# Computational Neuroscience: From Biology to Cognition

Intermediate article

*Randall C O'Reilly*, University of Colorado, Boulder, Colorado, USA
*Yuko Munakata*, University of Denver, Colorado, USA

**CONTENTS**

*Computational neuroscience involves the construction of explicit computational models that implement neural mechanisms to simulate cognitive functions such as perception, learning and memory, motor function, and language.*

## INTRODUCTION

This article describes computer models that simulate the neural networks of the brain, with the goal of understanding how cognitive functions (perception, memory, thinking, language, etc.) arise from their neural basis. Many neural network models have been developed over the years, focused at many different levels of analysis, from engineering, to low-level biology, to cognition. Here, we consider models that try to bridge the gap between biology and cognition. Such models deal with real cognitive data, using mechanisms that are related to the underlying biology.

## THE RELATIONSHIP BETWEEN COGNITIVE AND NEURAL THEORIES

Computational models provide an important tool for linking data across multiple levels of analysis. The cognitive implications of cellular and network properties of neurons are often not immediately apparent: there are too many factors at many different levels interacting in complex ways. Trying to develop behavioral predictions that capture the complexity of the neural level can be like trying to predict the weather from a number of satellite measurements. A computational model, of the weather or of the brain, can help by formalizing information and relating it through complex, emergent dynamics. Cognitive properties can thus be understood as the product of a number of lower-level interactions, and neural properties can be understood in terms of their functional role in cognitive processes. Further, the effects of manipulations of lower-level interactions (e.g. through genetic knockouts or lesions) can be simulated and reconciled with the observed behavioral effects. Importantly, these simulations can make sense of much more subtle behavioral effects than the generic impairment of behavior on a cognitive task.

Although models thus have the potential to clarify brain–behavior relations, they do not always do so. Models can be underconstrained by neural and behavioral data, and thus be of questionable value in showing how the brain actually subserves behavior. Moreover, models may be devised merely as demonstrations that a behavior can be simulated, but this is insufficient for understanding why the models behave as they do. Thus, models must be evaluated in a balanced way for whether they advance understanding of specific phenomena, provide general principles, and make useful links between brain and behavior.

This article reviews a number of neuroscience-based computational models of various cognitive phenomena, with an emphasis on the general principles embodied by these models and their implications for understanding the general nature of cognition. Specifically, we examine models of: vision, including topography and receptive fields

in primary visual cortex and spatial attention emerging from interactions between parietal and temporal streams of processing; episodic memory subserved by the hippocampus; conditioning and skill learning subserved by the basal ganglia and cerebellum; working memory and cognitive control subserved by the prefrontal cortex; and language processing guided by neuropsychological cases. For a more comprehensive treatment of many of these models and the ideas behind them, see O'Reilly and Munakata (2000).

## COMPUTATIONAL MODELS OF VISION GUIDED BY NEUROSCIENCE

Vision is one of the best-studied domains in cognitive neuroscience, having a long tradition of integrating biological and psychophysical levels of analysis. Computational models of vision have been influential in both the vision and computational research communities. We review two areas of visual modeling here: topography and receptive fields in primary visual cortex; and spatial attention and the effects of parietal lobe damage. Other major areas of visual processing that have been modeled include object recognition, motion processing, and figure–ground segmentation.

## Topography and Receptive Fields in Primary Visual Cortex

The primary visual cortex (V1) provides an interesting target for computational models, because it has a complex but relatively well-understood organization of visual feature detectors (a 'representational structure') subject to considerable experience-based developmental plasticity (Hubel and Wiesel, 1962; Gilbert, 1996). Thus, the major question behind many of the V1 models has been: can we reproduce the complex representational structure of V1 through principled learning mechanisms exposed to realistic visual inputs?

First, we summarize the complex representational structure of V1. V1 neurons are generally described as edge detectors, where an edge is simply a roughly linear separation between regions of relative light and dark. These detectors differ in their orientation, size, position, and *polarity* (i.e. whether they detect transitions from light to dark or dark to light, or dark–light–dark or light–dark–light). The different types of edge detectors (together with other neurons that appear to encode visual surface properties) are packed into the two-dimensional sheet of the visual cortex according

to a particular topographic organization. The large-scale organization is a 'retinotopic map', which preserves the topography of the retinal image in the cortical sheet. At the smaller scale are 'hypercolumns' (see Figure 1), containing smoothly varying progressions of oriented edge detectors, among other things (Livingstone and Hubel, 1988). The hypercolumn also contains 'ocular dominance columns', in which V1 neurons respond preferentially to input from one eye or the other.

Many computational models have emphasized only one or a few aspects of the many detailed properties of V1 representations; for reviews, see Swindale (1996) and Erwin *et al*. (1995). For example, models have demonstrated how ocular dominance columns can develop based on a Hebbian learning mechanism, with greater local correlations in the neural firing coming from within one eye than from across eyes (Miller *et al*., 1989). Hebbian learning encodes correlational structure by strengthening the weights between neurons that fire together, and decreasing the weights between those that do not. (See Oja (1982) and Linsker (1988) for mathematical analyses of Hebbian correlational learning.)

Several models have demonstrated how a realistic set of oriented edge-detector representations can develop in networks presented with natural visual scenes, preprocessed in a manner consistent with the contrast-enhancement properties of the retina (e.g. Olshausen and Field, 1996; Bell and Sejnowski, 1997; van Hateren and van der Schaaff, 1997; O'Reilly and Munakata, 2000). The Olshausen and Field (1996) model demonstrated that sparse
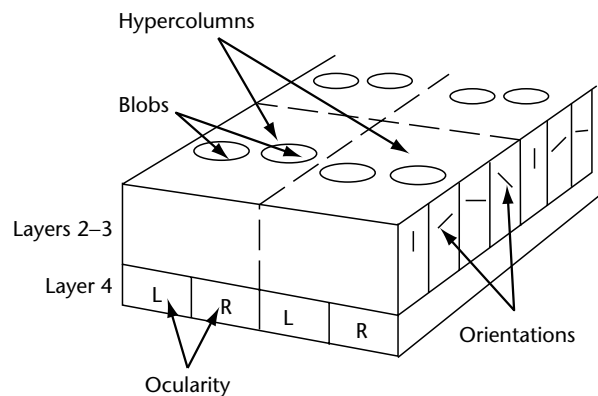


**Figure 1.** Structure of a cortical hypercolumn, representing a full range of orientations (in layers 2–3), ocular dominance columns (in layer 4, one for each eye), and surface features (in the blobs). Each such hypercolumn is focused within one region of retinal space, and neighboring hypercolumns represent neighboring regions.

representations (with relatively few active neurons) provide a useful basis for encoding real-world (visual) environments, but this model was not based on known biological principles. Subsequent work has shown how biologically-based models can develop oriented receptive fields, through a Hebbian learning mechanism with sparseness constraints in the form of inhibitory competition between neurons (a known property of cortex) (O'Reilly and Munakata, 2000). Furthermore, lateral excitatory connections within this network (another known property of cortex) produced a topographic organization consistent with several aspects of the hypercolumn structure (e.g. gradients of orientation, size, polarity, and phase tuning and pinwheel discontinuities: see Figure 2).

To summarize, these V1 models demonstrate how Hebbian learning mechanisms exposed to naturalistic stimuli, with certain kinds of biological prestructuring (e.g. connectivity patterns and inhibition), can produce aspects of the observed representational structure of V1. However, many complex aspects of early visual processing remain to be addressed, including motion, texture, and color sensitivity of different populations of V1 neurons.
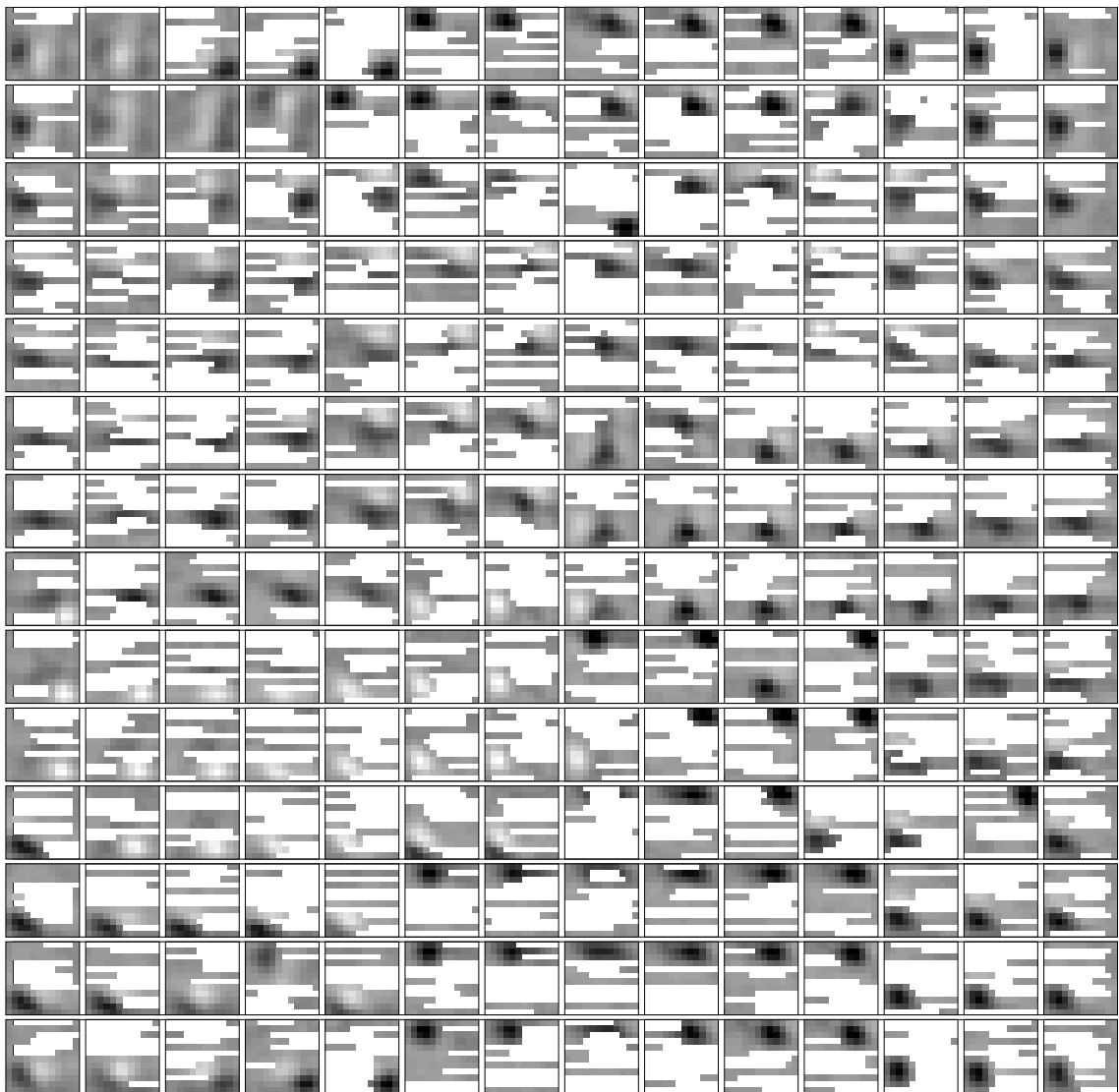


**Figure 2.** The receptive fields of model V1 neurons (O'Reilly and Munakata, 2000). Lighter shades indicate areas of on-center response, and darker shades indicate areas of off-center response. Individual units are shown by smaller grids (showing weights into those units from different locations in the retinally-organized input). These are organized into a larger grid representing the location of each unit within the simulated V1 hypercolumn.

## Spatial Attention and the Effects of Parietal Lobe Damage

Many computational models of higher-level vision have explored object recognition (e.g. Mozer, 1991; Fukushima, 1988; LeCun *et al.*, 1989) and spatial processing (e.g. Pouget and Sejnowski, 1997; Mozer and Sitton, 1998; Vecera and O'Reilly, 1998). Here we describe a model of spatial attention (Cohen *et al.*, 1994) that demonstrates how biologically-based computational models can provide alternative interpretations of cognitive phenomena. Spatial attention has traditionally been operationalized according to the Posner spatial cuing task (Posner *et al.*, 1984: see Figure 3). When attention is drawn, or cued, to one region of space, participants are faster to detect a target in that region (a validly cued trial) than a target elsewhere (an invalidly cued trial). Patients with damage to the parietal lobe have particular difficulty with invalidly cued trials.

According to the standard account of these data, spatial attention involves a 'disengage' module associated with the parietal lobe (Posner *et al.*, 1984). This module typically allows one to disengage from an attended location to attend elsewhere. This process of disengaging takes time; hence the slower detection of targets in unattended locations. Further, the disengage module is impaired with parietal damage, so that patients have difficulty disengaging from attention drawn to one side of the space.

Biologically-based computational models, based on recurrent excitatory connections and competitive inhibitory connections, provide an alternative explanation for these phenomena (Cohen *et al.*, 1994; O'Reilly and Munakata, 2000). In this framework, the facilitative effects of drawing attention to one region of space result from excitatory connections between spatial and other representations of that region: this excitatory support makes it easier to process information in that region. The slowing
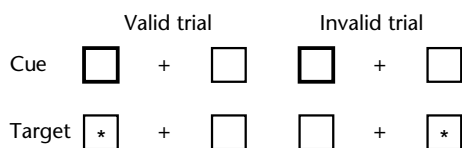
observed in the invalid trials results from inhibitory competition between different spatial regions. Under this model, damage to the parietal lobe simply impairs the ability of the corresponding region in space to have sufficient excitatory support to compete effectively with other regions.

The two models make distinct predictions (Cohen *et al.*, 1994; O'Reilly and Munakata, 2000). For example, following bilateral parietal damage, the disengage model predicts disengage deficits on both sides of space (Posner *et al.*, 1984), but the competitive inhibition model predicts reduced attentional effects (smaller valid and invalid trial effects). Data support the latter model (Coslett and Saffran, 1991; Verfaellie *et al.*, 1990), demonstrating the utility of biologically-based computational models for alternative theories of cognitive phenomena.

## COMPUTATIONAL MODELS OF EPISODIC MEMORY AND THE HIPPOCAMPUS

Damage to the hippocampus, in the medial temporal lobe, can produce severe memory deficits, while leaving unimpaired certain kinds of learning and memory (Scoville and Milner, 1957; Squire, 1992). Many computational models have been developed to explore the exact contribution of the hippocampus, and these models have had a major influence (e.g., Marr, 1971; Treves and Rolls, 1994; Hasselmo and Wyble, 1997; Moll and Miikkulainen, 1997; Alvarez and Squire, 1994; Levy, 1989; Burgess *et al.*, 1994; Samsonovich and McNaughton, 1997).

One framework has combined known biological features of the hippocampal formation with computationally motivated principles about learning and memory to further clarify the unique contributions of the hippocampus in memory (McClelland *et al.*, 1995; O'Reilly and Rudy, 2000, 2001; O'Reilly and McClelland, 1994; O'Reilly *et al.*, 1998). The central idea is that there are two basic types of learning that an organism must engage in – learning about specifics and learning about generalities – and that because the computational mechanisms for achieving these types of learning are in direct conflict, the brain has evolved two separate brain structures to achieve them. The hippocampus appears to be specialized for learning about specifics, while the neocortex is good at extracting generalities.

Learning about specifics requires keeping representations separated (to avoid interference), whereas learning about generalities requires overlapping representations that encode shared



**Figure 3.** The Posner spatial attention task. The cue is a brightening or highlighting of one of the boxes, which focuses attention on that region of space. Reaction times to detect the target are faster when this cue is valid (the target appears in that region) than when it is invalid (the target appears elsewhere).

structure across many different experiences. Furthermore, learning about generalities requires a slow learning rate to gradually integrate new information with existing knowledge, while learning about specifics can occur rapidly. This rapid learning is particularly important for episodic memory, where the goal is to encode the details of specific events as they unfold.

These computational principles provide a satisfying and precise characterization of the division of labor between the hippocampus and the neocortex. The models that implement these principles have been shown to account for a wide range of specific learning and memory findings, including nonlinear discrimination, incidental conjunctive encoding, fear conditioning, and transitive inference in rats (O'Reilly and Rudy, 2001) and human recognition memory (O'Reilly et al., 1998). However, these models fail to incorporate important aspects of the hippocampal formation (e.g. the subiculum and the mossy cells in the hilus), and many more complex behaviors that depend on the hippocampus (and its interactions with other brain areas) remain to be addressed.

## COMPUTATIONAL MODELS OF CONDITIONING AND SKILL LEARNING IN THE BASAL GANGLIA AND CEREBELLUM

A convergence between biological, behavioral and computational approaches has been achieved in the domain of conditioning (learning to associate stimuli and actions with rewards). In the computational domain, reinforcement learning mechanisms can change the behavior of a simulated animal according to reward contingencies in the environment (Sutton and Barto, 1998). Such learning mechanisms, including the 'temporal differences' algorithm (Sutton, 1988), not only work well mathematically (e.g. Dayan, 1992), but also correspond with aspects of neural recordings made in the reward-processing area of the brain (Montague et al., 1996; Schultz et al., 1997).

Specifically, a straightforward neural implementation of the temporal differences algorithm involves a systematic transition of reward-related neural firing similar to that observed in dopamine neurons in the midbrain. During a simple conditioning task where a sensory stimulus (e.g. a tone) reliably predicts a subsequent reward (e.g. orange juice), these neurons initially fire in response to the reward, but then after some trials of learning they respond to the sensory stimulus that predicts the

reward and no longer to the reward itself (Schultz et al., 1993; Schultz et al., 1995). This transfer of reward-related firing from the actual reward to predictors of the reward is an essential property of the temporal differences mechanism as implemented by Montague et al. (1996), which thus provides a principled, provably effective explanation for why the brain appears to learn in this manner.

Models of motor performance and skill learning have been developed based on the biological properties of the relevant underlying brain areas, including the basal ganglia (which includes the striatum, globus pallidus, substantia nigra, subthalamic nucleus, and nucleus accumbens) and the cerebellum (e.g. Beiser et al., 1997; Wickens, 1997; Houk et al., 1995; Berns and Sejnowski, 1996; Schweighofer et al., 1998a, b; Contreras-Vidal et al., 1997). These models accord well with detailed neural properties of these areas, but tend to focus on simpler aspects of motor performance: complex motor skills remain to be addressed.

## COMPUTATIONAL MODELS OF WORKING MEMORY, COGNITIVE CONTROL, AND PREFRONTAL CORTEX

The prefrontal cortex is important for a range of cognitive functions, which can be described generally as higher level cognition, in that they go beyond basic perceptual, motor, and memory functions. For example, frontal cortex has been implicated in problem-solving tasks like the Tower of Hanoi (e.g. Shallice, 1982; Baker et al., 1996; Goel and Grafman, 1995), which requires executing a sequence of moves to achieve a subsequent goal. Many theoretical perspectives summarize the function of frontal cortex in terms of 'executive control', 'controlled processing', or a 'central executive' (e.g. Baddeley, 1986; Shallice, 1982; Gathercole, 1994; Shiffrin and Schneider, 1977), without explaining at a mechanistic level how such functionality could be achieved. Computational models provide an important tool for exploring specific mechanisms that might achieve 'executive-like' functionality.

### Working Memory and Active Maintenance

One proposal is that the fundamental mechanism underlying frontal function is 'active maintenance', which then enables all the other 'executive' functionality ascribed to the frontal cortex (Cohen et al., 1996; Goldman-Rakic, 1987; Munakata, 1998;

O'Reilly *et al.*, 1999; O'Reilly and Munakata, 2000; Roberts and Pennington, 1996). For example, a flexible, adaptive, active maintenance system can meet information processing challenges by using the strategic activation and deactivation of representations (activation-based processing) instead of weight changes (weight-based processing) (O'Reilly and Munakata, 2000). There are trade-offs between these types of processing (e.g. activations can be more rapidly switched than weights, but they are also transient); so using both kinds of processing is better than using either alone.

There is considerable direct biological evidence that the frontal cortex subserves the active maintenance of information over time, as encoded in the persistent firing of frontal neurons (e.g. Fuster, 1989; Goldman-Rakic, 1987; Miller *et al.*, 1996). Many computational models of this basic active maintenance function have been developed (Braver *et al.*, 1995; Dehaene and Changeux, 1989; Zipser *et al.*, 1993; Seung, 1998; Durstewitz *et al.*, 2000; Camperi and Wang, 1997). Some models have further demonstrated that active maintenance can account for frontal involvement in a range of different tasks that might otherwise appear to have nothing to do with maintaining information over time.

## Inhibition, Flexibility, and Perseveration

For example, several models have demonstrated that frontal contributions to 'inhibitory' tasks can be explained in terms of active maintenance instead of an explicit inhibitory function. Actively maintained representations can support (via bidirectional excitatory connectivity) correct choices, which will therefore indirectly inhibit incorrect ones via standard lateral inhibition mechanisms within the cortex. A model of the Stroop task provided an early demonstration of this point (Cohen *et al.*, 1990). In this task, color words (e.g. 'red') are presented in different colors, and people are instructed to either read the word or name the color in which the word is written. In the conflict condition, the ink color and the word are different. Because we have so much experience of reading, we naturally tend to read the word, even if instructed to name the color, so that responses are slower and more error-prone in the color-naming conflict condition than in the word-reading one. These color-naming problems are selectively magnified with frontal damage. This frontal deficit has usually been interpreted in terms of the frontal cortex helping to inhibit the dominant word-reading pathway. However, Cohen *et al.* (1990) showed that they could account for both normal and frontal-damage

data by assuming that the frontal cortex instead supports the color-naming pathway, which then collaterally inhibits the word-reading pathway. Similar models have demonstrated that, in infants, the ability to inhibit perseverative reaching (searching for a hidden toy at a previous hiding location rather than at its current location) can develop simply through increasing ability to actively maintain a representation of the correct hiding location (Dehaene and Changeux, 1989; Munakata, 1998). Again, such findings challenge the standard interpretation that inhibitory abilities *per se* must develop for improved performance on this task (Diamond, 1991).

The activation-based processing model of frontal function can also explain why frontal cortex facilitates rapid switching between different categorization rules in the Wisconsin card sorting task and related tasks. In these tasks, subjects learn to categorize stimuli according to one rule via feedback from the experimenter, and then the rule is changed. With frontal damage, patients tend to perseverate in using the previous rule. A computational model of a related intradimensional/extradimensional (ID/ED) categorization task demonstrated that the ability to rapidly update active memories in frontal cortex can account for detailed patterns of data in monkeys with frontal damage (O'Reilly *et al.*, 2002; O'Reilly and Munakata, 2000).

Computational models of frontal function can provide mechanistic explanations that unify the various roles of the frontal cortex, from working memory to cognitive control and planning and problem-solving. However, it remains to be shown whether complex 'intelligent' behavior can be captured using these basic mechanisms.

## COMPUTATIONAL MODELS OF LANGUAGE USE GUIDED BY NEUROPSYCHOLOGICAL CASES

Damage to language-related brain areas causes a wide variety of impairments. One class of such impairments, the dyslexias (also known as alexias), have been the subject of a series of influential computational models of the normal and impaired reading process (Seidenberg and McClelland, 1989; Plaut and Shallice, 1993; Plaut *et al.*, 1996). These models simulate the pathways between visual word inputs (orthography), word semantics, and verbal word outputs (phonology), and can account for different kinds of dyslexias in terms of differential patterns of damage to these pathways.

These models have been influential in part because they suggest an alternative, somewhat

counterintuitive, interpretation of how words are represented and how language processing works. Traditional models have assumed that the brain contains a 'lexicon', with distinct representations for different words. Furthermore, these models assume that reading a word aloud (i.e., mapping from orthography to phonology) can occur via two different routes: pronunciation rules (for 'regular' words like 'make'), or a mechanism like a lookup table (for 'exception' words like 'yacht') (Pinker, 1991; Coltheart *et al.*, 1993; Coltheart and Rastle, 1994). In contrast to these dual-route models, the neural network models posit a single pathway to process both regular and exception words, and they employ a distributed lexicon without centralized, discrete lexical representations. Lexical processing occurs in pathways that map between different aspects of word representations (see Figure 4).

In general, neural networks can learn all kinds of different mappings: fully regular ones, like the spelling-to-sound mapping of the 'a' in words like 'make' and 'bake', as well as irregular mappings that occur in words like 'yacht'. Nevertheless, networks are sensitive to both the degree of regularity and the frequencies of different mappings. Specifically, neural network models predict frequency-by-regularity interactions that would not be expected in dual-route models, and which are observed in behavioral tests (Plaut *et al.*, 1996). Furthermore, these network models can account for patterns of deficit with brain damage that would seem improbable under dual-route models. For example, people with surface dyslexia can read

non-words (e.g. 'nust'), but they are impaired at retrieving semantic information from written words, and have difficulty in reading exception words. It would be natural in neural network models to interpret this as damage in the pathway between orthography and semantics. Interestingly, surface dyslexics' difficulty with exception words is generally limited to low-frequency exceptions (e.g. 'yacht'); they do not have difficulty reading high-frequency exceptions (e.g. 'are'). This pattern suggests that the remaining 'direct' pathway between orthography and phonology can handle both regular words and high-frequency exceptions, as in the network models. This pattern of data is not easily explained by the dual-route models: with two pathways, either regular words or exceptions should be affected, but not both, and independently of frequency.

Neural network models of language can provide alternative, counterintuitive ways of explaining some of the complex patterns of deficits that occur with brain damage. Nevertheless, such models remain controversial: neural network accounts are challenged by revised versions of dual-route models, and by the complexity of different neuropsychological profiles associated with damage to different language areas.

## CONCLUSION

Computational models based on the neural networks of the brain can provide important insights. Many models have applied a set of basic principles to a range of phenomena, and arrived at explanations completely different from those based on purely verbal cognitive theories. Hence, these models have played an important role in guiding empirical research and theorizing in a number of domains.

Despite these successes, many researchers remain skeptical of models. A common concern is that different models may employ different sets of mechanisms to explain the same data, so that it may not be very significant that a given model can simulate a set of data. Several points have been made in response to this concern.

Firstly, it applies not only to computational models, but to scientific theorizing in general (several theories can account for the same data). Competing theories and models can be evaluated by many other criteria than simply accounting for a set of data, such as the accuracy of their predictions, the coherence of their theoretical framework, and the ease of accounting for new data (Munakata and Stedron, 2002).
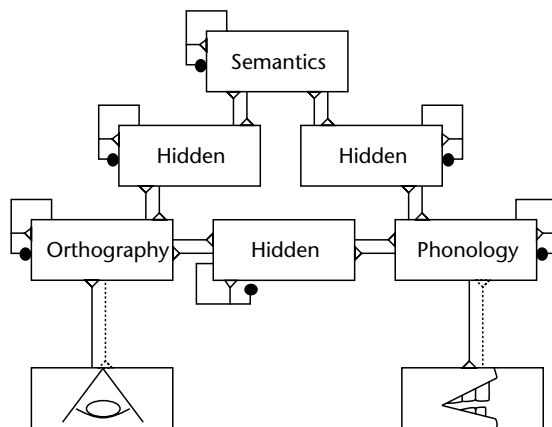


**Figure 4.** A neural network model of reading aloud. Words are represented in a distributed fashion across orthographic (visual word recognition), phonological (speech output), and semantic areas.

Secondly, mechanisms developed independently can turn out to be equivalent (e.g. O'Reilly, 1996), providing converging evidence for their utility, and indicating more coherence to principles than might otherwise be evident.

Thirdly, a common set of mechanisms appears to be emerging as the field continues to mature. For example, over 40 different phenomena (including most of what has been described above) have been modeled using a common set of mechanisms (O'Reilly and Munakata, 2000). This set of mechanisms was developed over many years by many different researchers, and has now been consolidated and integrated into one coherent framework (O'Reilly, 1998).

Therefore, there is a largely consistent set of ideas underlying many neural network models, and this framework provides an important way of understanding the connections between cognition and underlying neural systems.

## References

Alvarez P and Squire LR (1994) Memory consolidation and the medial temporal lobe: a simple network model. *Proceedings of the National Academy of Sciences* **91**: 7041–7045.

Baddeley AD (1986) *Working Memory*. New York, NY: Oxford University Press.

Baker SC, Rogers RD, Owen AM *et al.* (1996) Neural systems engaged by planning: a PET study of the Tower of London task. *Neuropsychologia* **34**: 515–526.

Beiser DG, Hua SE and Houk JC (1997) Network models of the basal ganglia. *Current Opinion in Neurobiology* **7**: 185–190.

Bell AJ and Sejnowski TJ (1997) The independent components of natural images are edge filters. *Vision Research* **37**: 3327–3338.

Berns GS and Sejnowski TJ (1996) How the basal ganglia make decisions. In: Damasio A, Damasio H and Christen Y (eds) *Neurobiology of Decision-Making*. Berlin, Germany: Springer-Verlag.

Braver TS, Cohen JD and Servan-Schreiber D (1995) A computational model of prefrontal cortex function. In: Touretzky DS, Tesauro G and Leen TK (eds) *Advances in Neural Information Processing Systems*, pp. 141–148. Cambridge, MA: MIT Press.

Burgess N, Recce M and O'Keefe J (1994) A model of hippocampal function. *Neural Networks* **7**: 1065–1083.

Camperi M and Wang XJ (1997) Modeling delay-period activity in the prefrontal cortex during working memory tasks. In: Bower J (ed.) *Computational Neuroscience*, chap. XLIV, pp. 273–279. New York, NY: Plenum Press.

Cohen JD, Dunbar K and McClelland JL (1990) On the control of automatic processes: a parallel distributed processing model of the Stroop effect. *Psychological Review* **97**: 332–361.

Cohen JD, Romero RD, Farah MJ and Servan-Schreiber D (1994) Mechanisms of spatial attention: the relation of macrostructure to microstructure in parietal neglect. *Journal of Cognitive Neuroscience* **6**: 377–387.

Cohen JD, Braver TS and O'Reilly RC (1996) A computational approach to prefrontal cortex, cognitive control, and schizophrenia: recent developments and current challenges. *Philosophical Transactions of the Royal Society, Series B* **351**: 1515–1527.

Coltheart M and Rastle K (1994) Serial processing in reading aloud: evidence for dual-route models of reading. *Journal of Experimental Psychology: Human Perception and Performance* **20**: 1197–1211.

Coltheart M, Curtis B, Atkins P and Haller M (1993) Models of reading aloud: dual route and parallel-distributed-processing approaches. *Psychological Review* **100**: 589–608.

Contreras-Vidal JL, Grossberg S and Bullock D (1997) A neural model of cerebellar learning for arm movement control: cortico-spino-cerebellar dynamics. *Learning and Memory* **3**: 475–502.

Coslett HB and Saffran E (1991) Simultanagnosia. To see but not two see. *Brain* **114**: 1523–1545.

Dayan P (1992) The convergenece of TD($\lambda$) for general $\lambda$. *Machine Learning* **8**: 341–362.

Dehaene S and Changeux JP (1989) A simple model of prefrontal cortex function in delayed-response tasks. *Journal of Cognitive Neuroscience* **1**: 244–261.

Diamond A (1991) Neuropsychological insights into the meaning of object concept development. In: Carey S Gelman R (eds) *The Epigenesis of Mind*, chap. III, pp. 67–110. Mahwah, NJ: Lawrence Erlbaum.

Durstewitz D, Seamans JK and Sejnowski TJ (2000) Dopamine-mediated stabilization of delay-period activity in a network model of prefrontal cortex. *Journal of Neurophysiology* **83**: 1733–1750.

Erwin E, Obermayer K and Schulten K (1995) Models of orientation and ocular dominance columns in the visual cortex: a critical comparison. *Neural Computation* **7**: 425–468.

Fukushima K (1988) Neocognitron: a hierarchical neural network capable of visual pattern recognition. *Neural Networks*: **1**: 119–130.

Fuster JM (1989) *The Prefrontal Cortex: Anatomy, Physiology and Neuropsychology of the Frontal Lobe*. New York, NY: Raven Press.

Gathercole SE (1994) Neuropsychology and working memory: a review. *Neuropsychology* **8**: 494–505.

Gilbert CD (1996) Plasticity in visual perception and physiology. *Current Opinion in Neurobiology* **6**: 269–274.

Goel V and Grafman J (1995) Are the frontal lobes implicated in 'planning' functions? Interpreting data from the Tower of Hanoi. *Neuropsychologia* **33**: 623–642.

Goldman-Rakic PS (1987) Circuitry of primate prefrontal cortex and regulation of behavior by representational memory. In: Brookhart JM and Mountcastle VB (eds) *Handbook of Physiology. The Nervous System*, vol. V, pp. 373–417. Baltimore, MD: American Physiological Society.

Hasselmo ME and Wyble B (1997) Free recall and recognition in a network model of the hippocampus: simulating effects of scopolamine on human memory function. *Behavioural Brain Research* **67**: 1–27.

van Hateren JH and van der Schaaff A (1997) Independent component filters of natural images compared with simple cells in primary visual cortex. *Proceedings of the Royal Society, Series B* **265**: 359–366.

Houk JC, Davis JL and Beiser DG (eds) (1995) *Models of Information Processing in the Basal Ganglia*. Cambridge, MA: MIT Press.

Hubel D and Wiesel TN (1962) Receptive fields, binocular interaction, and functional architecture in the cat's visual cortex. *Journal of Physiology* **160**: 106–154.

LeCun Y, Boser B, Denker JS *et al.* (1989) Backpropagation applied to handwritten zip code recognition. *Neural Computation* **1**: 541–551.

Levy WB (1989) A computational approach to hippocampal function. In: Hawkins RD and Bower GH (eds) *Computational Models of Learning in Simple Neural Systems*, pp. 243–304. San Diego, CA: Academic Press.

Linsker R (1988) Self-organization in a perceptual network. *Computer* **21**(3): 105–117.

Livingstone M and Hubel D (1988) Segregation of form, color, movement, and depth: anatomy, physiology, and perception. *Science* **240**: 740–749.

Marr D (1971) Simple memory: a theory for archicortex. *Philosophical Transactions of the Royal Society, Series B* **262**: 23–81.

McClelland JL, McNaughton BL and O'Reilly RC (1995) Why there are complementary learning systems in the hippocampus and neocortex: insights from the successes and failures of connectionst models of learning and memory. *Psychological Review* **102**: 419–457.

Miller EK, Erickson CA and Desimone R (1996) Neural mechanisms of visual working memory in prefontal cortex of the macaque. *Journal of Neuroscience* **16**: 5154–5167.

Miller KD, Keller JB and Stryker MP (1989) Ocular dominance column development: analysis and simulation. *Science* **245**: 605–615.

Moll M and Miikkulainen R (1997) Convergence-zone episodic memory: analysis and simulations. *Neural Networks* **10**: 1017–1036.

Montague PR, Dayan P and Sejnowski TJ (1996) A framework for mesencephalic dopamine systems based on predictive Hebbian learning. *Journal of Neuroscience* **16**: 1936–1947.

Mozer MC (1991) *The Perception of Multiple Objects: A Connectionist Approach*. Cambridge, MA: MIT Press.

Mozer MC and Sitton M (1998) Computational modeling of spatial attention. In: Pashler H (ed.) *Attention*, pp. 341–393. London, UK: UCL Press.

Munakata Y (1998) Infant perseveration and implications for object permanence theories: a PDP model of the A-not-B task. *Developmental Science* **1**: 161–184.

Munakata Y and Stedron JM (forthcoming). Memory for hidden objects in early infancy. In: Fagen J and Hayne H (eds) *Advances in Infancy Research*, vol. XIV. Norwood, NJ: Ablex.

Oja E (1982) A simplified neuron model as a principal component analyzer. *Journal of Mathematical Biology* **15**: 267–273.

Olshausen BA and Field DJ (1996) Emergence of simple-cell receptive field properties by learning a sparse code for natural images. *Nature* **381**: 607–609.

O'Reilly RC (1996) Biologically plausible error-driven learning using local activation differences: the generalized recirculation algorithm. *Neural Computation* **8**: 895–938.

O'Reilly RC (1998) Six principles for biologically-based computational models of cortical cognition. *Trends in Cognitive Sciences* **2**: 455–462.

O'Reilly RC and McClelland JL (1994) Hippocampal conjunctive encoding, storage, and recall: avoiding a tradeoff. *Hippocampus* **4**: 661–682.

O'Reilly RC and Munakata Y (2000) *Computational Explorations in Cognitive Neuroscience: Understanding the Mind by Simulating the Brain*. Cambridge, MA: MIT Press.

O'Reilly RC and Rudy JW (2000) Computational principles of learning in the neocortex and hippocampus. *Hippocampus* **10**: 389–397.

O'Reilly RC and Rudy JW (2001) Conjunctive representations in learning and memory: principles of cortical and hippocampal function. *Psychological Review* **108**: 311–345.

O'Reilly RC, Norman KA and McClelland JL (1998) A hippocampal model of recognition memory. In: Jordan MI, Kearns MJ and Solla SA (eds) *Advances in Neural Information Processing Systems*, vol. X, pp. 73–79. Cambridge, MA: MIT Press.

O'Reilly RC, Braver TS and Cohen JD (1999) A biologically based computational model of working memory. In: Miyake A and Shah P (eds) *Models of Working Memory: Mechanisms of Active Maintenance and Executive Control*, pp. 375–411. New York, NY: Cambridge University Press.

O'Reilly RC, Noelle D, Braver TS and Cohen JD (2002) Prefrontal cortex and dynamic categorization tasks: representational organization and neuromodulatory control. *Cerebral Cortex* **12**: 246–257.

Pinker S (1991) Rules of language. *Science* **253**: 530–535.

Plaut DC and Shallice T (1993) Deep dyslexia: a case study of connectionist neuropsychology. *Cognitive Neuropsychology* **10**: 377–500.

Plaut DC, McClelland JL, Seidenberg MS and Patterson KE (1996) Understanding normal and impaired word reading: computational principles in quasi-regular domains. *Psychological Review* **103**: 56–115.

Posner MI, Walker JA, Friedrich FJ and Rafal RD (1984) Effects of parietal lobe injury on covert orienting of visual attention. *Journal of Neuroscience* **4**: 1863–1874.

Pouget A and Sejnowski TJ (1997) Spatial transformations in the parietal cortex using basis functions. *Journal of Cognitive Neuroscience* **9**: 222–237.

Roberts RJ and Pennington BF (1996) An interactive framework for examining prefrontal cognitive processes. *Developmental Neuropsychology* **12**(1): 105–126.

Samsonovich A and McNaughton BL (1997) Path integration and cognitive mapping in a continuous attractor neural network model. *Journal of Neuroscience* **17**: 5900–5920.

Schultz W, Apicella P and Ljungberg T (1993) Responses of monkey dopamine neurons to reward and conditioned stimuli during successive steps of learning a delayed response task. *Journal of Neuroscience* **13**: 900–913.

Schultz W, Apicella P, Romo R and Scarnati E (1995) Context-dependent activity in primate striatum reflecting past and future behavioral events. In: Houk JC, Davis JL and Beiser DG (eds) *Models of Information Processing in the Basal Ganglia*, pp. 11–28. Cambridge, MA: MIT Press.

Schultz W, Dayan P and Montague PR (1997) A neural substrate of prediction and reward. *Science* **275**: 1593–1599.

Schweighofer N, Arbib M and Kawato M (1998a) Role of the cerebellum in reaching quickly and accurately. I: A functional anatomical model of dynamics control. *European Journal of Neuroscience* **10**: 86–94.

Schweighofer N, Arbib M and Kawato M (1998b) Role of the cerebellum in reaching quickly and accurately. II: A detailed model of the intermediate cerebellum. *European Journal of Neuroscience* **10**: 95–105.

Scoville WB and Milner B (1957) Loss of recent memory after bilateral hippocampal lesions. *Journal of Neurology, Neurosurgery, and Psychiatry* **20**: 11–21.

Seidenberg MS and McClelland JL (1989) A distributed, developmental model of word recognition and naming. *Psychological Review* **96**: 523–568.

Seung HS (1998) Continuous attractors and oculomotor control. *Neural Networks* **11**: 1253–1258.

Shallice T (1982) Specific impairments of planning. *Philosophical Transactions of the Royal Society, Series B* **298**: 199–209.

Shiffrin RM and Schneider W (1977) Controlled and automatic human information processing. II: Perceptual learning, automatic attending, and a general theory. *Psychological Review* **84**: 127–190.

Squire LR (1992) Memory and the hippocampus: a synthesis from findings with rats, monkeys, and humans. *Psychological Review* **99**: 195–231.

Sutton RS (1988) Learning to predict by the method of temporal diferences. *Machine Learning* **3**: 9–44.

Sutton RS and Barto AG (1998) *Reinforcement Learning: An Introduction*. Cambridge, MA: MIT Press.

Swindale NV (1996) The development of topography in the visual cortex: a review of models. *Network: Computation in Neural Systems* **7**: 161–247.

Treves A and Rolls ET (1994) A computational analysis of the role of the hippocampus in memory. *Hippocampus* **4**: 374–392.

Vecera SP and O'Reilly RC (1998) Figure–ground organization and object recognition processes: an interactive account. *Journal of Experimental Psychology: Human Perception and Performance* **24**: 441–462.

Verfaellie M, Rapcsak SZ and Heilman KM (1990) Impaired shifting of attention in Balint's syndrome. *Brain and Cognition* **12**: 195–204.

Wickens J (1997) Basal ganglia: structure and computations. *Network: Computation in Neural Systems* **8**: 77–109.

Zipser D, Kehoe B, Littlewort G and Fuster J (1993) A spiking network model of short-term active memory. *Journal of Neuroscience* **13**: 3406–3420.