# Generalisation benefits of output gating in a model of prefrontal cortex

Trent Kriete [a] & David C. Noelle [a]

[a] Cognitive and Information Sciences, University of California,
Merced, 5200 North Lake Road, Merced, CA, 95343, USA

Available online: 27 May 2011

PLEASE SCROLL DOWN FOR ARTICLE

Taylor & Francis
Taylor & Francis Group

# Generalisation benefits of output gating in a model of prefrontal cortex

Trent Kriete and David C. Noelle*

*Cognitive and Information Sciences, University of California, Merced, 5200 North Lake Road, Merced, CA 95343, USA*

The prefrontal cortex (PFC) plays a central role in flexible cognitive control, including the suppression of habitual responding in favour of situation-appropriate behaviours that can be quite novel. PFC provides a kind of working memory, maintaining the rules, goals, and/or actions that are to control behaviour in the current context. For flexible control, these PFC representations must be sufficiently componential to support systematic generalisation to novel situations. The anatomical structure of PFC can be seen as implementing a componential 'slot-filler' structure, with different components encoded over isolated pools of neurons. Previous PFC models have highlighted the importance of a dynamic gating mechanism to selectively update individual 'slot' contents. In this article, we present simulation results that suggest that systematic generalisation also requires an 'output gating' mechanism that limits the influence of PFC on more posterior brain areas to reflect a small number of representational components at any one time.

**Keywords:** prefrontal cortex; dopamine; gating; componential representation; systematic generalisation

## 1. Introduction

We have studied the ability of a leading model of the working memory circuits of the prefrontal cortex (PFC) (O'Reilly and Frank 2006) to learn componential representations from limited experience. This model contains three main features that support the learning of componential representations. First, based on anatomical studies of frontal connectivity, the model contains isolated 'stripes' of neurons that act as generic 'slot-filler' pools for the encoding of independent representational components. Second, a reinforcement learning mechanism, based on the midbrain dopamine (DA) system, learns to control a separate 'input gate' on each PFC stripe, allowing the system to learn when maintained patterns of activity should be updated. Third, a synaptic plasticity mechanism that incorporates both Hebbian and error-correction learning is used to shape representations within each PFC stripe, as well as in circuits that use the contents of the PFC working memory to generate responses. This model was trained to represent, actively maintain, and apply a sequence of stimulus–response rule instructions (e.g. 'when asked to generate the

second action, clap your hands'), and the ability of the network systematically to generalise to novel sequences was assessed.

While the architecture of this model supports the formation of componential representations in the PFC, systematic generalisation is hindered by the need to *decode* these representations to generate responses. In order to support systematic generalisation further, we have examined the utility of incorporating an 'output gate' on each PFC stripe, causing the activity in a given stripe to only be propagated outside of the PFC if a specific midbrain gating signal is provided. This mechanism is similar to the output gating mechanism in long short-term memory networks (Hochreiter and Schmidhuber 1997), but it differs in that (1) gating acts on full stripes rather than on individual units and (2) the gating signal is learned through a DA-based reinforcement learning mechanism rather than using the biologically implausible backpropagation-through-time algorithm. We show that the incorporation of an 'output gating' mechanism both speeds learning and greatly improves systematic generalisation.

## 2.  Background

The PFC has long been implicated in working memory and cognitive control (Fuster and Alexander 1971). Neurons in this area of the brain are seen as actively maintaining task-relevant information through sustained firing, and this sustained activity is seen as modulating processing in more posterior brain circuits, instantiating 'top-down control' by causing those circuits to respond in a manner that differs from their unmodulated, usual, more automatic pattern of responding (Miller and Cohen 2001). Data from electrophysiological and brain imaging studies have supported the special role of this brain region in active memory, specifically revealing *delay period activity*: neural firing patterns that persist during delays between the brief presentation of a cue and the generation of a response based on that cue. Delay period activity in the PFC neurons can encode many different kinds of task-relevant information, including spatial locations (Funahashi, Bruce, and Golman-Rakic 1989), non-spatial stimulus features (Cohen *et al.* 1994), and even abstract rules (Wallis, Anderson, and Miller 2001). This view of PFC function has led to an array of computational models, refining our understanding of how this neural system supports flexible cognition as well as how it might develop and adapt.

Early computational models of the PFC offered an explanation for the sustained delay period activity in terms of attractor dynamics in recurrent neural networks (Amit 1989). Recurrent excitation provided a mechanism through which firing rates could be maintained, even after the removal of the initiating stimulus and even in the face of other distracting inputs. This view of active maintenance in the PFC is also well supported by anatomical data, which shows increased recurrent excitation in frontal cortex when compared with more posterior cortical areas (Lewis and Gonzalez-Burgos 2000). In order to act as a useful working memory, however, the PFC must flexibly toggle between a state of active maintenance and a state of 'updating' – loading a new pattern of activation. Computational models have suggested that recurrent excitation, alone, is insufficient to intelligently manage this selective updating process, leading to a focus on special maintenance currents found in the PFC neurons (Camperi and Wang 1998). This naturally leads to the question of how the brain uses these currents to adaptively update the contents of the PFC.

Working memory contents should be updated so as to retain that information which is most important for task success in the current context. Since the ability to perform many critical cognitive tasks is learned, the utility of information to be considered for active maintenance in the PFC must also be learned from experience. This focus on learning led to the incorporation of the midbrain DA system into models of PFC function. Inspired by seminal work on the effect of tonic DA levels on stabilising activation patterns in the PFC (Cohen and Servan-Schreiber 1992) and

on more recent successes in modelling the role of phasic DA bursts in the learning of sequential action from sparse reinforcement (Montague, Dayan, and Sejnowski 1996) formal accounts of PFC updating were developed, in which, frontal maintenance currents were manipulated by a DA-based reinforcement learning circuit implementing the method of temporal differences (Sutton 1988). By learning, from experience, when to actively maintain the current contents of the PFC and when to update this pattern of activation, this approach was able to model human performance in working memory and cognitive control domains such as the *AX continuous performance task* (Braver and Cohen 2000).

While this approach to modelling PFC function proved fruitful, open problems remained. One such problem involved explaining how appropriate PFC representations developed. While some DA-based updating models explored issues of frontal representation (O'Reilly, Noelle, Braver, and Cohen 2002), these models invariably encoded actively maintained information in a format that was directly specified by the model designer. It was not clear that standard neural learning algorithms could produce the kind of componential representations that would allow for the encoding of highly novel working memory contents, such as an abstract behavioural rule (e.g. 'when the string of letters on the computer screen is an English word, depress the keyboard space bar as quickly as possible'). While some connectionist models had been shown to develop componential attractor dynamics (Plaut and McClelland 1993), these models stipulated the componential structure of their representations through direct teaching signals. The development of such componential representations at internal neural layers, as would be the case for the PFC, proved to be elusive without strong innate constraints (Noelle and Cottrell 1996; Noelle and Zimdars 1999). In the face of these difficulties, it was shown that a DA-based updating mechanism, coupled with rich training on a variety of related cognitive tasks, could support the development of componential PFC representations (Rougier, Noelle, Braver, Cohen, and O'Reilly 2005). The resulting computational model displayed good qualitative and quantitative fits to both healthy and frontally damaged human performance on multiple clinical measures of working memory updating and cognitive control (i.e. the Wisconsin Card Sort Task and the Stroop task).

A major limitation of these DA-based updating models involved the global nature of the PFC updating signal. Either the entire contents of the PFC working memory system were actively maintained or the entire contents were updated. There was no mechanism in place to allow for the maintenance of some information in the PFC while some other information was updated. Inspired by anatomical and physiological evidence suggesting that dense recurrent excitation in the PFC is organised into relatively isolated sets of 'stripes' (Lewis and Gonzalez-Burgos 2000), potentially supporting many separate attractor networks, models were constructed that allowed for the independent updating of isolated 'stripes' of the PFC neurons. Early models of this kind reverted to the use of hand-designed PFC representations (Frank, Loughry, and O'Reilly 2001), but later models demonstrated that appropriate representations could develop from experience in a stripe-structured PFC (O'Reilly and Frank 2006). These models also introduced some further anatomical detail, with the PFC updating controlled by loop-like projections between the PFC and the basal ganglia (BG), most notably through striatal matrisomes (Matrix). The DA contribution to learning in the BG is modelled using a reinforcement learning algorithm called *PVLV* (primary value/learned value), which allows the dorsal striatum to learn, from experience, when it is appropriate to update each PFC stripe. In this way, separate components of the current, actively maintained, PFC representation may be updated separately.

It is important to note that while these stripe-based models provide architectural support for encoding different independent components of a memory trace (e.g. the condition and the action of a condition/action rule) in different pools of the PFC neurons, the discovery of isolated components and their assignments to stripes must emerge from a learning process. Thus, these networks still require a rich collection of training experiences in order to discover the different components that might need to be independently maintained. Simulation results suggest that the BG-based

PFC updating mechanism can do a fairly good job of segregating representational components across the PFC stripes, but generalisation to novel combinations of components is still limited. This limitation arises because the full contents of the PFC, across all stripes, are transmitted in parallel to more posterior brain areas in order to produce responses. These posterior circuits are shaped by standard neural learning mechanisms, making them susceptible to the same componential generalisation problems observed in other connectionist networks. Despite the fact that the PFC representations modulating these posterior circuits come to have a relatively clean componential structure, the posterior circuits can come to depend upon spurious correlations observed during training, causing them to fail when the PFC actively maintains a novel combination of components unlike those seen during the development of the posterior system.

This limitation suggests an additional computational mechanism that might improve the use of stripe-based componential representations in the PFC. The idea is to restrict the transmission of actively maintained PFC patterns to posterior brain areas so that the posterior circuits only process the contents of one (or a small number) of the PFC stripes at any one time. In this way, posterior circuits do not need to learn to 'decode' the full componential PFC representation, but only need to learn to process individual components (or a small number of components) at a time. We refer to such a mechanism as a stripe-based *output gating mechanism*, following the metaphor of a gated enclosure. When the *output gate* on a given PFC stripe is closed, actively maintained patterns of firing in that stripe are *not* communicated downstream to more posterior brain areas. When the *output gate* on a given PFC stripe is open, in contrast, the contents of that stripe are made visible to posterior circuits. Exactly when to open and close the various output gates needs to be learned. In the model presented here, the same reinforcement learning algorithm that is used to learn when to update a PFC stripe is also used to determine when to open the output gate on each stripe. The goal of the initial simulations reported here is to show that such a learned stripe-based output gating mechanism can improve network generalisation to novel combinations of components stored in the PFC working memory system.

## 3. The model

### 3.1. *An output gating network*

In order to investigate the computational benefits of an output gating mechanism, we constructed a PFC model that extended a current stripe-based updating model (O'Reilly and Frank 2006). Both the previous model and the one reported here made use of the *Leabra* computational cognitive neuroscience modelling framework (O'Reilly and Munakata 2000). This biologically grounded framework uses a rate-coded model of neural firing, supports extensive bidirectional excitation and rapid pooled inhibition, and incorporates synaptic plasticity driven by both correlated firing (Hebbian) and error correction (contrastive Hebbian learning).

A schematic diagram of the output gating network is shown in Figure 1. In that diagram, the *PFC maint* box represents a collection of isolated pools, or stripes, of excitatory PFC pyramidal cells, with each stripe equipped with recurrent connections and special maintenance currents that allow distributed patterns of activity to be actively maintained. Pooled lateral inhibition within these stripes limits the total amount of activity, constraining runaway excitation that might arise from positive feedback. The *matrix maint* box represents the circuits in the BG that drive the updating of the corresponding PFC stripes in PFC maint. The matrix maint circuit learns when to update the PFC contents from the PVLV reinforcement learning algorithm, simulating the contribution of the DA system in modifying synaptic strengths in the BG. Inputs to PFC maint and matrix maint encode information that may, at certain times, need to be stored in the PFC. The nature of these inputs in our simulations is discussed below.
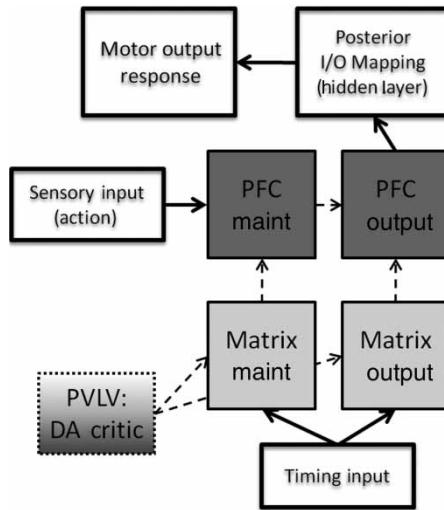
Figure 1. Output gating network diagram.

In previous stripe-based PFC models (O'Reilly and Frank 2006), the activation levels of all of the units in PFC maint were continuously sent to a hidden layer that represented relevant *posterior* brain areas. This hidden layer then drove an explicit *motor output response*, which determined if reward was provided to reinforcement learning aspects of the network. As shown in Figure 1, our output gating network does *not* include projections directly from the PFC maint to posterior brain areas. Instead, the neural stripes in PFC maint are mirrored in *PFC output*, with hard-wired direct connections allowing the pattern of activity in PFC maint to be directly reflected in the firing of the PFC output units. Individual stripes in PFC output only come to reflect their corresponding PFC maint inputs when they are updated by signals from the BG in *matrix output*. Just as the matrix maint circuit uses PVLV reinforcement learning to learn when to update the PFC maint stripes, the matrix output circuit uses the same reinforcement learning mechanism to learn when to update the PFC output stripes, effectively learning which stripes in PFC maint to make visible to posterior brain areas at any given point in time. This is how output gating is implemented in this model. It is worth noting that the same updating mechanism is used for both PFC maint and PFC output, with the only difference being the nature of the inputs provided to the two different collections of the PFC stripes. Mapping this model onto anatomy, the PFC output stripes might be seen as stripes that are closer to the frontal motor areas, e.g. SMA (supplementary motor area), but we prefer to interpret them as the output neurons (i.e. cortical layers 5 and 6) in the patches corresponding to the PFC maint stripes (i.e. neurons in cortical layers 2 and 3).

In the simulations reported here, we compare the generalisation performance of this output gating network with that of a network without output gating. The network without output gating has the same basic architecture as that shown in Figure 1, except it lacks PFC output and matrix output, and it projects the contents of the PFC maint stripes directly, in parallel, to the posterior hidden layer.

## 3.2. *A generalisation task*

The primary question of interest is whether or not an output gating mechanism would support improved generalisation when novel combinations of PFC representational components were actively maintained. In order to investigate this question, a simple sequence learning task was

used. The network was expected to remember, using active maintenance in the PFC, a presented sequence of actions and then produce that sequence when prompted. Sequence elements were presented as paired associates, with each element including its ordinal location in the sequence (e.g. '1', '2', '3'), provided at the *timing input* layer, and the action to be taken at that position in the sequence (e.g. 'clap', 'jump', 'spin'), provided at the *sensory input (action)* layer (Figure 1).[1] Both inputs were encoded in a localist fashion, with one unit becoming active for each ordinal position and each action. When a sequence element was presented to the network, it was expected to echo the current sensory input (action) pattern at its motor output response layer. After the presentation of several sequence elements, the network's memory was queried by presenting a timing input without an action, with the expectation that the network would produce the action corresponding the given sequence position at its motor output response layer. Every correct output produced a reward signal, provided to PVLV in order to learn when to update the PFC stripes. Synaptic strengths were also modified using the standard Leabra neural learning mechanisms, with target responses provided at the motor output response layer.

In order to examine generalisation performance, a testing set of sequences was pre-generated, and the network was trained on sequences uniformly sampled from the remaining space of possible sequences. New sequences were sampled, with replacement, throughout training, allowing the network to experience an increasingly broad range of sequences as training continued. For convenience, one 'epoch' of training was recorded for every 500 sequence presentations. Generalisation was assessed in terms of the network's performance on the novel testing sequences, which were never experienced during training. Since the size of the space of sequences grows dramatically with the action *vocabulary* of the network (i.e. the number of different possible actions), we systematically varied the vocabulary size, using collections of 10, 20, and 30 actions for different networks. For the results reported here, the sequence length was set at three actions.

We additionally examined a training regime in which the training patterns were specifically designed to include spurious correlations. Five of the actions were selected, and training sequences were not allowed to contain more than one of the actions from this set of five. Thus, during training, the network never experienced sequences that involved any pair of actions from this set of five. The testing set then involved sequences that made use *only* of these five actions, ensuring their appearance together in individual sequences. This 'restricted pairs' training process was designed to provide a stronger test of the ability of these networks to handle the novel recombination of the PFC components.

## 4. Simulation results

Networks were trained on uniformly sampled training sequences until they reached a stringent performance criterion on those sequences, making fewer than one action output error per 500 sequence presentations. The networks were then tested on a separate testing set of 100 sequences, and the number of output errors, as a percentage of the total number of outputs produced while processing the 100 testing sequences, was measured. This process was repeated, with randomised initialisation and randomised selection of training and testing sequences, for 20 networks in each condition. The average generalisation performance results are shown in Figure 2, with error bars displaying one standard error of the mean. Note that while both networks with output gating and networks without display a relatively low generalisation error, the networks with an output gating mechanism perform significantly better, and the benefit of the output gating mechanism grows with the size of the action vocabulary.

The advantage of output gating is more strongly evident in the 'restricted pairs' training case, where spurious correlations were purposely included in the set of training sequences. These results,
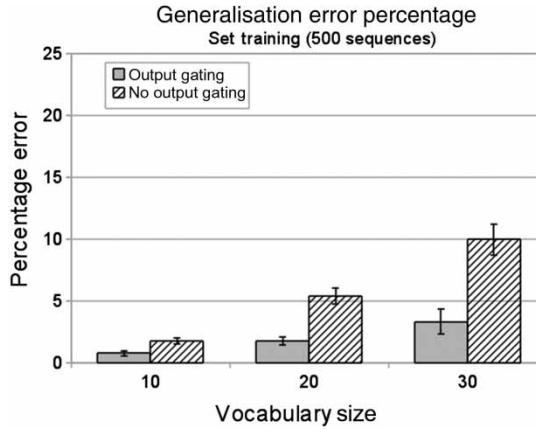
Figure 2. Generalisation results for training to a criterion on uniformly sampled training and testing sequences.
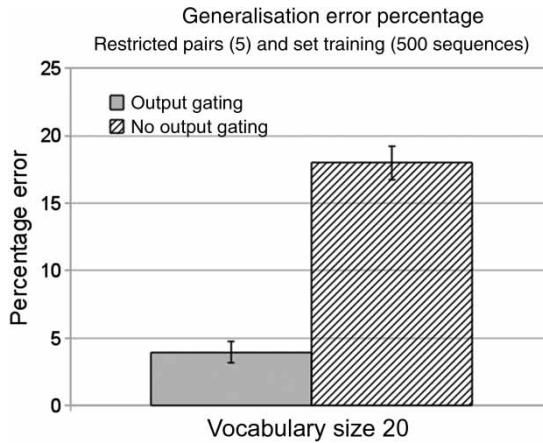


Figure 3. Generalisation results for training to a criterion on 'restricted pairs' training sequences.

using an action vocabulary size of 20, are summarised in Figure 3. Once again, error bars reflect one standard error of the mean. Note the more than four-fold increase in error when output gating is removed.

In addition to examining generalisation performance once a stringent criterion had been reached on training sequences, we measured generalisation performance during the training process, allowing us to obtain a sense of the speed of learning. As shown in Figure 4, the use of an output gating mechanism speeded the learning process, particularly for larger action vocabularies. This difference was particularly salient when 'restricted pairs' training was used, as shown in Figure 5.

## 5. Discussion

The referents of thought are often naturally described as having a componential structure. For example, an explicit stimulus–response rule combines a condition for rule application with an action to take. Componential descriptions are most apt when the parts are somewhat independent, producing a combinatorial space of different well-formed objects. Many cognitive processes
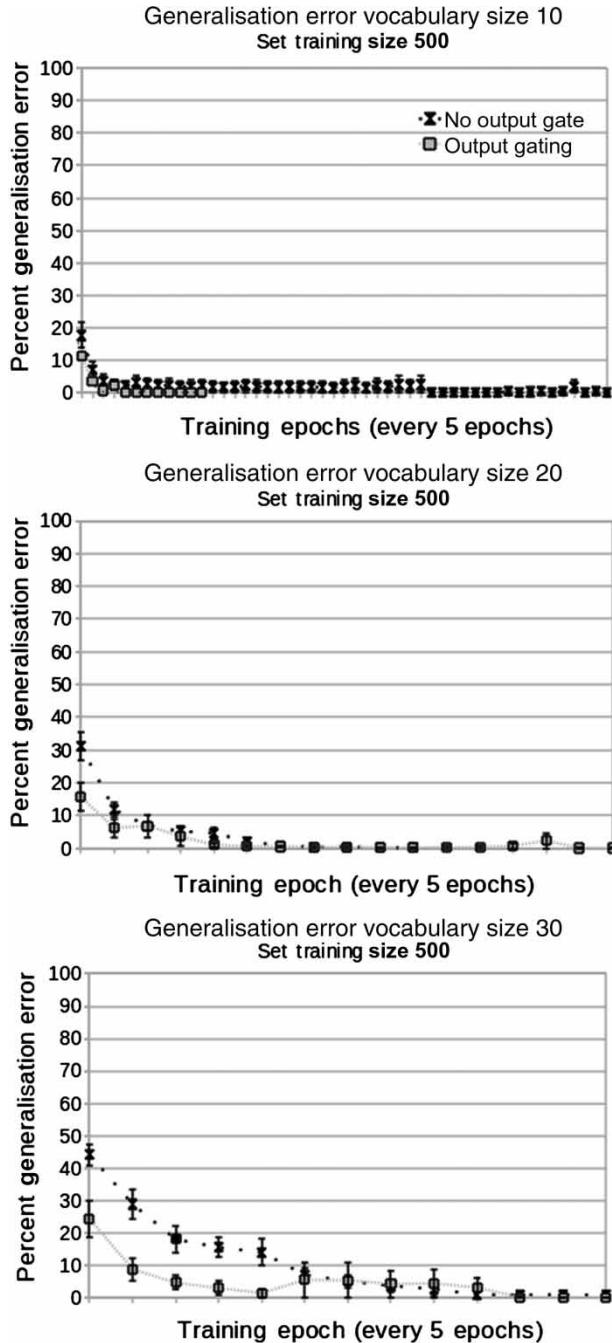
*T. Kriete and D.C. Noelle*

Figure 4.   Generalisation during training on uniformly sampled training and testing sequences (vocabulary size varying across graphs).

appear to leverage this structure, exhibiting substantial systematicity, matching the quasi-regular structure of the environment (Plaut, McClelland, Seidenberg, and Patterson 1996).

Thus, in some way, neural representations must reflect the componential structure of the world. Many computational models use a kind of 'slot-filler' approach, with a separate pool of neural processing elements dedicated to each part of a componential representation. Alternatively, part
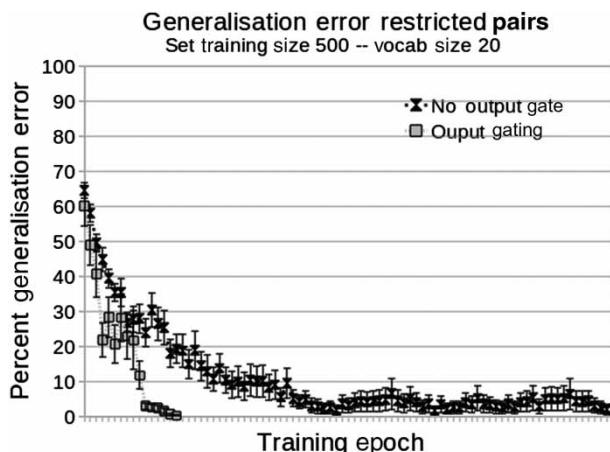
Figure 5. Generalisation during training on 'restricted pairs' training sequences.

representations may be superimposed across a common set of neural processing elements (Pollack 1990; Plate 1995). Learning must play a central role in the shaping of componential representations, since only experience reveals the independently varying parts of objects, as well as constraints on their combination. This gives rise to the question of how the brain learns such componential representations from experience.

Standard statistical learning methods often fail in this regard. The problem stems from spurious correlations that arise from limited experience. Because the space of possible combinations of parts is often large, some partial combinations will never be observed, just by chance, but this does not entail that those combinations are ill-formed or meaningless. For example, the fact that one has never been asked to press the 'J key' whenever a red object appears on the screen does not prohibit the representation of this explicit rule in working memory or its application. Most statistical learning techniques would transform such a perfect anti-correlation into a constraint, however, limiting systematic generalisation.

This learning problem is particularly pronounced in recurrent attractor networks, like those used in models of human PFC. Spurious correlations are amplified by the recurrent flow of activity in such networks, resulting in a failure to form componential representations (Noelle and Cottrell 1996) unless a strong inductive bias is employed (Noelle and Zimdars 1999). The learning of componential representations is particularly important in the PFC, however, as the working memory circuits in this brain area play a central role in allowing us to flexibly enact novel behaviours.

Previous work on the stripe-based PFC models has suggested that the striped organisation of PFC provides a natural slot-filler structure upon which componential representations can develop. The use of a DA-based reinforcement learning mechanism to selectively update the contents of the PFC stripes can give rise to appropriate representations. The work reported here suggests, however, that the development of such representations is not sufficient to produce flexible generalisation. The problem arises when more posterior brain areas must decode the componential representations of the PFC in order to generate behaviour. Without a strong innate inductive bias to respect the stripe-based structure of the PFC, posterior circuits fall prey to spurious correlations and fail to properly respond to novel PFC contents. Here, we suggest that the brain handles this problem by limiting the influence of the PFC on posterior areas to that encoded by a small number of PFC stripes at any one time. This limitation encourages posterior areas to learn to process the individual components of the PFC representations, rather than arbitrary combinations

of components, with the full impact of the PFC contents unfolding over time. In essence, this is the same solution suggested by Elman (1990), where a common set of connection weights traverse componential structures over time. Thus, we suggest that systematic generalisation arises from (1) distributed neural codes within the PFC stripes and in more posterior brain areas, (2) slot-filler structures built from isolated pools of neural units that encode different representational components, supported by the isolated neural stripes in PFC anatomy, and (3) the temporally extended and sequential application of selective PFC stripe contents to the rest of the brain, controlled by a learned output gating mechanism. In this way, both distributed representations and more localist mechanisms are combined in order to support processing of novel structured representations. Distributed codes provide an efficient means of encoding fillers, as well as support some 'interpolation-based' generalisation to slightly novel filler values. The slots, or stripes, composing structured representations in the PFC, are processed in a more localist fashion, with a dedicated pool of neurons encoding each slot and an adaptive output gating mechanism supporting the processing of each slot in isolation. Unrolling the processing of structural components in time, using output gating, biases neural learning mechanisms away from capturing and enforcing correlations between the PFC stripe contents, including spurious correlations, thus supporting the learning of structured representations that are robust to the novel recombination of slot-filler components, allowing for systematic generalisation in performance.

## Acknowledgement

## Note

1. One might imagine the two network inputs as having different sources, with the timing input coming from some form of internal clock, perhaps in the cerebellum, and the sensory input (action) coming from high-level perceptual circuits.

## References

Amit, D.J. (1989), *Modeling Brain Function: The World of Attractor Neural Networks*, Cambridge, UK: Cambridge University Press.

Braver, T.S., and Cohen, J.D. (2000), 'On the Control of Control: The Role of Dopamine in Regulating Prefrontal Function and Working Memory', in *Control of Cognitive Processes: Attention and Performance XVIII* (Chap. 31), eds. S. Monsell and J. Driver, Cambridge, MA: MIT Press, pp. 713–737.

Camperi, M., and Wang, X.J. (1998), 'A Model of Visuospatial Working Memory in Prefrontal Cortex: Recurrent Network and Cellular Bistability', *Journal of Computational Neuroscience*, 5, 383–405.

Cohen, J.D., and Servan-Schrieber, D. (1992), 'Context, Cortex, and Dopamine: A Connectionist Approach to Behavior and Biology in Schizophrenia', *Psychological Review*, 99(1), 45–77.

Cohen, J.D., Forman, S.D., Braver, T.S., Casey, B.J., Servan-Schreiber, D., and Noll, D.C. (1994), 'Activation of Prefrontal Cortex in a Nonspatial Working Memory Task with Functional MRI', *Human Brain Mapping*, 1, 293–304.

Elman, J.L. (1990), 'Finding Structure in Time', *Cognitive Science*, 14(2), 179–211.

Frank, M.J., Loughry, B., and O'Reilly, R.C. (2001), 'Interactions Between Frontal Cortex and Basal Ganglia in Working Memory: A Computational Model', *Cognitive, Affective, and Behavioral Neuroscience*, 1(2), 137–160.

Funahashi, S., Bruce, C.J., and Golman-Rakic, P.S. (1989), 'Mnemonic Coding of Visual Space in the Monkey's Dorsolateral Prefrontal Cortex', *Journal of Neurophysiology*, 61, 331–349.

Fuster, J.M., and Alexander, G.E. (1971), 'Neuron Activity Related to Short-Term Memory', *Science*, 173, 652–654.

Hochreiter, S., and Schmidhuber, J. (1997), 'Long Short Term Memory', *Neural Computation*, 9, 1735–1780.

Lewis, D.A., and Gonzalez-Burgos, G. (2000), 'Intrinsic Excitatory Connections in the Prefrontal Cortex and the Pathophysiology of Schizophrenia', *Brain Research Bulletin*, 52, 309–317.

Miller, E.K., and Cohen, J.D. (2001), 'An Integrative Theory of Prefrontal Cortex Function', *Annual Review of Neuroscience*, 24, 167–202.

Montague, P.R., Dayan, P., and Sejnowski, T.J. (1996), 'A Framework for Mesencephalic Dopamine Systems Based on Predictive Hebbian Learning', *Journal of Neuroscience*, 16, 1936–1947.

Noelle, D.C., and Cottrell, G.W. (1996), 'In Search of Articulated Attractors', in *Proceedings of the 18th Annual Conference of the Cognitive Science Society*, ed. G.W. Cottrell, La Jolla, CA: Lawrence Erlbaum, pp. 329–334.

Noelle, D.C., and Zimdars, A.L. (1999), 'Methods for Learning Articulated Attractors over Internal Representations', in *Proceedings of the Twenty First Annual Conference of the Cognitive Science Society*, eds. M. Hahn and S.C. Stones, Vancouver, BC: Lawrence Erlbaum, pp. 480–485.

O'Reilly, R.C., and Frank, M.J. (2006), 'Making Working Memory Work: A Computational Model of Learning in the Frontal Cortex and Basal Ganglia', *Neural Computation*, 18(2), 283–328.

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

O'Reilly, R.C., Noelle, D.C., Braver, T.S., and Cohen, J.D. (2002), 'Prefrontal Cortex and Dynamic Categorization Tasks: Representational Organization and Neuromodulatory Control', *Cerebral Cortex*, 12(3), 246–257.

Plate, T.A. (1995), 'Holographic Reduced Representations', *IEEE Transactions on Neural Networks*, 6(3), 623–641.

Plaut, D.C., and McClelland, J.L. (1993), 'Generalization with Componential Attractors: Word and Nonword Reading in an Attractor Network', in *Proceedings of the 15th Annual Conference of the Cognitive Science Society*, Boulder, CO: Lawrence Erlbaum, pp. 824–829.

Plaut, D.C., McClelland, J.L., Seidenberg, M.S., and Patterson, K. (1996), 'Understanding Normal and Impaired Word Reading: Computational Principles in Quasi-regular Domains', *Psychological Review*, 103(1), 56–115.

Pollack, J.B. (1990), 'Recursive Distributed Representations', *Artificial Intelligence*, 46(1–2), 77–105.

Rougier, N.P., Noelle, D.C., Braver, T.S., Cohen, J.D., and O'Reilly, R.C. (2005), 'Prefrontal Cortex and Flexible Cognitive Control: Rules Without Symbols', *Proceedings of the National Academy of Sciences*, 102(20), 7338–7343.

Sutton, R.S. (1988), 'Learning to Predict by the Methods of Temporal Differences', *Machine Learning*, 3, 9–44.

Wallis, J.D., Anderson, K.C., and Miller, E.K. (2001), 'Single Neurons in Prefrontal Cortex Encode Abstract Rules', *Nature*, 411, 953–956.