

Simulation and Analysis of Adaptive Agents: an Integrative Modelling Approach

Tibor Bosse¹, Catholijn M. Jonker², and Jan Treur¹

¹Vrije Universiteit Amsterdam, Department of Artificial Intelligence

Tel. +31.20.5987750

{tbosse, treur}@cs.vu.nl

<http://www.cs.vu.nl/~{tbosse, treur}>

²Division Cognitive Engineering,

Nijmegen Institute for Cognition and Information, University of Nijmegen

C.Jonker@nici.ru.nl

<http://www.nici.ru.nl/~catholj>

Abstract. Agent-based simulation methods are a relatively new way to address complex systems. Usually the idea is that the agents used are rather simple, and the complexity and adaptivity of such a system are modelled by the interaction between these agents. However, another way to exploit agent-based simulation methods is by use of agents that themselves also have certain forms of learning or adaptation. In order to simulate adaptive agents with abilities matching those of their real-world biological or societal counterparts, a natural approach is to incorporate certain adaptation mechanisms such as classical conditioning into agent models. Existing models for adaptation mechanisms are usually based on quantitative, numerical methods, and more in particular, differential equations. Since agent-based simulation is usually based on qualitative, logical languages, these quantitative models are often not directly appropriate as an input in the context of agent-based simulation. To deal with this problem, this paper puts forward an integrative approach to simulate and analyse the dynamics of a conditioning process of an adaptive agent, integrating quantitative, numerical and qualitative, logical aspects within one expressive temporal specification language. To obtain a simulation model, an executable sublanguage of this language is used to specify the agent's adaptation mechanism in detail. For analysis and validation, in the proposed approach both properties characterising the externally observable adaptive behaviour and properties characterising the dynamics of internal intermediate states have been identified, formally specified and automatically checked on the generated simulation traces. As part of the latter, an approach to (formally) specify and check representational relations for intermediate, internal agent states is put forward. This enables verification of whether the representational content of an intermediate state a modeller has in mind indeed is in accordance with the agent model's internal dynamics. For a biological agent with known neural mechanisms, such as *Aplysia*, the modelling approach incorporates high-level modelling of neural states occurring as intermediate states and relates them to their representational content specification. This provides the possibility to validate not only the resulting observable behaviour of a simulation model against the observable behaviour of the agent in the real world, but also the intermediate states of the agent in the model against the intermediate states of the agent in the world.

1 Introduction

Agent-based modelling techniques are often used to model and simulate (natural or artificial) agent systems that have to deal with dynamic and uncertain environments. Usually agent-based simulation methods use agents that are rather simple; the complexity and adaptivity of such a system are modelled by the interaction between these agents. However, another way to exploit agent-based simulation methods is by use of agents that themselves also have certain forms of learning or adaptation. Therefore, an important challenge for the area of agent-based modelling is the notion of adaptive agent. An example of a basic mechanism for adaptation that can be found in many organisms is classical conditioning [19]. In order to create agent-based simulations with adaptive abilities matching those of their biological counterparts, a natural approach is to integrate such adaptation mechanisms into agent-based simulation models, e.g., [1].

In the literature adaptation mechanisms such as classical conditioning are usually described and analysed informally. If formalisation is used, this is often based on mathematical models using differential equations, e.g., Dynamical Systems Theory (DST) [20]. In contrast, agent-based simulation models traditionally make use of qualitative, logical languages, such as Golog [21], MetatheM [7], or 3APL [6]. Most of these languages are appropriate for expressing qualitative relations, but less suitable to work with more complex numerical structures as, for example, in differential equations. Therefore, integrating such mathematical models within the design of agent-based simulation models is difficult. To achieve this integration, it is necessary to bridge the gap between the quantitative nature of existing adaptation models and the type of languages typically used in agent-based simulation.

In the area of simulation, a formalised model is used to compute the simulation steps. Languages and software environments are available to support this modelling process. Validation of a model is usually not formally supported; it is considered a different issue. Often validation is done informally, by hand (or eye), based on comparison of a simulation trace with an empirical trace. In addition, sometimes specific (e.g., statistical) techniques are used to support certain aspects of validation. Usually in the domain that is modelled, *global properties* that should hold for the behaviour of a simulation model can be identified. As the languages used to specify a simulation model are directed to *local properties* (the steps between successive states), such global properties cannot be formalised in these languages. To obtain more support, also for validation of a simulation model, it is needed to integrate the modelling of such global properties in a formal manner as well, so that their specification and automated checking on simulation traces also can be supported by the modelling environment.

In accordance with the findings mentioned above, this paper introduces an approach for simulation and analysis of adaptive agent behaviour and underlying mechanisms that is integrative in two ways:

- (1) It combines in one modelling framework both *qualitative, logical* and *quantitative, numerical* aspects
- (2) It enables modelling dynamics both at a *local level* (internal mechanisms of the agent) and at a *global level* (externally observable agent behaviour, and representation relationships between internal and external states)

Modelling dynamics at a local level concerns expressing temporal relationships between pairs of successive states, such as described by direct causal relations, or, for example, by the basic steps within an *adaptation mechanism*. A difference or differential equation is an example of a local level specification of dynamics. From a local perspective, the dynamics of the actual underlying (e.g., neural) mechanisms that play a role in the real world can be investigated. Local level specifications are the basis for the computation steps for a simulation model.

From the global perspective, more complex relationships over time can be used to model dynamics for adaptive agents. For example, the dynamics of *observed adaptive agent behaviour* can be analysed, i.e., how during a history of (learning) experiences, the behaviour is changing. For example, the performance of actions depending on a stimulus in the present and a certain training history (series of training stimuli in the past) can be modelled. This can take the form of a temporal relationship (an *input-output correlation*) involving a longer time duration and several agent input and output states over time.

Besides input-output correlations describing adaptive agent behaviour as just discussed, also from a global perspective *representation relations* for intermediate, internal agent states can be modelled. During modelling, for an internal or intermediate state of the agent as introduced in the model, often a modeller has in mind a certain representational content, i.e., how it relates to other concepts outside the agent. To take a simple example, it may be expected that the internal belief that a horse is nearby correlates to the actual presence of this horse. Such expected representation relations may or may not be inspired by knowledge of how the agent's adaptation mechanism is realised in Nature. The approach put forward includes ways to (formally) specify such representation relations and verify them against simulation traces, showing whether this representational content is in accordance with the agent model's internal dynamics. In this way the modelling approach can also address the issue of realism of internal or intermediate states in a simulated agent. For example, for an adaptive biological agent with known neural mechanisms, such as *Aplysia*, the modelling approach incorporates the modelling of neural states occurring as intermediate states and relating them to other states in the world according to their representational relation specification. This provides the possibility to validate not only the

resulting observable behaviour of an agent simulation model against the observable behaviour of the agent in the real world, but also the internal, intermediate states of the agent in the model against the internal, intermediate states of the agent in the world. Thus it can be verified to what extent the model satisfies *internal realism* in addition to external realism.

As both the adaptation mechanism and the externally observable behaviour are modelled in the form of temporal relationships, within the modelling approach it is also possible to logically relate the dynamics of internal agent models involving a (neural) adaptation mechanisms to the model for the dynamics of the externally observable adaptive behaviour. Such *interlevel relations* often take the hierarchical form of an AND-tree (or a number of them), with the most global property at the top (root) and the most local at the leaves. Such a hierarchical structure can be useful in the analysis of, in case, why a global property fails on a certain simulation trace. By going down in the tree and at each level checking the properties under the failing node, finally the leaf or leaves that fail(s) can be found, thus pinpointing the (local) cause of the failure. This can be useful in debugging a model, but also in the analysis of the circumstances under which a model will function well and under which not, and the reasons why.

If the actual underlying neural mechanisms are included in the analysis of adaptive behaviour, the sea hare *Aplysia* is an appropriate species to study, since its neural mechanisms have been well-investigated; cf. [10]. In this paper it will be shown how the proposed modelling approach for adaptive agents can be used to simulate and analyse both *Aplysia's* adaptive behaviour and the underlying neural mechanisms.

An overview of the paper is as follows. In Section 2 the high-level modelling approach is briefly introduced. Section 3 introduces the case study and the state properties for this case study. In Section 4 the executable local dynamic properties describing basic mechanisms for the case study are presented; simulations on the basis of these local dynamic properties are discussed in Section 5. In Section 6 the interlevel relations between dynamic properties of the externally observable behaviour and the local properties describing the internal mechanisms are discussed. In Section 7 different approaches to representational content are explored and formalised. Section 8 discusses how all these dynamic properties have been checked against the simulation traces. Section 9 is a discussion.

2 Modelling Approach

To formally specify dynamic properties that express criteria for representational content from a temporal perspective an expressive language is needed. Dynamics will be described in the next section as evolution of *states* over time. The notion of state as used here is characterised on the basis of an ontology defining a set of *state properties* that do or do not hold at a certain point in

time. Examples of state properties are ‘the agent is hungry’, ‘the agent observes rain’, ‘the agent has internal state s ’, or ‘the environmental temperature is 7°C ’. Real value assignments to variables are also considered as possible state property descriptions. For example, in a quantitative modelling approach (such as [20]), based on variables x_1, x_2, x_3, x_4 , that are related by differential equations over time, value assignments such as

$$\begin{aligned} x_1 &\leftarrow 0.06 \\ x_2 &\leftarrow 1.84 \\ x_3 &\leftarrow 3.36 \\ x_4 &\leftarrow -0.27 \end{aligned}$$

are considered state descriptions. State properties are described by ontologies that specify the concepts used.

Based on such state properties, *dynamic properties* can be formulated that relate a state at one point in time to one or more states at other points in time. A simple example is the following dynamic property specification:

‘at any point in time t_1 if the agent observes rain at t_1 , then there exists a point in time t_2 after t_1 such that at t_2 the agent has internal state property s ’

Here, for example, s can be viewed as a sensory representation of the rain. To express such dynamic properties, and other, more sophisticated ones, the temporal trace language TTL is used [14]. This language can be classified as a reified predicate-logic based temporal language; see [8], [9].

Within this language, explicit references can be made to time points and traces. Here a fixed *time frame* \mathcal{T} is assumed which is linearly ordered. Depending on the application, it may be continuous (e.g., the real numbers), or discrete (e.g., the set of integers or natural numbers or a finite initial segment of the natural numbers), or any other form, as long as it has a linear ordering. Moreover, a *trace or trajectory* over an ontology Ont is a time-indexed sequence of states over Ont . The sorted predicate logic temporal trace language TTL is built on atoms referring to, e.g., traces, time and state properties. State properties are denoted by terms in the language TTL. For example,

in the internal state of agent A in trace γ at time t property s holds

is formalised by $\text{state}(\gamma, t, \text{internal}(A)) \models s$. Here \models is a predicate symbol in the language, usually used in infix notation, which is comparable to the `Holds`-predicate in situation calculus. Dynamic properties are expressed by temporal statements built using the usual logical connectives and quantification (for example, over traces, time and state properties).

To be able to perform some simulation experiments, a simpler temporal language has been used to specify executable models in a declarative manner. This language (the *leads to* language [4]) enables to model direct temporal dependencies between two state properties in successive states. This executable format is defined as follows. Let α and β be state properties of the form ‘conjunction of atoms or negations of atoms’, and e, f, g, h non-negative real numbers. Then the notation $\alpha \rightarrow_{e, f, g, h} \beta$, means:

If state property α holds for a certain time interval with duration g

then after some delay (between e and f) state property β will hold for a certain time interval of length h .

For a precise definition of the *leads to* format in terms of the language TTL, see [14]. A specification of dynamic properties in *leads to* format has as advantages that it is executable and that it can often easily be depicted graphically. The *leads to* format has shown its value especially when temporal or causal relations in the (continuous) physical world are modelled and simulated in an abstract, non-discrete manner; for example, the intracellular chemistry of *E. coli* [12].

3 The Aplysia Case Study

To illustrate the proposed approach for modelling and simulation of adaptive agents, it is applied in a case study. As the topic of the case study, the sea hare *Aplysia* was chosen. The motivation for this choice is two-fold. First, *Aplysia* is a clear example of an adaptive agent. Second, the internal neural mechanisms of *Aplysia* are relatively simple, and therefore well understood. This enables the modeller to (formally) describe *Aplysia*’s behaviour both from an *internal* perspective (i.e., at a *local level*, considering neural mechanisms of the agent) and from an *external* perspective (i.e., at a *global level*, considering externally observable agent behaviour). As a result, both *interlevel relations* (see Section 6) and *representation relations* (see Section 7) can be established between both types of descriptions. First, in Section 3.1, *Aplysia*’s behaviour will be described from an external perspective. In Section 3.2, *Aplysia*’s behaviour will be described from an internal perspective.

3.1 External Perspective

Aplysia is a sea hare that is often used to do experiments. It is able to learn on the basis of classical conditioning. In this section, a simplified description is given of this learning behaviour (viewed from an external perspective), based on [10], pp. 155-156.

Behaviour before learning phase

Initially the following behaviour is shown:

- a tail shock leads to a response (contraction)
- a light touch on its siphon is insufficient to trigger such a response

Learning phase

Now suppose the following experimental protocol is undertaken. In each trial the subject is touched lightly on its siphon and then, shocked on its tail (as a consequence it responds).

Behaviour after a learning phase

It turns out that after a number of trials (three in the example) the behaviour has changed:

- the animal also responds (contracts) on a siphon touch.

Note that, to characterise behaviour, there is a difference between the *learned* behaviour (which is simply an *adapted* stimulus-response behaviour) and the *learning* behaviour, which is a form of *adaptive* behaviour, no stimulus-response behaviour. To specify such behaviours the following sensor and effector states are used: tail_shock, siphon_touch, contraction. In terms of these state properties the following global dynamic properties can be specified in *leads to* format:

GP1 (Contraction Upon Tail Shock)

At any point in time t,

if a tail shock occurs then it will contract

Formally:

tail_shock $\rightarrow_{e,f,g,h}$ contraction

GP2 (Contraction Upon Siphon Touch)

At any point in time t,

if a siphon touch occurs then it will contract

Formally:

siphon_touch $\rightarrow_{e,f,g,h}$ contraction

The latter property specifies the behaviour that is the *result* of the learning process. However, the behaviour shown by the learning *process* itself is not expressed here, and as this process involves more complex temporal relationships, is even not expressable in *leads to* format. However, it is expressable in TTL format:

GP3 (Learning to Contract Upon Siphon Touch)

At any point in time t,

if a siphon touch occurs

and at three different earlier time points t1, t2, t3,

a siphon touch occurred, directly followed by a tail shock

then it will contract

Formally:

$$\begin{aligned}
& \forall \gamma \forall t \text{ state}(\gamma, t) \models \text{siphon_touch} \ \& \\
& \exists t1, t2, t3, u1, u2, u3 \ t1 < u1 < t2 < u2 < t3 < u3 < t \ \& \\
& \text{state}(\gamma, t1) \models \text{siphon_touch} \ \& \text{state}(\gamma, u1) \models \text{tail_shock} \ \& \\
& \text{state}(\gamma, t2) \models \text{siphon_touch} \ \& \text{state}(\gamma, u2) \models \text{tail_shock} \ \& \\
& \text{state}(\gamma, t3) \models \text{siphon_touch} \ \& \text{state}(\gamma, u3) \models \text{tail_shock} \\
& \Rightarrow \exists t' \geq t \ \text{state}(\gamma, t') \models \text{contraction}
\end{aligned}$$

As can be seen, the temporal complexity of the learning behaviour specification is much higher than that of the learned behaviour.

3.2 Internal Perspective

This section describes *Aplysia*'s behaviour from an internal perspective. The internal neural mechanism for *Aplysia*'s conditioning can be depicted as in Figure 1; cf. [10].

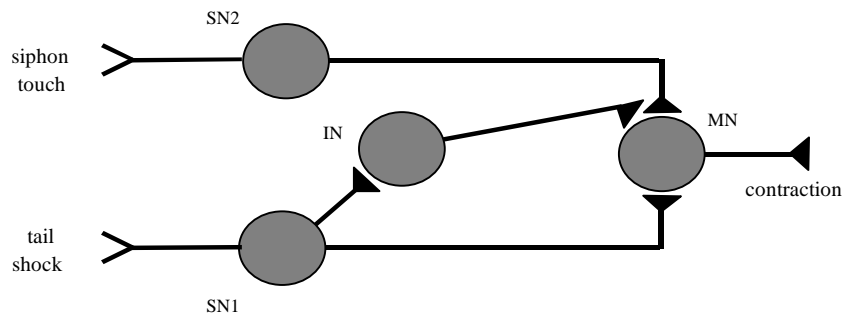


Fig. 1. Neural mechanisms

A tail shock activates a sensory neuron SN1. Activation of this neuron SN1 activates the motoneuron MN; activation of MN makes the sea hare move. A siphon touch activates the sensory neuron SN2. Activation of this sensory neuron SN2 normally does not have sufficient impact on MN to activate MN. After learning, activation of SN2 has sufficient impact to activate MN. In addition, activation of SN1 also leads to activation of the intermediary neuron IN. If both SN2 and IN are activated simultaneously, this changes the synapse between SN2 and MN: it causes this synapse to produce more neurotransmitter if SN2 is activated. After a number of times this leads to the situation that also activation of SN2 yields activation of MN.

To model the example the following internal state properties are used:

- SN1 sensory neuron 1 is activated
- SN2 sensory neuron 2 is activated
- IN intermediary neuron IN is activated

MN motoneuron MN is activated
S(r) the synapse between SN2 and MN is able to produce an amount r of
 neurotransmitter

The dynamics of these internal state properties involve temporal *leads to* relationships, which are analysed in more detail in the next section.

4 Local Dynamic Properties

To model the internal dynamics of the example, the following local properties (in *leads to* format) are considered. They describe the basic steps or mechanisms of the process.

LP1 (SN1 Activation)

At any point in time,

if a tail shock occurs then SN1 will be activated

Formally:

$$\text{tail_shock} \rightarrow_{e,f,g,h} \text{SN1}$$

LP2 (SN2 Activation)

At any point in time,

if a siphon touch occurs then SN2 will be activated

Formally:

$$\text{siphon_touch} \rightarrow_{e,f,g,h} \text{SN2}$$

LP3 (IN and MN Activation by SN1)

At any point in time,

if activation of SN1 occurs then IN and MN will be activated

Formally:

$$\text{SN1} \rightarrow_{e,f,g,h} \text{IN} \wedge \text{MN}$$

LP4 (Synaps Adaptation)

At any point in time,

if activation of SN2 occurs

and activation of IN occurs

and the synaps has strenght r with $r < 4$

then the synaps will have strenght r+1

Formally:

$$\text{S}(r) \wedge \text{SN2} \wedge \text{IN} \wedge r < 4 \rightarrow_{e,f,g,h} \text{S}(r+1)$$

LP5 (MN Activation by SN2)

At any point in time,
 if activation of SN2 occurs
 and the synaps has strenght 4
 then MN will be activated

Formally:

$$S(4) \wedge SN2 \xrightarrow{e,f,g,h} MN$$

LP6 (Contraction by MN)

At any point in time,
 if activation of MN occurs then it will contract

Formally:

$$MN \xrightarrow{e,f,g,h} \text{contraction}$$

LP7 (Synaps State Persistence)

At any point in time,
 if the synaps has strenght r with $r < 4$
 and the synaps has not strenght r+1
 then the synaps will have strenght r

Formally:

$$S(r) \wedge \text{not } S(r+1) \wedge r < 4 \xrightarrow{e,f,g,h} S(r)$$

LP8 (Synaps State Persistence)

At any point in time,
 if the synaps has strenght 4 then the synaps will have strenght 4

Formally:

$$S(4) \xrightarrow{e,f,g,h} S(4)$$

LP9 (Initialisation)

At the start, the synaps has strenght 1

Formally:

$$\text{start} \xrightarrow{e,f,g,h} S(1)$$

In Figure 2 an overview of these properties is given in a graphical form. Here, the circles denote state properties and the arrows denote dynamic properties.

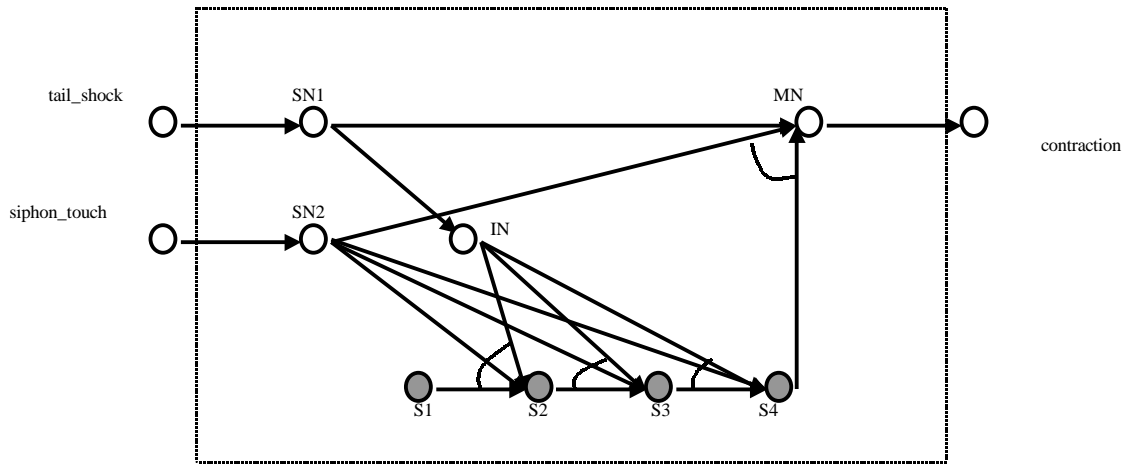


Fig. 2. Overview of the basic dynamics of the simulation model

Note that this model is based on a number of simplifications. For example, it is assumed that after exactly 4 steps the strength of the synapse between SN2 and MN is maximal, and that there is no extinction. However, since our modelling approach supports the use of quantitative concepts (such as real numbers and mathematical operations), it is easy to incorporate such features in the model. A rather straightforward way to do this is by replacing LP4 through LP8 by the following local properties. Here, β indicates the learning rate, K indicates the maximal strength of the synapse between SN2 and MN (e.g., 4), ϵ indicates the extinction rate, and t indicates the minimum threshold of S needed to have SN2 influence MN. For all values, real numbers can be used.

LP4 (Synaps Adaptation)

At any point in time,

if activation of SN2 occurs

and activation of IN occurs

and the synaps has strenght r

then the synaps will have strenght $\beta*(K-r)+(r*\epsilon)$

Formally:

$$S(r) \wedge SN2 \wedge IN \xrightarrow{e,t,g,h} S(\beta*(K-r)+(r*\epsilon))$$

LP5 (MN Activation by SN2)

At any point in time,

if activation of SN2 occurs
 and the synaps has strenght $r > t$
 then MN will be activated

Formally:

$$S(r) \wedge SN2 \wedge r > t \rightarrow_{e,f,g,h} MN$$

LP7 (Synaps State Decay)

At any point in time,

if the synaps has strenght r
 and SN2 is not activated
 then the synaps will have strenght $r * \epsilon$

Formally:

$$S(r) \wedge \text{not } SN2 \rightarrow_{e,f,g,h} S(r * \epsilon)$$

LP8 (Synaps State Decay)

At any point in time,

if the synaps has strenght r
 and IN is not activated
 then the synaps will have strenght $r * \epsilon$

Formally:

$$S(r) \wedge \text{not } IN \rightarrow_{e,f,g,h} S(r * \epsilon)$$

Another extension to the model would be to introduce real-valued arguments for the state properties SN1, SN2, IN and MN as well, indicating the strength of their activation. This would allow the model to distinguish between, for example, tail shocks of different strengths. Although these extensions are relatively easy to perform, for reasons of presentation in the remainder of this paper the simplified model is used.

5 Simulation

As mentioned in the Introduction, local level specifications are the basis for the computation steps for a simulation model. Thus, special software environments can be created to enable the simulation of local level specifications, as long as these are in an executable format. For the executable language *leads to*, such a software environment has indeed been built, see [4] for details. Based on an input consisting of dynamic properties in *leads to* format, this software

environment generates simulation traces. An example of such a trace can be seen in Figure 3. Here, time is on the horizontal axis, the state properties are on the vertical axis. A dark box on top of the line indicates that the property is true during that time period, and a lighter box below the line indicates that the property is false. This trace is based on all local properties identified in Section 4. In property LP1 and LP2 the values (0,0,1,3) have been chosen for the timing parameters e, f, g, and h. In all other properties, the values (0,0,1,1) have been chosen. As can be seen in Figure 3, at the beginning of the trace the organism has not performed any conditioning. The initial siphon touch it receives does lead to the activation of sensory neuron SN2, but the synapse between SN2 and motoneuron MN does not produce much neurotransmitter yet (indicated by internal state property S(1)). Thus, the activation of SN2 does not yield an activation of MN, and consequently no external action follows. In contrast, it is shown that a shock of the organism' s tail does initially lead to the external action of contraction. This can be seen in Figure 3 between time point 10 (when the tail shock occurs) and time point 13 (when the animal contracts). After that, the actual learning phase starts. This phase consists of a sequence of three trials where a siphon touch is immediately followed by a tail shock. As a result, the sensory neuron SN2 is activated at the same time as the intermediary neuron IN, which causes the synapse to change so that it can produce an increased amount of neurotransmitter each time SN2 is activated. Such a change in the synapse is indicated by a transition from one internal state property to another (first from S(1) to S(2), then to S(3), and finally to S(4)). As soon as internal state property S(4) holds (see time point 44), the conditioning process has been performed successfully. From that moment, *Aplysia*' s behaviour has changed: it also contracts on a siphon touch.

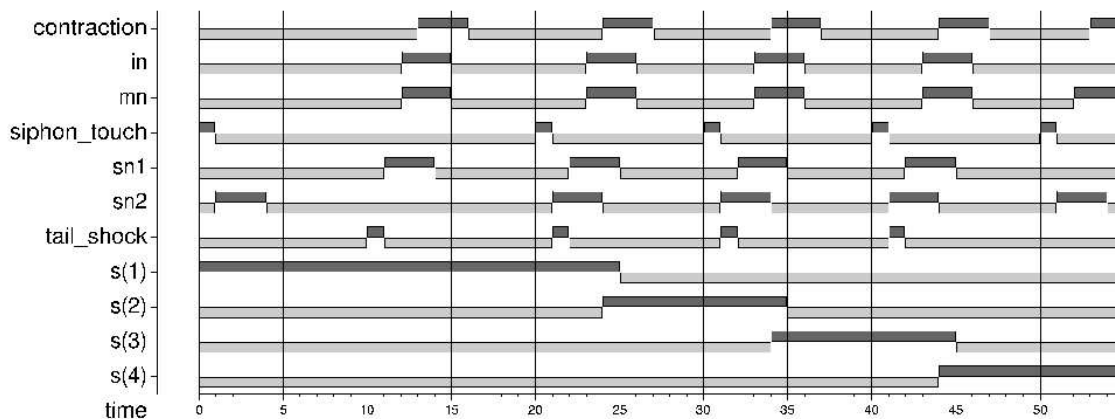


Fig. 3. Example simulation trace

For the purposes of this example, the amount of trials is kept low (three). However, similar experiments have been performed with a case of 1000 learning steps. Since the abstract way of modelling used for the simulation is not computationally expensive, also these simulations took no more than 90 seconds. In addition, our simulation approach has possibilities to incorporate real numbers in state properties, and to perform complex mathematical operations with these numbers. This makes it more expressive than more traditional forms of temporal logic.

6 Interlevel Relations

In the previous sections, both the internal (neural) adaptation mechanism and the externally observable behaviour of *Aplysia* were modelled in the form of temporal relationships. Within the presented modelling approach, this implies that it is also possible to logically relate the dynamics of both models. This section outlines these *interlevel* connections between dynamic properties at different levels. It will be shown how the description at the level of the neurological mechanisms (the local dynamic properties LP1 through LP9) can be logically related to the description at the level of the overall behaviour (the global property GP3). This way, a formalisation is obtained of the (interlevel) reduction relation between the two levels. To be precise, this relation is described by the following implication:

$$(1) \quad \text{LP1 through LP9 \& CWA} \quad \Rightarrow \quad \text{GP3}$$

This equation states that the local properties together imply the global property GP3 (which expresses that experiencing the combination of a tail shock and a siphon touch three times results in a response to the siphon touch alone). Moreover, one additional property is introduced, i.e., CWA. This second-order property that is commonly known as the Closed World Assumption expresses that at any point in time a state property that is not implied by a specification to be true is taken to be false. Let \mathcal{T}_h be the set of all local properties LP1 through LP9, then the formalisation is:

Closed World Assumption (CWA)

$$\forall P \in \text{At}(\text{ONT}) \quad \forall \gamma \quad \forall t: \mathcal{T}_h \quad \text{state}(\gamma, t) \models P \Rightarrow \text{state}(\gamma, t) \models \text{not } P$$

The Closed World Assumption is needed to ensure that the intermediate results as indicated by the $S(r)$ state properties can only hold as a result of the local properties LP1 through LP9, and not because of some other (mysterious) cause.

Essential milestones in the proof of relationship (1) are that subsequently $S(1)$, $S(2)$, $S(3)$, and $S(4)$ will hold. These milestones can be seen as the result of a learning process. Therefore, an additional lemma is introduced. This lemma describes the effect of a learning step on the synapse, showing the increase of parameter r in state property $S(r)$, given that the siphon is

touched, directly followed by a tail shock. In this case study the effect we are interested in is already reached at $r=4$. The lemma can easily be adapted for more lengthy learning processes. Formally, the lemma is specified as:

M(g, h, r) Learning step

$$\begin{aligned}
& \forall \gamma \quad \forall t_1, t_2, u_1 \\
& \quad t_1 < u_1 < t_1 + g \ \& \ t_1 < t_2 < t_1 + g \ \& \ r < 4 \ \& \\
& \quad \forall t \quad [t_1 \leq t < t_1 + h \Rightarrow \text{state}(\gamma, t) \models \text{siphon_touch}] \ \& \\
& \quad \forall t \quad [u_1 \leq t < u_1 + h \Rightarrow \text{state}(\gamma, t) \models \text{tail_shock}] \ \& \\
& \quad \forall t \quad [t_2 \leq t < t_2 + h \Rightarrow \text{state}(\gamma, t) \models S(r)] \\
& \Rightarrow \exists t_3 \ [t_3 \geq t_2 \ \& \ \forall t \ [t_3 \leq t < t_3 + h \Rightarrow \text{state}(\gamma, t) \models S(r+1)]]
\end{aligned}$$

Property M(g,h,r) can be proved for $g=1$, $h=1$, and r varying from 1 to 4 from LP1, LP2, LP3, LP4, LP7, and CWA, taking (0, 0, 1, 3) as timing parameters in LP1 and LP2, and (0, 0, 1, 1) for the timing parameters of the other local properties.

$$(2) \quad \text{LP1} \ \& \ \text{LP2} \ \& \ \text{LP3} \ \& \ \text{LP4} \ \& \ \text{LP7} \ \& \ \text{CWA} \quad \Rightarrow \quad M(1, 1, r)$$

The introduction of property M(1,1,r) allows one to reduce relationship (1) to the following, simpler implication:

$$(3) \quad \text{LP2} \ \& \ \text{LP5} \ \& \ \text{LP6} \ \& \ \text{LP7} \ \& \ \text{LP8} \ \& \ \text{CWA} \ \& \ M(1,1,r) \quad \Rightarrow \quad \text{GP3}$$

Figure 4 provides a visualisation of relationship (1) through (3). The semantics of this tree is as follows: if a certain trace satisfies all lower-level properties connected to a certain higher-level property, then this trace also satisfies the higher-level property.

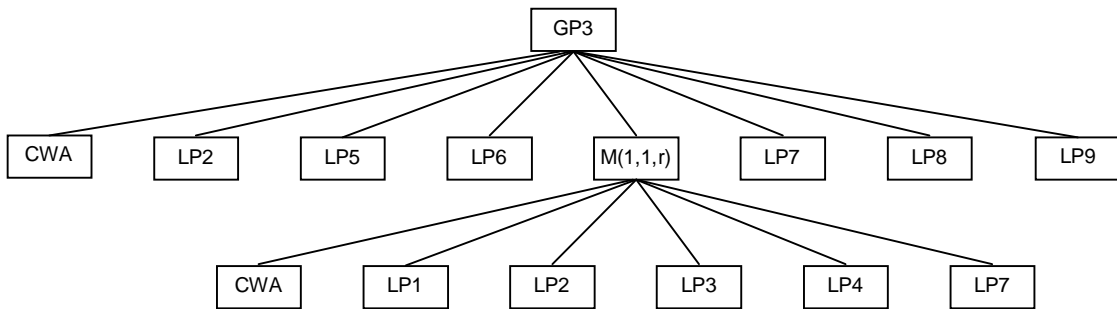


Fig. 4. Interlevel relationships for property GP3

The full proof of these relationships is a difficult issue, and is left out of this paper. Instead only a sketch of the proof is given, in which some initialisation issues are ignored. First a proof sketch of implication (2) is provided. For this proof, the crucial points are that the siphon touch

and the tail shock are coordinated in time such that SN2 (by application of LP2) exists long enough for it to co-exist with IN. Given LP7 and CWA it becomes clear that S(r) persists long enough for LP4 to have effect. The following sketch is illustrated by Figure 5. First, assume that:

$$\begin{aligned} & t1 < u1 < t1 + g \ \& \ t1 < t2 < t1 + g \ \& \ r < 4 \ \& \\ \forall t \quad & [t1 \leq t < t1 + 1 \Rightarrow \text{state}(\gamma, t) \models \text{siphon_touch}] \ \& \\ \forall t \quad & [u1 \leq t < u1 + 1 \Rightarrow \text{state}(\gamma, t) \models \text{tail_shock}] \ \& \\ \forall t \quad & [t2 \leq t < t2 + 1 \Rightarrow \text{state}(\gamma, t) \models S(r)] \end{aligned}$$

Then CWA can be applied to derive that:

$$\forall t \quad [t \leq t2 \Rightarrow \text{state}(\gamma, t) \models \text{not } S(r+1)]$$

In addition LP1 can be applied to derive:

$$\forall t \quad [u1 + 1 \leq t < u1 + 4 \Rightarrow \text{state}(\gamma, t) \models \text{SN1}]$$

Using this information, and LP3, the following is derived:

$$\forall t \quad [u1 + 2 \leq t < u1 + 5 \Rightarrow \text{state}(\gamma, t) \models \text{IN}]$$

Similarly, LP2 can be applied on the time duration of the siphon touch to derive:

$$\forall t \quad [t1 + 1 \leq t < t1 + 4 \Rightarrow \text{state}(\gamma, t) \models \text{SN2}]$$

In order to apply LP4 on the intersection interval of the periods during which both SN2 and IN hold (i.e., $[u1 + 2, t1 + 4>$, which has a duration > 1), it must be ascertained that S(r) also holds long enough (at least 1) in that interval. Since $t1 < u1 < t1 + g \ \& \ t1 < t2 < t1 + g$, the absolute difference $|t2 - u1| < 1$.

For the other case LP7 needs to be applicable to derive a persistence of S(r), to the extent that it overlaps long enough SN2 and IN. Given that LP4 cannot be applied in the time period $[t2, t2 + 1>$, and the fact that there is no local property that derives S(r+1) during that same interval, CWA is applicable. Even stronger, CWA is applicable until LP4 is applicable, deriving:

$$\forall t \quad [t2 \leq t < u1 + 3 \Rightarrow \text{state}(\gamma, t) \models S(r)]$$

Note that given that $t1 < u1 < t1 + g$ (with $g=1$), the interval $[t2, u1 + 3>$ overlaps with $[u1 + 2, t1 + 4>$ for the interval $[u1 + 2, u1 + 3>$. Therefore, LP4 can be applied on this interval to derive the result of M(1,1,r), thus proving relationship (2):

$$\forall t \quad [u1 + 3 \leq t < u1 + 4 \Rightarrow \text{state}(\gamma, t) \models S(r+1)]$$

As can be seen from the proof sketch, the timing issues make proofs complex. In Figure 5, the timing information is left out. The tree only gives insight into which local properties (or CWA) are applied on which state properties.

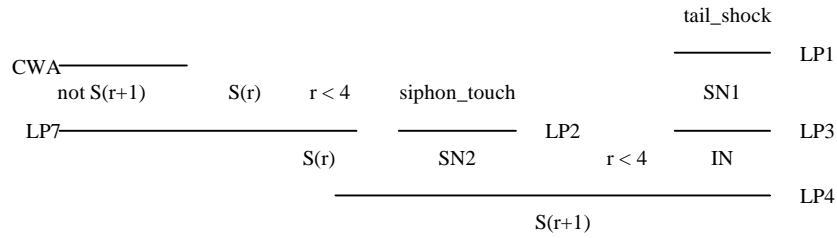


Fig. 5. Sketch of interlevel relationship (2)

The proof sketch for interlevel relationship (3) (as illustrated in Figure 6, again with all timing elements left out) takes all the siphon touches, tail shocks as given in the precondition of the implication in GP3 as hypotheses and shows how to derive a contraction of *Aplysia*. The initial assumption LP9 provides that $S(1)$ holds in the beginning. By the CWA it is possible to apply LP7 ensuring that $S(1)$ persists until the first sequence of siphon touch and tail shock have taken place and $M(1,1,1)$ can be applied. The pattern of applying CWA and LP7 to ensure persistence is repeated for every new occurrence of the sequence of siphon touch and tail shock, until $S(4)$ holds. Because of LP8, $S(4)$ persists until a new siphon touch occurs, and LP2 has been applied leading to SN2. The persisting of $S(4)$ and the existence of SN2 make LP5 applicable, leading to MN. Finally, the application of LP6 on MN leads to a contraction, thus completing the proof of (3). Combining the proofs of relationship (2) and (3) eventually results in the proof of (1).

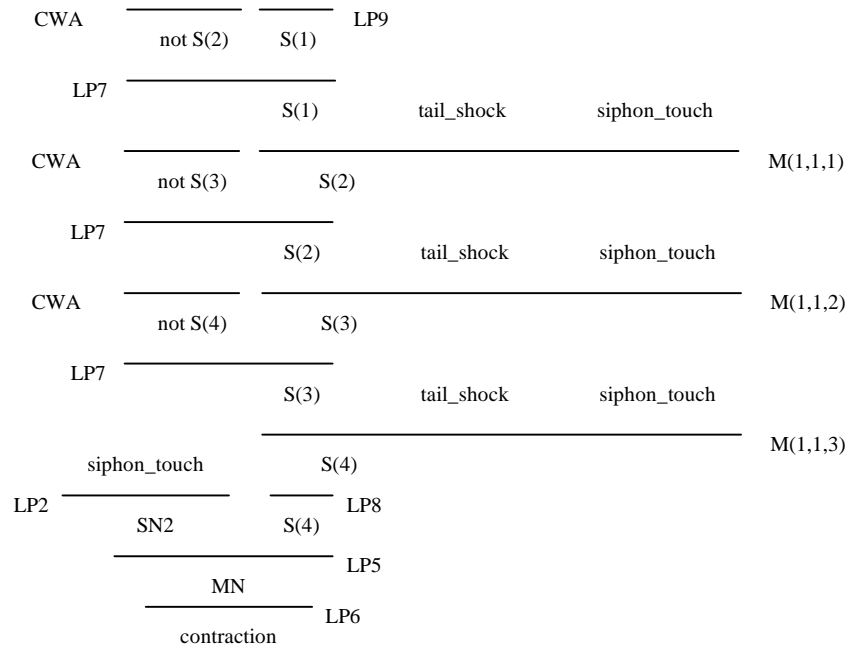


Fig. 6. Sketch of interlevel relationship (1)

Establishing interlevel relations such as represented in Figure 4 can be useful in the analysis of simulation traces. To illustrate this, assume that, in a given simulation trace, a certain global property (e.g., GP3) does not hold. Then by a refutation process it can be concluded that one of the lower level properties does not hold either (i.e., CWA, LP2, LP5, LP6, LP7, LP8, LP9 or M(1,1,r) does not hold). If, after checking these properties, it turns out that M(1,1,r) does not hold, then either CWA, LP1, LP2, LP3, LP4 or LP7 does not hold. Thus, by this example refutation analysis eventually the cause of the unsatisfactory behaviour can be reduced to the failure of a local property.

7 Representational Relations

In the literature on Philosophy of Mind different types of approaches to representational content of an internal state property have been put forward, for example the causal/correlational, interactionist and relational specification approach; cf. [2]; [16], pp. 191-193, 200-202. These approaches to representational content have in common that the occurrence of the internal state property at a specific point in time is related to the occurrence of other state properties, at the same or at different time points. The ‘other state properties’ can be of two types:

- A. *external world state properties*, independent of the agent
- B. the agent's sensor state and effector state properties, i.e. the agent's *interaction state properties* (interactivist approach)

Furthermore, the type of relationships can be (1) purely functional *one-to-one correspondences*, (e.g., the correlational approach), or (2) they can involve more *complex relationships* with a number of states at different points in time in the past or future, (e.g., the interactivist approach). So, four types of approaches to representational contents are distinguished, that can be indicated by codings such as A1, A2, and so on. Below, examples of such approaches are given.

According to the *causal/correlational approach* (see [16], pp. 191-193), the representational content of a certain internal state is given by a one-to-one correlation to another (in principle external) state property: type A1. For example, the internal belief that a horse is nearby is correlated to the actual presence of this horse, which is an external state property. Such an external state property may exist backward as well as forward in time. Hence, for the current example, in order to define the representational content of an internal state property, one should try if this can be related to a world state property that either existed in the past or will exist in the future. For example, the representational content for internal state property SN1 can be defined as world state property tail_shock, by looking backward in time. However, for some of the other internal state properties the representational content cannot be defined adequately according to the causal/correlational approach. In these cases, reference should not be made to one single state in the past or in the future, but to a temporal sequence of inputs or output state properties, which is not considered to adequately fit in the correlational approach. This shows that especially in cases where the agent learns from a number of trials extending over time, a classical approach to representational content is insufficient. Some authors even claim that it is a bad idea to aim for a notion of representation in such cases; e.g., [15], [23].

As an alternative, Bickhard's *Temporal-Interactivist approach* [2,9] relates the occurrence of internal state properties to sets of past and future interaction traces: type B. In this paper the focus is on the B2 type, which is the more advanced case.

The *Relational Specification approach* to representational content is based on a specification of how the occurrence of an internal state property relates to properties of states distant in space and time; cf. [16], pp. 200-202. In this paper it is used in conjunction with the temporal-interactivist approach. Thus, the representational content of a certain internal state can be defined by specifying a temporal relation of the internal state property to sensor and action states in the past and future. An overview for the content of all internal state properties of the case study, according to the temporal relational specification approach is given, in an informal

notation, in Table 1. Note that these relationships in fact are defined at a semantic level, not syntactically specified in a modelling language. Different interaction state properties, separated by commas, should be read as the temporal sequence of these states.

Table 1. Temporal-Interactivist Representation Relation (sketch)

Internal State Property	Content (backward)	Content (forward)
S(2)	siphon_touch, tail_shock	
S(3)	siphon_touch, tail_shock, siphon_touch, tail_shock	
S(4)	siphon_touch, tail_shock, siphon_touch, tail_shock, siphon_touch, tail_shock	any siphon_touch is followed by contraction

Table 2 and 3 describe the same information as Table 1, but this time syntactically, expressed by TTL formulae. The following abstractions are used to describe training periods:

$$\text{training_up_to}(\gamma, t1, u1, 1) \equiv$$

$$u1 = t1 + 1 \ \& \ \text{state}(\gamma, t1) \models \text{siphon_touch} \ \& \ \text{state}(\gamma, u1) \models \text{tail_shock}$$

$$\text{training_up_to}(\gamma, t1, u2, 2) \equiv$$

$$\exists u1, t2 [u1 < t2 \ \& \ u2 = t2 + 1]$$

$$\text{training_up_to}(\gamma, t1, u1, 1) \ \& \ \text{state}(\gamma, t2) \models \text{siphon_touch} \ \& \ \text{state}(\gamma, u2) \models \text{tail_shock}$$

$$\text{training_up_to}(\gamma, t1, u3, 3) \equiv$$

$$\exists u2, t3 [u2 < t3 \ \& \ u3 = t3 + 1]$$

$$\text{training_up_to}(\gamma, t1, u2, 2) \ \& \ \text{state}(\gamma, t3) \models \text{siphon_touch} \ \& \ \text{state}(\gamma, u3) \models \text{tail_shock}$$

Table 2. Temporal-Interactivist Representation Relation (specification, backward)

I.s.p.	Content (backward)
S(2)	$\forall t_1, u_1$ [training_up_to(γ , t_1 , u_1 , 1) & $\neg\exists t_0$ [training_up_to(γ , t_0 , u_1 , 2)] $\Rightarrow \exists t_2 > u_1$ [state(γ , t_2) \models S(2)]] $\forall t_1, u_2$ [training_up_to(γ , t_1 , u_2 , 2) & $\neg\exists t_0$ [training_up_to(γ , t_0 , u_2 , 3)] $\Rightarrow \exists t_3 > u_2$ [state(γ , t_3) $\not\models$ S(2)]]
S(3)	$\forall t_1, u_2$ [training_up_to(γ , t_1 , u_2 , 2) & $\neg\exists t_0$ [training_up_to(γ , t_0 , u_2 , 3)] $\Rightarrow \exists t_3 > u_1$ [state(γ , t_3) \models S(3)]] $\forall t_1, u_3$ [training_up_to(γ , t_1 , u_3 , 3) & $\neg\exists t_0$ [training_up_to(γ , t_0 , u_3 , 4)] $\Rightarrow \exists t_4 > u_3$ [state(γ , t_4) $\not\models$ S(3)]]
S(4)	$\forall t_1, u_3$ [training_up_to(γ , t_1 , u_3 , 3) & $\neg\exists t_0$ [training_up_to(γ , t_0 , u_3 , 4)] $\Rightarrow \exists t_4 > u_3$ [state(γ , t_4) \models S(4)]]

Table 3. Temporal-Interactivist Representation relation (syntactic level, forward)

I.s.p.	Content (forward)
S(4)	$\exists t' \geq t$ [state(γ , t') \models siphon_touch & $\forall t' \geq t$ [state(γ , t') \models siphon_touch $\Rightarrow \exists t'' \geq t'$ state(γ , t'') \models contraction]

Consider, for example, the backward representational content of state property S(2). According to Table 2, the occurrence of exactly one learning trial (indicated by the fact that at u_1 , a training period up to 1 but not up to 2 has passed) eventually leads to a time point where S(2) holds. In addition, to make the content more precise, it is specified that the occurrence of exactly two learning trials eventually causes S(2) not to hold.

As stated earlier, representational relations such as the ones specified here may correspond to certain expectations that the modeller has about the behaviour of the model. By (formally) specifying such expected representation relations and verifying them against simulation traces, it can be shown whether they are in accordance with the agent model's internal dynamics. This provides the possibility to validate not only the resulting observable behaviour of an agent simulation model against the observable behaviour of the agent in the real world, but also the internal, intermediate states of the agent in the model against the internal, intermediate states of the agent in the world. Thus it can be verified to what extent the model satisfies *internal realism* in addition to external realism.

8 Checking Dynamic Properties

In addition to the simulation software, a software environment has been developed that enables to check dynamic properties specified in `TTL` against simulation traces. This software environment takes a dynamic property and one or more (empirical or simulated) traces as input, and checks whether the dynamic property holds for the traces. Traces are represented by sets of Prolog facts of the form

```
holds(state(m1, t(2)), a, true).
```

where `m1` is the trace name, `t(2)` time point 2, and `a` is a state formula in the ontology of the component's input. It is indicated that state formula `a` is true in the component's input state at time point `t2`. The program for temporal formula checking basically uses Prolog rules for the predicate `sat` that reduce the satisfaction of the temporal formula finally to the satisfaction of atomic state formulae at certain time points, which can be read from the trace representation. Examples of such reduction rules are:

```
sat(and(F,G)) :- sat(F), sat(G).
sat(not(and(F,G))) :- sat(or(not(F), not(G))).
sat(or(F,G)) :- sat(F).
sat(or(F,G)) :- sat(G).
sat(not(or(F,G))) :- sat(and(not(F), not(G))).
```

Using automatic checks of this kind, many of the properties presented in this paper have been checked against traces such as the one depicted in Figure 3. In particular, dynamic property GP3 (expressing the learning behaviour) has been checked successfully against all generated traces. Furthermore, the representation relations denoted in Table 2 have been checked. The duration of these checks varied from 1 to 3 seconds, depending on the complexity of the formula. They all turned out to be successful, which validates (for the given traces at least) our choice for the representational content of the internal state properties. However, note that these checks are only an empirical validation, they are no exhaustive proof as, e.g., model checking is. Currently, the possibilities are explored to combine `TTL` with existing model checking techniques; cf. [5], [18], [22].

9 Discussion

This paper introduces an integrative modelling approach for simulation and analysis of adaptive agent behaviour and underlying mechanisms. The approach is integrative in two ways. First, it combines both *qualitative, logical* and *quantitative, numerical* aspects in one modelling framework. Second, it allows to model both dynamics at a *local level* (internal neural mechanisms of the agent; cf. [11]) and dynamics at a *global level* (externally observable agent behaviour, and representation relationships between internal and external states).

The neural processes of the *Aplysia* case study (cf. [10]) have been formalised by identifying executable local dynamic properties for the basic dynamics of *Aplysia*'s neural conditioning mechanism. On the basis of these local properties simulations have been made. Moreover, it is shown how the descriptions at these two levels (i.e., the level of the neurological mechanisms and of the overall behaviour) can be logically related to each other, which can be considered as a formalisation of the (interlevel) reduction relations between the two levels. Such interlevel relations can be useful in the analysis of simulation traces, because they allow the modeller to reduce the failure of a global behavioural property to the failure of a local internal property of the model. This can be useful in debugging a model, but also in the analysis of the circumstances under which a model will function well and under which not, and the reasons why.

Finally, the presented approach allows the modeller to (formally) specify and check representation relations, which relate internal or intermediate states of the agent simulation model to other states of the model, possibly at different time points. In this paper, it was explored how representation relations can be defined for adaptive agents, using approaches such as in [2], [13]; [16], pp. 200-202. The specifications of the representational content of the internal (neural) state properties for *Aplysia* have been validated by automatically checking them on the traces generated by the simulation model. As a result, not only the resulting observable behaviour of the agent simulation model has been validated against the observable behaviour of the agent in the real world, but also the internal states of the agent in the model have been validated against the internal states of the agent in the world. Thus it was verified to what extent the model satisfies *internal realism* in addition to external realism.

Concerning related work, in [3] another formal model is described of the dynamics of conditioning processes, using a similar modelling approach. However, that paper focuses on human conditioning, based on existing literature such as [17]. Instead, the current paper focuses on the specific case of *Aplysia*, of which the neural mechanisms are much simpler and therefore better understood. As a consequence, the model presented in the current paper is at a neural

level, whereas the model of [3] is at a functional level. Another difference is that their model concentrates more on the temporal aspects of the conditioning.

References

1. Balkenius, C. and Morén, J. Dynamics of a classical conditioning model. *Autonomous Robots*, 7, 1999, 41-56.
2. Bickhard, M.H. Representational Content in Humans and Machines. *Journal of Experimental and Theoretical Artificial Intelligence*, vol. 5, 1993, pp. 285-333.
3. Bosse, T., Jonker, C.M., Los, S.A., Torre, L. van der, and Treur, J. Formalisation and Analysis of the Temporal Dynamics of Conditioning. In: Mueller, J.P. and Zambonelli, F. (eds.), *Proceedings of the Sixth International Workshop on Agent-Oriented Software Engineering, AOSE'05*. Lecture Notes in Artificial Intelligence, vol. 3950. Springer Verlag, 2006, pp. 54-68. Extended version to appear in *Cognitive Systems Research Journal*, 2007.
4. Bosse, T., Jonker, C.M., Meij, L. van der, and Treur, J. LEADSTO: a Language and Environment for Analysis of Dynamics by SimulaTiOn. In: Eymann, T. et al. (eds.), *Proc. of the 3rd German Conference on Multi-Agent System Technologies, MATES'05*. LNAI 3550. Springer Verlag, 2005, pp. 165-178.
5. Clarke, E.M., O. Grumberg, D.A. Peled, (1999). *Model Checking*, MIT Press, Cambridge Massachusetts, London England.
6. Hindriks, K.V., F.S. de Boer, W. van der Hoek and J.-J.Ch. Meyer, Agent programming in 3APL. *Autonomous Agents and Multi-Agent Systems*, vol. 2, 1999, pp. 357-401.
7. Fisher, M. (2005). Temporal Development Methods for Agent-Based Systems, *Journal of Autonomous Agents and Multi-Agent Systems*, vol. 10, pp. 41-66.
8. Galton, A. (2003). Temporal Logic. *Stanford Encyclopedia of Philosophy*, URL: <http://plato.stanford.edu/entries/logic-temporal/#2>.
9. Galton, A. (2006). Operators vs Arguments: The Ins and Outs of Reification. *Synthese*, vol. 150, 2006, pp. 415-441.
10. Gleitman, H. *Psychology*, W.W. Norton & Company, New York, 1999.
11. Hawkins, R.D., and Kandel, E.R. Is There a Cell-Biological Alphabet for Simple Forms of Learning? *Psychological Review*, vol. 91, 1984, pp. 375-391.
12. Jonker, C.M., Snoep, J.L., Treur, J., Westerhoff, H.V., and Wijngaards, W.C.A. BDI-Modelling of Intracellular Dynamics. In: A.B. Williams and K. Decker (eds.), *Proc. of the First International Workshop on Bioinformatics and Multi-Agent Systems, BIXMAS'02*, 2002, pp. 15-23.
13. Jonker, C.M., and Treur, J. A Temporal-Interactivist Perspective on the Dynamics of Mental States. *Cognitive Systems Research Journal*, vol. 4, 2003, pp. 137-155.
14. Jonker, C.M., Treur, J., and Wijngaards, W.C.A. A Temporal Modelling Environment for Internally Grounded Beliefs, Desires and Intentions. *Cognitive Systems Research Journal*, vol. 4, 2003, pp. 191-210.

15. Keijzer, F. Representation in Dynamical and Embodied Cognition. *Cognitive Systems Research Journal*, vol. 3, 2002, pp. 275-288.
16. Kim, J., *Philosophy of Mind*. Westview Press, 1996.
17. Machado, A. Learning the temporal Dynamics of Behaviour. *Psychological Review*, vol. 104, 1997, pp. 241-265.
18. McMillan, K.L., (1993). *Symbolic Model Checking: An Approach to the State Explosion Problem*. PhD thesis, School of Computer Science, Carnegie Mellon University, Pittsburgh, 1992. Published by Kluwer Academic Publishers, 1993.
19. Pavlov, I.P. *Conditioned reflexes*. Oxford: Oxford University Press, 1927.
20. Port, R.F. and van Gelder, T.J. *Mind as Motion: Explorations in the Dynamics of Cognition*. MIT Press, Cambridge, Mass., 1995.
21. Reiter, R. (2001). *Knowledge in Action: Logical Foundations for Specifying and Implementing Dynamical Systems*. MIT Press, 2001.
22. Sharpanskykh, A., and Treur, J., Verifying Interlevel Relations within Multi-Agent Systems. In: *Proc. of the 17th European Conference on Artificial Intelligence, ECAI'06*. IOS Press, to appear, 2006.
23. Sun, R. Symbol grounding: a new look at an old idea. *Philosophical Psychology*, Vol.13, No.2, 2000, pp. 149-172.