

# A Computational Model of Emotional Learning in the Amygdala

**Jan Morén**

Lund University Cognitive Science  
Kungshuset, Lundagård  
S-222 22 LUND, Sweden  
jan.moren@lucs.lu.se

**Christian Balkenius**

Lund University Cognitive Science  
Kungshuset, Lundagård  
S-222 22 LUND, Sweden  
christian.balkenius@lucs.lu.se

## Abstract

The amygdala has repeatedly been implicated in emotional reactions and in learning of new emotionally significant stimuli. The system forms an important part of motor learning as well as attention. This paper presents a neurologically inspired computational model of the amygdala and the orbitofrontal cortex that aims to partially reproduce the same characteristics as the biological system. This model has been tested in simulations, the results of which are presented.

## 1. Introduction

The amygdala is a small structure in the medial temporal lobe that is thought to be responsible for the emotional evaluation of stimuli (Rolls, 1995). This evaluation is in turn used as a basis for emotional states, for emotional reactions and is used to signal attention and laying down long-term memories. As a brain structure, the amygdala is old (it can be readily identified in fish, for example), and fairly uniform in large-scale structure across species (McDonald, 1992).

We believe that the amygdalo-orbitofrontal system implements part of a two-process model of learning, as described by Mowrer (Mowrer, 1973). This approach separates learning into first a stimulus-emotional system that evaluates incoming stimuli, and a second learning system that uses this evaluation as a reinforcer for stimulus-response learning. Among the advantages of this approach is that the motivation to respond and the response itself are cleanly separated (Rolls, 1986)

As part of our investigation, we have implemented a computational model of the amygdala and the orbitofrontal cortex, and are testing this in simulation.

This is not a detailed physiological model, even though it shares its larger-scale structure with that of the real amygdalo-orbitofrontal system; instead, our aim is to make use of neurophysiological data to construct a *functional* model of emotional processing as part of a general learning system.

Another paper on this model, with a stronger emphasis on the neurophysiological side, can be found in (Balkenius and Morén, 2000b). For an overview of the amygdalo-orbitofrontal system as a part of a larger motivational system, see (Morén and Balkenius, 2000).

## 2. Physiology

The prevailing view of the amygdala is that it is responsible for emotional processing (LeDoux, 1995; Rolls, 1995). It has extensive interconnections with many other areas, especially higher sensory cortical areas, smell and taste, and the basal ganglia (assumed to handle reinforcement of motor actions (Gray, 1995; Houk et al., 1995)). This makes this structure very well placed to handle complex evaluation of multi-modal combinations of stimuli.

It is believed that the amygdala has several functions. There is evidence both for a role in directing attention as well as for laying down long-term memories. Our interest in this structure is currently focused on its function as an evaluator of emotionally significant stimuli, however.

The orbitofrontal cortex is thought to inhibit inappropriate responses from the amygdala, based on the context given by the hippocampus.

### 2.1 The amygdala

The amygdala is organized mainly in a feed-forward fashion, i.e. there is extensive flow from

the lateral and basal areas to the medial and central areas, and less information flow in the opposite direction (figure 1). The lateral and basal areas are cortex-like and receive inputs from many higher sensory cortices, such as the prefrontal temporal cortex and the olfactory bulb. These areas have extensive interconnections, and project into the central and medial area the central and medial areas that are noncortex-like. The central and the medial nuclei receive partially different projections from the lateral amygdala and project into the brainstem and the hypothalamus (McDonald, 1992).

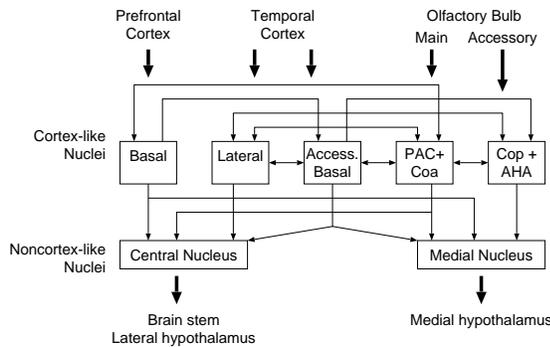


Figure 1: The large-scale structure of the amygdala. Of interest is the extensive interconnections between the cortex-like nuclei and the feed-forward flow between the cortex-like and noncortex-like nuclei. Redrawn from (McDonald, 1992).

The amygdala reacts to a number of innate stimuli that have an a priori emotional charge, such as hunger, pain, certain smells and other stimuli that are intrinsically important. Whenever such a stimulus is encountered, the amygdala will elicit a set of responses that are used in motor learning and in the attention systems. One particular set of stimuli is received directly from the thalamus, rather than from the sensory cortices (McIntosh and Gonzalez-Lima, 1998; LeDoux, 1995). This works as an early, fast sensory classification system (see section 5).

Thus far the system is not much more than a generator of simple punishments and rewards. What makes the amygdala important is that it is instrumental in learning new associations between emotionally charged and neutral stimuli. Whenever a neutral stimulus occurs in close temporal proximity to an emotionally charged stimuli, this stimulus will itself subsequently be able to elicit that same or a similar reaction.

There is ample evidence that learning does occur in the amygdala. Experiments by Weiskrantz (described in (Rolls, 1986)) shows that animals with

a bilateral lesion in the amygdala can no longer acquire an emotional reaction to novel stimuli in the presence of an aversive stimuli. (Rolls, 1992) reviews evidence that bilateral damage produces effects similar to the Klüver-Bucy syndrome in monkeys, resulting in tameness, lack of emotional responses, loss of appetite and consuming of previously rejected foods (Kolb and Wishaw, 1990). Especially interesting, the monkeys are unable to learn an active avoidance task, where a light would signal the imminent onset of a shock unless the subject actively avoided it, a clear indication that they were unable to associate the light with fear.

## 2.2 The orbitofrontal cortex

We view this evaluative system as composed of at least two interacting parts: the amygdala and the orbitofrontal cortex. Whereas the amygdala learns appropriate associations between neutral and emotionally charged stimuli, the orbitofrontal cortex inhibits the expression of these associations as needed depending on context and other factors.

The orbitofrontal cortex is known to inhibit areas it is connected to, and it has extensive interconnections with the hippocampus. There are also connections both to and from the amygdala. (Shimamura, 1995) has argued that the role of the cortical areas is to inhibit posterior areas whenever their reactions are deemed inappropriate due to changing context or reinforcement.

Some evidence for this comes from patients with a damaged orbitofrontal cortex. Shimamura tested patients on the Wisconsin card sorting test. Patients are asked to sort cards according to an unknown criteria such as color or value; the subjects would only be told if a given card has been correctly sorted or not. When the patient has figured out the rule, it is changed by the experimenter. Subjects with these lesions are then unable to change their own sorting rule, even when they can verbalize that the rule has changed.

More evidence is found in (Dias et al., 1997). Dias *et al.* shows that inhibitory processes can be found in many regions within the prefrontal cortex, and that the orbitofrontal cortex is especially important in the inhibition of emotional reactions.

## 3. The Model

There seems little doubt that emotional learning occurs in the Amygdala, with first-order conditioning

between primary and secondary stimuli, and perhaps also between secondary stimuli. Also, evidence suggests that these functions are partially separated, with the orbital cortical areas implicated in secondary conditioning.

We have attempted to capture these features in a computational model suitable for comparisons between neurophysiological data and simulations. We hope that this approach will enable us to attain a clearer understanding both of the functions of the amygdala and of the limitations of the model; this would have been difficult to accomplish with a model that is not testable in simulation.

The model is divided into two parts, very roughly corresponding to the amygdala and the orbital frontal cortex, respectively. Of course, these areas are complex, and we have not in any way attempted to capture all of their functionality. The amygdaloid part receives inputs from the thalamus and from cortical areas, while the orbital part receives inputs from the cortical areas and the amygdala only.

The system also receives a reinforcing signal. This signal has been left unspecified, as it is still unclear from where it comes.

Let's take a look at the model in figure 2. There is one  $A$  node for every stimulus  $S$  (including one for the thalamic stimulus). There is also one  $O$  node for each of the stimuli (except for the thalamic node). There is one output node in common for all outputs of the model, called  $E$  above. The  $E$  node simply sums the outputs from the  $A$  nodes, then subtracts the inhibitory outputs from the  $O$  nodes. The result is the output from the model.

The thalamic connection is calculated as the maximum over all stimuli  $S$  and becomes another input to the amygdaloid part:

$$A_{th} = \max(S_i)$$

Unlike other inputs to the amygdala, the thalamic input is not projected into the orbitofrontal part and can not by itself be inhibited. This will be addressed in section 5.

For each  $A$  node, there is a plastic connection weight  $V$ . Any input is multiplied with this weight to become the output from the node. The  $O$  nodes behave analogously, with a connection weight  $W$  applied to the input signal to create an output.

The connection weights  $V_i$  are adjusted proportionally to the difference between the reinforcer and the

activation of the  $A$  nodes. The  $\alpha$  term is a constant used to adjust the learning speed:

$$\Delta V_i = \alpha(S_i \max(0, Rew - \sum_j A_j))$$

This is an instance of a simple associative learning system, not unlike the Rescorla-Wagner model of learning (Rescorla and Wagner, 1972). The real difference is in the fact that this weight-adjusting rule is monotonic, i.e. the weights  $V$  can not decrease. This may at first seem like a fairly substantial drawback; however, there are good reasons for this design choice. Once an emotional reaction is learned, this should be permanent. It is the task of the orbitofrontal part to inhibit this reaction when it is inappropriate. As seen in section 2, there is experimental evidence that this is in fact the correct approach.

The reinforcer for the  $O$  nodes is calculated as the difference between the previous output  $E$  and the reinforcing signal  $Rew$ . In other words, the  $O$  nodes compare the expected and received reinforcer and inhibits the output of the model should there be a mismatch:

$$\Delta W_i = \beta(S_i \sum_j (O_j - Rew))$$

The orbitofrontal learning rule is very similar to the amygdaloid rule. The only – but crucial – difference is that the orbitofrontal connection weight can both increase and decrease as needed to track the required inhibition.  $\beta$  is another learning rate constant.

The node values are then calculated as:

$$A_i = S_i V_i,$$

$$O_i = S_i W_i,$$

and

$$E = \sum_i A_i - \sum_i O_i$$

The  $A$  nodes give outputs proportionally to their contribution in predicting the reward  $Rew$ , while the  $O$  nodes inhibit the output of  $E$  as necessary.

This system works at two levels: the amygdaloid part learns to predict and react to a given reinforcer. This subsystem can never unlearn a connection; once learned, it is permanent, giving the system the ability to retain emotional connections for as long

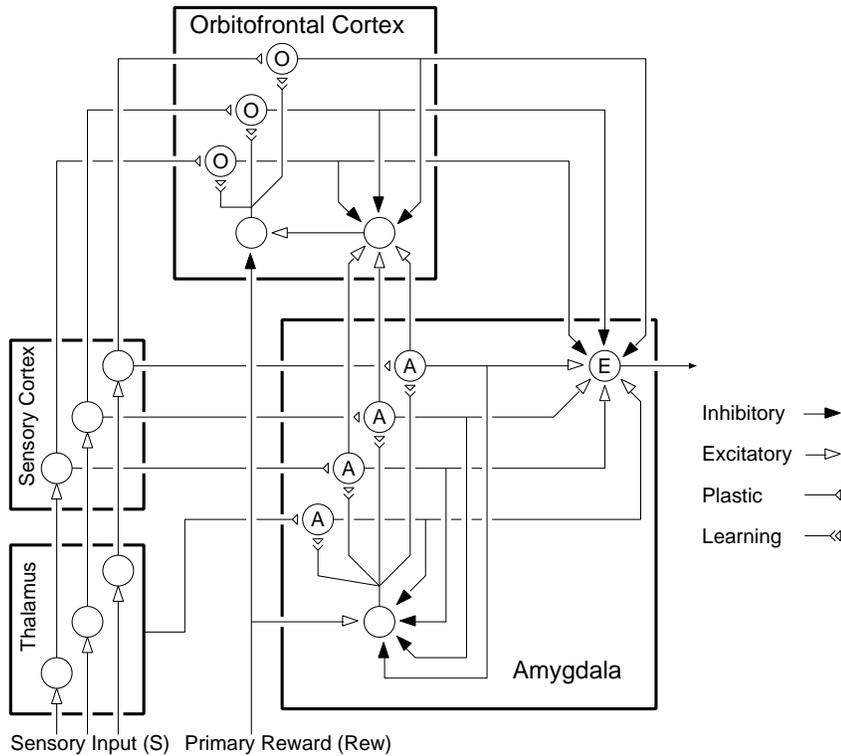


Figure 2: A graphical depiction of the model. At the top is the rudimentary orbitofrontal part (here without an external context), at the bottom right is the amygdaloid part and at left are the thalamic and sensory-cortical modules. The thalamic and sensory-cortical parts are just place-holders in this version of the model. The sensory inputs  $S$  enter the thalamic part, where a thalamic input to the amygdala is calculated as the maximum over all inputs. A primary reward signal  $Rew$  enters both the amygdaloid and orbitofrontal parts.

as necessary. The orbitofrontal system tracks mismatches between the base systems predictions and the actual received reinforcer and learns to inhibit the system output in proportion to the mismatch.

These subsystems receive partially different inputs; the base system receives finely discriminated inputs from the sensory cortex, and a coarse signal  $Th$  from Thalamus. Sensory cortex receives its inputs from thalamus also, and it is assumed that it is responsible for the subdividing and discrimination of the coarse input from thalamus. The thalamic input is a low-level stimulus signal that will be present even in the absence of higher cortical areas (LeDoux, 1995). Although this input seems unnecessary – and indeed harmful – at present, there is evidence that this path does exist (see section 5 for further discussion about this).

The orbitofrontal system currently receives almost the same input as does the amygdaloid system. These inputs work as a quick-and-dirty substitute for a proper context representation. In the near future, however, this system will receive a context representation (that will indirectly include the stim-

uli from sensory cortex) that will enable the model to handle contextual cues properly.

## 4. Simulations

We have run a basic set of simulations to verify some assumptions about the workings of this model. The basic features we have tested are acquisition-extinction-reacquisition, simple blocking and conditioned inhibition.

We have previously run another set of simulations using the model both with and without the orbitofrontal or cortical parts and compared its performance with data from animal studies. These results are available in (Balkenius and Morén, 2000b).

### 4.1 Acquisition

This is a basic learning experiment, where the model is expected to associate a stimulus with a reward/reinforcer, disassociate the stimulus once the reinforcer is absent, then reassociate them again.

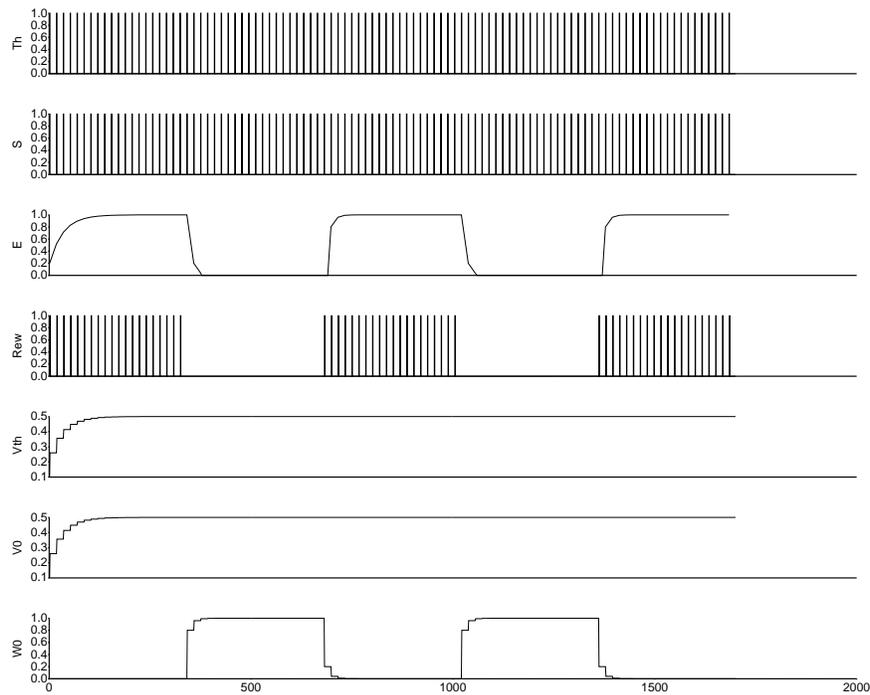


Figure 3: The result of Acquisition, extinction and reacquisition. From top to to bottom, the graphs are:  $Th$ , the thalamic input;  $S$ , the stimulus input;  $E$ , the output of the model (smoothed);  $Rew$ , the reinforcing signal;  $Vth$ , the amygdaloid connection weight for the thalamic input;  $V_0$ , the thalamic weight for  $S$ ; and  $W_0$ , the orbitofrontal connection weight for the stimulus. The learning parameters are  $\alpha = 0.2$  and  $\beta = 0.8$ .

This represents a minimal functionality of any associative learning model.

In figure 3 we see the inputs ( $S$  and  $Th$ , the thalamic connection), the reinforcement signal ( $Rew$ ) and the output ( $E$ ). The acquisition-extinction cycle is repeated three times to see how the system reacts during reacquisition.

As seen in the figure, the system manages to learn this simple association well; the output tracks the reinforcer without any problems.

In this figure, we have included the connection weights for both the amygdaloid and orbitofrontal parts of the model.

The stimulus  $S$  and the thalamic input  $Th$  occur simultaneously and with the same intensity, resulting in  $Vth$  and  $V_0$  sharing the responsibility for the association to  $Rew$  – the highest attained level of these weights are both 0.5.

When the reinforcer disappears, the amygdaloid weights are not affected; instead the orbitofrontal weight  $W_0$  rapidly increases and inhibits the output. As soon as the reinforcer reappears,  $W_0$  decreases to zero, allowing the amygdala to express the previously learned association.

The output  $E$  increases to its full level faster on subsequent trials than it did the first time. This effect is well established in the literature; see cf. (Mackintosh, 1983).

## 4.2 Blocking

In this blocking simulation we show the ability of the model not to associate a stimulus with the reinforcer if there is already an established association that can explain the contingency.

A blocking schedule is run in three phases: first associate  $S_0$  with the reinforcer, then present both  $S_0$  and  $S_1$  together with the reinforcement, and last, test  $S_1$  to see whether it has been associated with the reinforcer. There should be no response to  $S_1$ . This is consistent with experimental data, and is explained by the principle of parsimony: do not associate a reinforcer with several stimuli, when one is enough to explain the association.

Looking at figure 4, presentation of  $S_1$  alone does give a lower response than that of  $S_0$ ; the reason the response is not zero is due to the ever present thalamic input  $Th$ , that picks up half of the reinforcement during the initial association with  $S_0$ . As

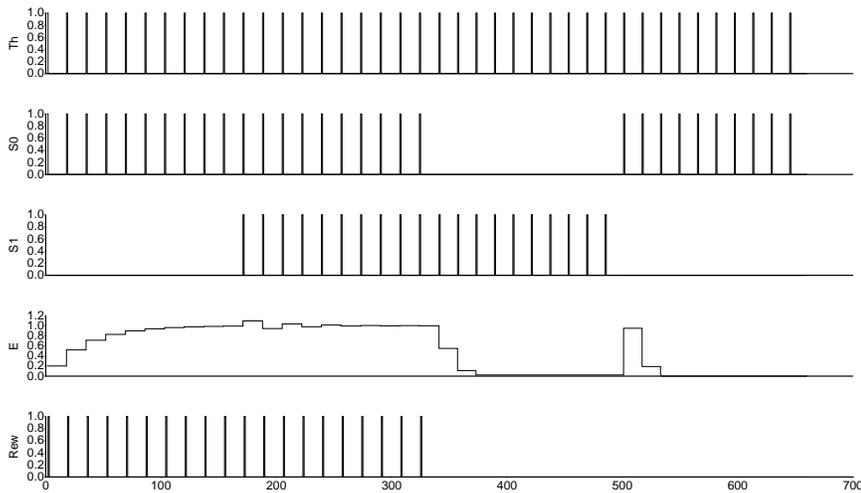


Figure 4: The result of a blocking experiment.  $S_0$  and  $S_1$  are stimulus inputs,  $Rew$  is the reward and  $E$  is the output from the model.  $Th$  is the thalamic input. As in the acquisition simulation,  $\alpha$  and  $\beta$  are set at 0.2 and 0.8, respectively.

the figure shows, the first response – once  $S_0$  and  $Rew$  is no longer present – is at 0.5, which is the response level of  $Th$ .

When the system is presented with  $S_1$  and  $Th$  in the absence of a reward, the orbitofrontal part will learn to inhibit a response through the connection weight for  $S_0$ , even though it is  $Th$  that is really responsible for the response.

The last part of this simulation shows that the response to  $S_0$  has been unaffected by the testing phase of the blocking experiment. Of course, the response is again quickly inhibited due to the absence of a reinforcer.

### 4.3 Conditioned Inhibition

In a conditioned inhibition schedule, the aim is to show that inhibition is an active process, not merely a decrease in associative strength. A stimulus can be given inhibitory properties, that can actively inhibit the response of other stimuli. Again, there is experimental evidence that this effect is common in animals (Mackintosh, 1983).

The schedule for conditioned inhibition is somewhat involved. We want to establish an inhibitory association with a stimulus, then test it with another stimulus that already has an association with the reinforcer. Creating an inhibitory association can be done by explicitly omitting the expected reinforcer whenever the stimulus is present.

First, associate  $S_2$  with the reinforcer; this is the stimulus that will be used for testing. Next, asso-

ciate  $S_0$  with reinforcement, and  $S_0+S_1$  with no reinforcement. This should give  $S_1$  inhibitory properties as  $S_0$  predicts the presence of the reinforcer. To test the result,  $S_2$  (the test stimulus) and  $S_1$  are presented together, and should give little or no response. Last, we also present  $S_2$  alone, to show that it has not been affected by the inhibitory associating stage.

The results are as expected:  $S_1$  and  $S_2$  give only a small, immediately decaying response, while  $S_2$  alone gives a satisfactory response. This result is due to the fact that the orbitofrontal part actively learns to inhibit responses in the presence of  $S_1$ , rather than the amygdaloid part unlearning anything.

## 5. Discussion

The model presented in this paper has several attractive characteristics. It can handle most common associative learning experiments and it is easy to implement and use as part of a larger system.

Currently, the model is not a complete learning system. As it is an emotional evaluator of stimuli, it needs several components to handle “real” learning tasks. The two most important missing parts are a context model and some form of motor learning system that can use the output of this model.

The context system (thought to be residing in the hippocampus) is currently being investigated and we already have a model up and running; see (Balkenius and Morén, 2000a) for details. A proper

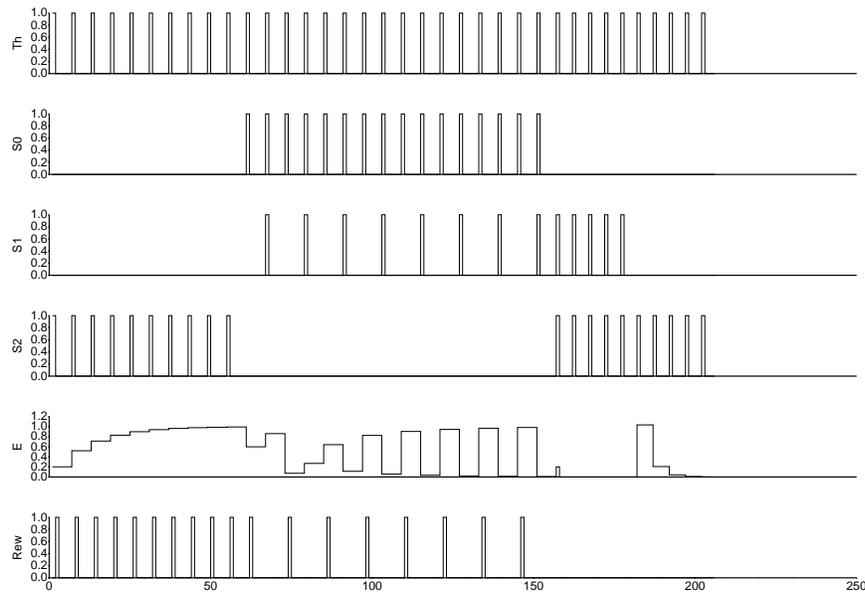


Figure 5: The result of conditioned inhibition on the model.  $S_0$ ,  $S_1$  and  $S_2$  are the stimulus inputs,  $Rew$  is the reward,  $Th$  is the thalamic input and  $E$  is the output from the model.  $\alpha$  is 0.2 and  $\beta$  is 0.8 .

context model will in the simplest case be replacing the orbitofrontal inputs from the sensory cortex in the model with context representations, and will hopefully work “out of the box” together with this model. This will allow the system to achieve real context-dependent learning, and will resolve any problems with the thalamic connection currently present in this model.

A motor learning system would use the output of this model as a reinforcing signal for learning motor sequences. There are any number of models capable of being used, such as Q-learning (Watkins and Dayan, 1992), HQ (Morén, 1998) and TD( $\lambda$ ) (Sutton and Barto, 1998), if one wanted to get a complete system.

A recurring question is the origin of the reinforcing signal. As the model currently stands, it appears from nowhere (or rather, from the simulator) to control the behaviour of the model. We believe this signal is the result of a reaction to the presentation of a primary stimulus or an emotionally charged stimulus. First order conditioning would be the association with a primary stimulus (a stimulus that has an intrinsic emotional charge), while second order conditioning would be the association with previously learned stimuli. This also has the implication that an inhibition of an emotional reaction also prevents the learning of new associations.

The presence of the thalamic input seems somewhat difficult to justify in this model, as it somewhat confuses the results of the simulations in this paper. As

it does not project into the orbitofrontal part, it can not be directly inhibited, and will confuse the inhibitory system. As the thalamic stimulus will not occur on its own (unless the sensory cortex is damaged), it will always be possible to at least partially inhibit its expression. Also, any effect the connection will have on an undamaged system will be less noticeable when there are many cortical inputs.

There seems little doubt that this path does exist in animal brains and has approximately the same function as it has in the model.

Two factors speak in favor of this interconnection: speed and fault tolerance. First, picking a rough estimate from the thalamus directly is a faster data path than going through the sensory cortex, allowing the system to react faster to broad classes of stimuli. Second, this path allows some emotional learning to proceed even if parts of the sensory cortex are damaged.

A difference between the thalamic interconnection in this model and in physical brains is that this interconnection is a single stimulus in the model, whereas the physical systems have a number of connections. These probably represent an early, rough classification and analysis of sensory stimuli performed by the thalamus.

## Conclusion

In this paper, we have presented a model of emotional learning in the amygdala. We have made use of neurophysiological data to guide the design of the model and have tested the model in three simulated experiments.

The model is divided into two main parts: one corresponding to the amygdala and one corresponding to the orbitofrontal cortex. Their internal structures are not closely modelled from their biological counterparts, but share their larger-scale organisation with the real system.

The simulated experiments are three classical conditioning experiments: Acquisition-extinction-reacquisition, blocking and conditioned inhibition. All three simulations performed well, indicating that the model at least has the basic features needed for associative learning.

## Acknowledgements

We are very grateful for the support from The Swedish Council for Research in the Humanities and Social Sciences (HSFR) and the Swedish Foundation for Strategic Research (SSF).

Other papers pertaining to this subject can be found at:

<http://www.lu.se/People/Christian.Balkenius/Conditioning.Habituation/index.html>

## References

- Balkenius, C. and Morén, J. (2000a). A computational model of context processing. Submitted.
- Balkenius, C. and Morén, J. (2000b). Emotional learning: A computational model of the amygdala. *Cybernetics and Systems*. In press.
- Dias, R., Robbins, T., and Roberts, A. (1997). Dissociable forms of inhibitory control within prefrontal cortex with an analog of the Wisconsin card sort test: Restriction to novel situations and independence from “on-line” processing. *The Journal of Neuroscience*, 17(23):9285–9297.
- Fanselow, M. S. and LeDoux, J. E. (1999). Why we think plasticity underlying pavlovian fear conditioning occurs in the basolateral amygdala. *Neuron*, 23:229–232.

- Gray, J. (1995). A model of the limbic system and basal ganglia: Application to anxiety and schizophrenia. In Gazzaniga, M., (Ed.), *The Cognitive Neurosciences*, pages 1165–1176, Cambridge, MA. MIT Press.
- Houk, J., Davis, J., and Beiser, D. (1995). *Models of information processing in the basal ganglia*. MIT Press, Cambridge, MA.
- Kolb, B. and Wishaw, I. (1990). *Fundamentals of human neuropsychology*. MIT Press, New York.
- LeDoux, J. (1995). In search of an emotional system in the brain: leaping from fear to emotion and consciousness. In Gazzaniga, M., (Ed.), *The Cognitive Neurosciences*, pages 1049–1061, Cambridge, MA. MIT Press.
- Mackintosh, N. (1983). *Conditioning and associative learning*. Oxford University Press, Oxford.
- McDonald, A. (1992). Cell types and intrinsic connections of the amygdala. In Aggleton, J., (Ed.), *The Amygdala: Neurobiological Aspects of Emotion, Memory and Mental Dysfunction*, pages 67–96, New York. Wiley-Liss.
- McIntosh, A. and Gonzalez-Lima, F. (1998). Large-scale functional connectivity in associative learning: Interrelations of the rat auditory, visual, and limbic systems. *Journal of Neurophysiology*, (80):3148–3162.
- Morén, J. (1998). Dynamic action sequences in reinforcement learning. In Pfeifer, R., Blumberg, B., Meyer, J., and Wilson, S., (Eds.), *From animals to animats 5: Proceedings of the Fifth International Conference on Simulation of Adaptive Behavior*, pages 366–271, Cambridge, MA. MIT Press.
- Morén, J. and Balkenius, C. (2000). Reflections on emotion. To appear in EMCSR2000, Vienna.
- Mowrer, O. (1973). *Learning theory and behavior*. Wiley, New York.
- Rescorla, R. and Wagner, A. (1972). A theory of pavlovian conditioning: variations in the effectiveness of reinforcement and nonreinforcement. In Black, A. and Prokasy, W., (Eds.), *Classical conditioning II: current research and theory*, pages 64–99, New York. Appleton-Century-Crofts.

- Rolls, E. (1986). A theory of emotion, and its application to understanding the neural basis of emotion. In Oomava, Y., (Ed.), *Neural and Chemical Control*, pages 325–344, Tokyo. Japan Scientific Societies Press.
- Rolls, E. (1992). Neurophysiology and functions of the primate amygdala. In Aggleton, J., (Ed.), *The Amygdala: Neurobiological Aspects of Emotion, Memory and Mental Dysfunction*, pages 143–165, New York. Wiley-Liss.
- Rolls, E. (1995). A theory of emotion and consciousness, and its application to understanding the neural basis of emotion. In Gazzaniga, M., (Ed.), *The Cognitive Neurosciences*, pages 1091–1106, Cambridge, MA. MIT Press.
- Shimamura, A. (1995). Memory and frontal lobe function. In Gazzaniga, M., (Ed.), *The Cognitive Neurosciences*, pages 803–813, Cambridge, MA. MIT Press.
- Sutton, R. and Barto, A. (1998). *Reinforcement Learning : An Introduction*. MIT Press, New York.
- Watkins, C. and Dayan, P. (1992). Q-learning. *Machine Learning*, 8:279–292.