy have had a mainly biological descriptive value. In recent times the goals of neuroanatomical research are becoming more focused on understanding the computational significance of the cortical neuronal circuits. Experimentalists should now be further encouraged to formulate their questions so as to resolve more abstract questions of learning in concrete anatomical and physiological terms. We need to understand whether and how brain structure is related to brain function. In particular, to what extent can we understand the function of a neural circuit from its anatomical structure? Indeed, to what extent can it be asserted that “the architecture is the algorithm”? An even broader question is how these different brain areas are integrated into an autonomous system. Can we exploit these principles for constructing artificial learning technologies?

The cerebellum, hippocampus, and neocortex, for example, each have different, but rather regular neuronal architectures. This regularity suggests that each region has a characteristic computational circuit, suited to the respective tasks that they implement. How are different learning competences embedded into these anatomically distinct architectures? Usually, learning is seen in terms of synaptic mechanisms, rather than as the operation of an entire learning/teaching subcircuit. It is possible that a significant fraction of (e.g.) cortical circuits is used to support learning/teaching of the local information processing circuit. How are these two processing and learning functions accomplished in real time?

Learning in the cerebral cortex takes place in highly recurrent internal and interareal circuits, with each part of the cortex processing different combinations of sensory and motor modalities and synthesizing novel perceptual, cognitive, affective and motor activities. What new computational principles do cortical circuits provide for the animal? How do these loopy circuits, whose processing is dominated by cortical rather than subcortical input, produce a consistent interpretation of the world rather than learn a hallucination. Indeed, how does normal cortical processing break down to lead to hallucinations and other signs of mental disorders? Is there active regulation of learning, so that only relevant patterns of activity are imprinted into the circuits? Is it possible that the cortical circuits have two basic sub-components; one performing processing, and the other exercising a teaching /learning role? In this way different neuronal architectures reflect in some large part the implementations of different learning algorithms. If architectures do reflect the learning algorithms, then one might ask of e.g. cortex, what are the different roles of the various laminae and/or neuron types in the implementation of learning.

The cerebellum and hippocampus are also laminated but have qualitatively different functions than the cerebral cortex. Even different areas of the cerebral cortex have different variations on the laminated 6-layer structure. Motor cortex does not have a layer 4. Primary visual cortex in primates has several sub-laminations in layer 4. Prefrontal cortex has more delayed period activity than found in the primary

sensory areas. What other specializations have evolved for other functions such as language, invariance (what) vs. location (where), and temporal (auditory) vs. spatial (visual)?

Although the cortex is highly structured, and a product of development, the cortex is, if not a *tabula rasa*, a structure with exceptional adaptability i.e. the cortical neuropil has a great deal of flexibility in the way that the connectivity and intrinsic properties of neurons can unfold during development while it receives inputs from the environment. To what extent does the cortex produce a circuit that is essentially free of bias, which is then configured by world data? Or is learning strongly constrained by the developing circuit organization, so that various kinds of appropriate data can be learned at different stages, as development unfolds? If the latter, then understanding the mechanisms of cortical self-construction and learning are inseparable. On a more general level, is it the case that learning is essentially intertwined with the self-organization of the multilevel structure of biological matter? Their stability depends on similar principles of coupling between scale levels. These need to be explored as general principles of multilevel organization, with particular implications for learning and adaptation, which are probably crucial processes in maintaining the integrity of organizations in the face of a non-stationary environment. These notions are also relevant for the development of artificial information processing systems. Often these systems are dominated by recurrence between many distributed modules, as found in biology. Understanding these principles in cortex has already begun to have benefit also in the development of artificial learning systems. Their systematic use may lead to solutions of many currently intractable problems in engineering and technology.

Within the cortical hierarchy, different levels have different functions. How is the cortical activity at these different levels used to control behavior? This is an issue in systems integration. How can multiple learning mechanisms (Hebbian, spike-timing dependent plasticity (STDP), homeostatic) and learning algorithms (supervised, unsupervised, reinforcement) be integrated? Many agree that the brain should be viewed as a “learning system“, but most existing models for learning in neural circuits/systems focus just on a single learning mechanism. Instead neural circuits and systems may be understood as support architectures for a variety of interacting learning algorithms. Understanding these neuronal network architectures and the interactions between learning algorithms is a new frontier in intelligent information processing and learning. They can provide new insights into how humans and other animals can be flexible in nonstationary environment through structures that are self-configuring. Learning clearly depends both on modification of synaptic connectivity (anatomy) and synaptic plasticity (physiology). Most models for learning have focused largely on the physiology of synapses rather than the anatomical growth processes. Those models that do combine both lead us to expect that large benefits for biology and technology can be anticipated by using learning rules that combine both these aspects.

How do the local interactions of neurons support the extraction of more general statistical properties of signals?

Behaving brains are exquisitely sensitive to environmental statistics and use principles of local computation that process huge amounts of non-stationary spatio-temporally distributed information. A key challenge is to understand how brains implicitly embody statistical constraints while learning incrementally in real time using only local computations in brain circuits.

What is the role of feedback in autonomous adaptive learning?

Anatomical studies make it clear that brain subsystems often interact closely with one another through bottom-up, horizontal, and top-down connections. A key issue for the future of biological and biomimetic computation is to understand whether and how these interactions support autonomous learning.

It is well-known that many *feedforward* learning systems are incapable of autonomous learning in a non-stationary world. In particular, they may exhibit catastrophic forgetting if learning is too fast or world statistics quickly change. In contrast, the brain can rapidly (often even with a single exposure) learn an important rare event for which there was no obvious prior. How does the brain avoid the problem of catastrophic forgetting in response to non-stationary data?

What is the role of the areal interactions, coherence, and synchrony in cortical learning?

Much behavioral and electrophysiological data, together with modeling studies, indicate that autonomous learning may use top-down feedback, notably learned top-down expectations that focus attention, to stabilize learning in response to a non-stationary world. In particular, bottom-up and top-down feedback exchanges generate dynamical states that can synchronize the neural activity of large brain regions that cooperate to represent coherent knowledge in the world, as they actively suppress signals that are not predictive in a given environment. Such dynamical states provide a way to understand how the brain copes with the vicissitudes of daily experience.

This raises the issue of what other useful properties are achieved by such dynamical states? In particular, how are these brain states linked to cognitive states? How can classical statistical learning approaches be modified or enhanced to incorporate coherent computation using dynamical states? How can such dynamical states be optimally represented in fast software and hardware systems in applications.

How does the failure of prediction influence the search for solutions?

The feedback processes in autonomous brain systems enable them to predict probable events. At the moment of a predictive failure, the correct answer is by definition unknown. How then, does an autonomous system use this predictive mismatch to drive a process of autonomous search, or hypothesis testing, to obtain a better answer? And how does this happen without an external teacher, and in a world filled with many distracters? How are spatial and object attention shifts controlled during autonomous search to discover predictive constraints that are hidden in world data? Much further work needs to be done to fully characterize how the brain regulates the balance between expected and unexpected events.

What properties are required for "autonomous" learning in a changing world?

In what ways are autonomous learners similar to, and differ from machine learning approaches? There are deep questions about the generality and optimality of learning principles. Are there clear optimal methods of learning that nature has selected and that we can discover? Or are there simply broad constraints common to the class of successful learners? Can artificial autonomous learners surpass biological systems?

Much experimental data and modeling indicates that the brain is capable of autonomously learning in real time in response to a changing world. In technology as well, many outstanding research problems concern how to achieve more autonomous control that can cope with unanticipated situations. Much

work in the design of increasingly autonomous mobile robots operating in increasingly unconstrained environments exhibit this trend, as exemplified by the DARPA autonomous vehicle Grand Challenges.

These parallel goals raise the following issues: First: How can we best understand how the brain achieves autonomous real-time learning in a changing world? That is, how can we design autonomous agents that learn in real time as they interact with a non-stationary environment that may not include any explicit teachers? Second: How can such autonomous agents flexibly combine environmental feedback that reflects the statistics of the world with feedback from other learners that may be more selective and goal-oriented?

How are the known different learning mechanisms combined in autonomous agents?

Much experimental and modeling evidence suggests that there are at least five functionally distinct types of learning:

- *Recognition learning* whereby we learn to recognize, plan, and predict objects and events in the world.
- *Reinforcement learning* whereby we evaluate and ascribe value to objects and actions in the world.
- *Adaptive timing* whereby we synchronize our expectations and actions to match world constraints.
- *Spatial learning* whereby we localize ourselves and navigate in the world.
- *Sensorimotor learning* whereby we carry out discrete actions in the world.

These may be called the *What*, *Why*, *When*, *Where*, and *How* learning processes, respectively. All of these types of learning interact. For example, visual object learning is a form of recognition learning that generates size-invariant and position-invariant representations of objects, whereas spatial learning generates representations of object position. Interactions across the *What* and *Where* cortical processing streams enable us to recognize valued objects and then reach towards them in space. Making intentional, goal-directed plans for the future potentially involves interactions among all five types of learning. How can these different types of learning systems be integrated into increasingly complete and autonomous system architectures, chips, and robots?

How do the interactions of social agents promote their learning?

Developmental science provides excellent examples of fast and powerful bidirectional learning. For example, parent-child interaction is prototypical case of teaching and learning effortlessly, efficiently, and adaptively for a changing world. For many tasks, it appears that humans may learn best from other social agents. However, the characteristics of ‘social agents’ are yet to be clearly defined. What are the cues or characteristics of social agents? If we could flesh out these characteristics we could develop learning/teaching technologies that captured those critical features and characteristics to enhance human education.

Studies of robotic (and educational software technologies) can also be used to explore the ‘social’ characteristics that support learning. For example, human interaction and learning from a robot may be enhanced if the robot:

- has human features
- reacts with correct timing (regardless of features)
- interactively adjusts to the needs of the learner
- have actions that seem to be goal-directed/planful.

What is the ‘social factor’ in learning? There are several hypotheses to be explored. Is the role of the teacher purely as a spotlight that highlights the thing to be learned? Attention and arousal work like this. Does the social nature/status of the teacher affect what and how things are learned? Do learners develop a portfolio of experts, learning different things from different people according to the different skills that those experts exhibit. Who do we learn from? Must they be trusted, and what does that mean? Do we learn preferentially from those with whom we share priors? Or from those agents that we can best simulate and understand? For example, learning slows down when there is a lack of understanding (and shared priors) between the participating agents. Does social/interactivity impact neural systems in particular ways (e.g. hormones; extension of the critical period in bird learning by social signals, shared neural circuits for coding actions of self and other)?

If the learner and teacher must be mutually adapted, then what are the cues used by learner to detect and select candidate teachers? We need also to study the strategies of a good teacher/tutor does: What cues do teachers use to modify their behavior in order to optimize child and student learning. Rapid and effective learning appears to depend on social scaffolding. How do we recognize and promote the assembly of that scaffold?

There are important advantages to be gained by harnessing a technology to promote learning on a long-time scale. These long-term relationships could be optimized by co-adaptation of the artifact and human partners. Endogenous neural compensatory learning on short and long time scales and physical constraints of interaction provide challenges to such synergistic learning.

Synergistic learning would be influenced by windows of opportunity that may be critical for induction of sustained learning. Ultimately, the technology becomes a training and educational tool that can be weaned away, after promoting synergistic learning in the biological system. Learning in the merged systems will have occurred when there are carry-over effects beyond the time-period when the technology is interacting with the biological systems. The synergistic learning platform could thus allow us to discover the principles governing activity dependent learning in living systems, to develop novel approaches to sense the dynamic changes in adaptive living system and the environment, and to deliver novel adaptive technology that encourages appropriate learning in biological systems.

OPEN QUESTIONS IN LEARNING BY ARTIFICIAL SYSTEMS

What are the major challenges in 'Machine' Learning?

The ability of machine-learning algorithms to use such forms of unsupervised reward signals holds the potential to increase the spectrum of problems they can solve. At the heart of all machine learning applications, from problems in object recognition to audio classification to navigation, is the issue of what data there is to learn from. Because unlabeled natural scene data is vastly easier to obtain than any labeled data, in the following, we focus on algorithms that are able to exploit unlabeled natural scene data. We pose the following questions:

What are effective shallow learning algorithms?

What shallow learning algorithms does the brain use? The engineering of machine learning has produced highly successful shallow algorithms for supervised learning, such as the support vector machine (SVM), logistic regression, generalized linear models, and many others, as well as “linear” unsupervised methods such as principal component analysis (PCA), independent component analysis (ICA), sparse coding, factor analysis, product of experts (PoEs), and restricted Boltzmann machines (RBMs). However, to date even the “shallow” levels of computation in the neocortex (for example, visual cortical area V1) are not fully understood. What learning principles result in the cortical organization of early processing, such as V1, A1, etc.? Can such algorithms be validated against what is known about V1, A1, etc. in the brain?

What are effective deep learning algorithms?

What deep learning algorithms does the brain use? To realize significant progress on developing human-like learning capabilities, it is of fundamental importance for us to develop effective learning algorithms for deep architectures. Indeed, many of the most difficult AI tasks of pressing national interest—such as computer vision object recognition—appear to be particularly suited to deep architectures. The visual pathway for object recognition in the visual cortex involves many layers of non-linear processing, most or all of which seem to support learning. How can effective deep learning algorithms be developed, and applied effectively to tasks such as visual object recognition? Many learning algorithms are defined in terms of explicit purely feedforward terms (such as ICA, PoE), purely feedback/generative terms (sparse coding, mixture of Gaussians), or a mix of both (RBM, stacked autoencoders, PCA); what role, if any, does feedback play in unsupervised and deep learning (ARTMAP)?

What is real-time autonomous learning?

As pointed out above, in most ML-applications, the intelligent aspects of learning are managed by the human supervisor (and not by the learning algorithm). Typically this human supervisor must:

- select the training examples
- choose the representation of the training examples
- choose the learning algorithm
- choose the learning rate
- decide when to stop learning
- choose the way in which the performance of the learning algorithm is evaluated.

This absolute requirement for a human expert supervisor precludes ubiquitous use of ML. It also means that ML has not yet captured the essence of learning.

By contrast, most current approaches to machine learning typically require the human supervisor to design the learning architecture, select the training examples, design the representation of the training examples, choose the learning algorithm, choose the learning parameters, decide when to stop learning, and choose the way in which the performance of the learning algorithm is evaluated. This strong dependence on human supervision has greatly held back the development and ubiquitous deployment artificial learning mechanisms. This also means that current ML algorithms have not yet encompassed the essence of learning, and so research on a variety of topics outlined below should be encouraged and supported.

We need to invest substantial efforts into the development of architectures and algorithms for autonomous learning systems. Most work on ML involves just a single learning algorithm, whereas architectures composed of several autonomously interacting learning algorithms and self-regulation mechanisms (each with a specific subtask) are needed. In particular, reinforcement learning systems need to interact with learning algorithms that optimize the classification of states required for particular tasks.

Is it possible to develop neuromorphic electronic learning systems?

Understanding the principles and the architecture of learning machines at the intersection of the disciplines of biology, physics and information is an exciting intellectual endeavor with enormous technological implications. The natural world is a tapestry of complex objects at different spatial and temporal scales that emerge as forces of nature transform and morph matter into animate and inanimate systems. Self-assembly at all scales is pervasive in nature where living systems of macroscopic organism dimensions organized in societal communities have evolved from hierarchical networks of nanoscale components i.e. molecules, into cells and tissue.

In the brain, optimization of functionality is heterarchically organized at all levels of the system *simultaneously*, in a seamless fashion. Understanding the physical constraints (cost functions) in the organization of the learning machinery will permit understanding of both brain function and theories that have the potential to bridge the physical scales from molecules to networks, individual behavior, and societies.

The synthesis of biologically-inspired synthetic structures that abstract the functional and developmental organization of the brain enables the rapid prototyping of machines where learning (acquiring knowledge) and functioning (using the knowledge to perform a specific task) is intricately intertwined. The computational complexity of modern learning algorithms and the engines that drive them are not

capable of this functionality today and even with advances in CMOS technological do not scale to large scale problems such as for example speaker independent large vocabulary speech recognition systems. Analog VLSI technology that fully embodies the style of learning in the mammalian neocortex should be based on a laminar architecture with at least six distinct and interacting processing layers (LAMINART).

What are electronic signals and information representations are suited for learning?

Research in several labs over the last decade has converged on data representation in synthetic neuromorphic systems called the Address Event Representation (AER). In AER, each ‘neuron’ on a sending device is assigned an address. When the neuron produces a spike its address is instantaneously put on an asynchronous digital bus. Event ‘collisions’ (cases in which sending nodes attempt to transmit their addresses at exactly the same time) are managed by on-chip arbitration schemes. AER allows for an encoding and processing that preserves the “analog” characteristics of real world stimuli, while at the same time allowing for the robust transmission of information over long distances using “spike” like stereotypical (digital) signals. Exploring learning algorithms/architectures in this representation allows for asynchronous machines that encode knowledge locally and globally as a sparse learnable network graph. The AER representation leads naturally to event-based, data-driven computational paradigms. This style of computation shares some properties with the one used by the nervous system, yet is still largely unexplored by the computer science and engineering community (consider for example the frame-based way that dominates computer vision). This representation is ideal for hardware implementations of neuromorphic system, and optimally suited for spike-based (be it from real or silicon neurons) learning mechanisms.

From a synthesis perspective, the engineering and computer science communities have developed architectural frameworks, CAD tools and efficient methods to design and manufacture microsystems at the “chip” level, moving up to the “board” level and down at the “micro” and “nano” levels. These require serial “pick and place” processes that are slow and expensive. As Complementary Metal Oxide Semiconductor (CMOS) VLSI technologies rapidly advances to deep sub-micron processes, the nanometer feature size is making the chasm between the micro/nanoscale device function and the macro scale system organization greater and greater. Developing tools and methodologies to accomplish this that mimic biology is crucial for further advances in the field; for example developing tools that automatically wire “neural like circuits” using algorithms based on the principles of gene networks. The broad research directions outlined above addresses fundamental questions at the interface of biological and physical systems as we strive to engineer new forms of complex informed matter. Our ultimate goal is the synthesis of networks at multiple physical scales in hybrid animate/inanimate technologies that can transduce, adapt, compute, learn and operate under closed loop. The outcome of this research effort impacts a diverse range of applications, from tissue engineering and rehabilitation medicine to biosensors for homeland defense.

What are the elementary units of biological learning?

Modern computers rely on the classical notion of a single “processor” or multiprocessor coupled to a memory hierarchy to process and maintain the states in the machine. Digital memories can be modified very rapidly and selectively, and with an arbitrarily large accuracy. These memories can then be preserved for arbitrarily long times, or at least until they are modified again. More obviously continuous valued physical systems such as neuromorphic electronic circuits, must rely on variables that are

encoded in some physical quantity like the charge across a capacitor. Such a quantity should be modifiable (plasticity) and it should be stable in time (memory preservation). Stability usually emerges from the interaction of the circuit elements that are responsible for implementing the memory element (for example a synapse). In such a system, the number of stable states is limited, and, as a consequence, memories have a dramatically short lifetime. Forgetting is not due to the passage of time, because each state is assumed to be inherently stable, but it is due to the overwriting of new memories. For example, when every memory element is bistable, every transition to a different state erases completely the memory of previous experiences. In order to improve the storage capacity, the memory devices should be smarter than a simple switch-like device, and the experience driven transitions from one stable state to another should depend on the previous history of modifications. This is called metaplasticity (see above), and it implies the transfer of much of the processing at the level of single memory elements. This may be one of the solutions adopted by biological synapses, for which the consolidation of a modification implies the activation of a cascade of biochemical processes working on multiple timescales. Metaplasticity can lead to a dramatic increase in memory performance, especially when the number of memory elements is large. What are the fundamental principles of metaplasticity, and how can they be used to leverage learning?

Unlike digital computers, brains often process distributed patterns of analog information. In many parts of the brain, such distributed patterns, rather than the activities of individual cells, are the units of information processing and learning. From this perspective, many properties of brain dynamics, such as the role of synchronous oscillations, become clearer. Thus one important goal of future research should be to understand how VLSI systems be designed for carrying out self-synchronizing processing of distributed patterns in laminar cortical circuits?

What is the relationship between physical self-assembly and learning?

In biology self-assembly and organization is dynamic. Many biological functions at the cellular and sub-cellular level are controlled by weak, non covalent interactions such as electrostatic, van der Waals forces, hydrogen bonds and metal coordination chemistry. Supramolecular chemistry is responsible for the intelligent function of animate matter, from the encoding of genetic information in basic amino-acids sequences at the sub-cellular level to the transport of ions and small molecules through cell membranes. Understanding how biological information processing systems employ *dynamic* matter and *learning* at all levels and time scales in networks of complex structures links the science of learning to emerging advances in materials science, chemistry and, specifically, nanotechnology.

How can the human interface to robots and other machines be enhanced?

Apprenticeship learning (also called imitation learning) has been applied with great success to a range of robotic and other artificial systems, ranging from autonomous cars and helicopters to intelligent text editors. Apprenticeship learning here refers either to situations where a separate (external) demonstration of a task is provided by a human to an artificial learning system—such as the human using her own hand to show a robot how to grasp an object. Alternatively teleoperation can be used to demonstrate directly through the robot that is attempting to learn to perform a task. For example, a human may demonstrate flying an aircraft, and the *same* aircraft may then try to learn to fly itself. Because the demonstration and the task to be learned took place using the same robotic hardware, this approach finesses the problem of having to find a mapping from the human's body parts/actions to the robot's body parts/actions. No doubt this will be a rapidly expanding field, as ever more complex machines must be taught more efficiently how to perform their required tasks.

There is already significant potential for cross-fertilization between development psychology and robotics, which have traditionally been two entirely separate research fields, even though both have converged to fairly similar classes of ideas in learning from teachers. Robotic apprenticeship learning today is extremely primitive compared to that studied in development psychology. Using insights from development psychology to develop robust apprenticeship learning methods holds the potential to revolutionize the capabilities of today's robots and computers. Similarly, insights from robotic apprenticeship learning—which has gained expertise over the past few decades about which classes of algorithms do and do not work on robots—will naturally further inform further developments in development psychology, and suggest novel theories and classes of experiments.

Some of the central questions and challenges facing apprenticeship learning are:

- Given a good demonstration of a task, what are effective inverse learning algorithms for inferring what goal the teacher was trying to attain? Similarly, given one or more noisy (or suboptimal) demonstrations of a task, what are effective inference algorithms for estimating the teacher's true goal?
- Many robots exist in exponentially large state spaces, which are infeasible to explore completely. How can demonstrations of a task be used to provide exploration information or to guide exploration?
- Given a multitude of demonstrations of many *different* tasks, what are effective strategies for retrieving the most appropriate piece of learned knowledge (or demonstration) when the robot faces a new, specific, task?
- Robots often reason about control tasks at different levels of abstraction (as in hierarchical control). How can demonstrations that are provided at one or more different levels of abstraction be combined and used effectively?
- If the robot observes an external demonstration of a task (i.e., if the demonstration was not via teleoperation), how can it find an appropriate mapping between the teacher's body parts/actions and the robot's own body parts/actions?
- What are the fundamental theoretical limits of apprenticeship learning, in terms of the number of demonstrations required, length of demonstration required, complexity of the task (and how do these interact with each other and prior learning)?
- What are effective principles for choosing how best to demonstrate tasks to an artificial or robotic system?

How can synergistic learning between humans and machines be enhanced?

A specific platform for investigating bidirectional (synergistic) learning is interaction between neural systems and intelligent machines. Two of the most important trends in recent technological developments are that:

- technology is increasingly integrated with biological systems
- technology is increasingly adaptive in its capabilities

The combination of these trends produces new situations in which biological systems and advanced technologies co-adapt. That is, each system continuously learns to interact with its environment in a manner directed at achieving its own objectives, yet those objectives may, or may not, coincide with those of its partner(s). The degree of 'success' in this learning process is thus highly dependent on the dynamic interaction of these organic and engineered adaptive systems. Optimizing the technology necessitates an approach that looks beyond the technology in isolation and looks beyond the technology as it interacts with the biological system in its current state. Here, the design of effective technology must consider its adaptive interaction with a biological system that is continuously learning.

Furthermore, often the objective of the technology is to shape or favorably influence the learning process. A set of scientific and technological challenges are emerging in the efforts to design engineered systems that guide and effectively utilize the complexity and elegance of biological adaptation.

The interaction between technological and biological systems could be improved by designing technological system to embody biological design principles wherever possible.

How can efficient communication across human/machine interface be learned?

A platform for addressing the future challenges of science and engineering of co-adaptive synergistic learning could be adaptive integration of technology with a person who has experienced traumatic injury that leads to neuromotor disability. Such a platform could be utilized to address fundamental issues regarding learning in biological systems, the design of adaptive engineered systems, and the dynamics of co-adaptive systems. The engineered system needs to access the patterns of activity of the nervous system. The patterns of activity of the biological system could be accessed using adaptive technology, soft and hardware that learns from a biological system that is nonstationary, dynamic, functions across multiple time and spatial scales and multiple modalities.

The adaptive technology that influences the biological system on short time-scales can be designed to be biomimetic, where the design of the control system is guided by the physical and programmatic constraints observed in biological systems, allows for real-time learning, stability, and error correction that accounts for the biological systems non-linearities and paucity of inputs to influence the biological system. Algorithms developed have to be adaptive, self-correcting and self-learning. Active learning on the part of the adaptive technology requires probing the living system in order to respond. Active teaching requires probing and interacting modifying the living system in some way.

How should portable human-machine learning systems be implemented physically?

To integrate synthetic learning technology with biological behaving systems, and especially people, requires learning machines that are commensurate with the constraints of human behavior. This means that learning machines must have the appropriate size, address issues of energy use, and be robust to environmental change. Donning and doffing of the learning machine as it interacts with the person will be of a paramount importance. In the event that the learning machine is implanted, additional constraints of material compatibility will have to be taken into account as will issues of communication across living and non-living matter. Similarly, ability to change architectural design of implanted systems are clear barriers and hence approaches that maximize functionality and perhaps included redundancy in design are necessary.