

# Temporal and Interactionist Perspectives on the Dynamics of Mental States

Catholijn M. Jonker and Jan Treur

*Department of Artificial Intelligence, Vrije Universiteit Amsterdam,  
De Boelelaan 1081a, 1081 HV Amsterdam, The Netherlands*

*Email: {jonker, treur}@cs.vu.nl URL: <http://www.cs.vu.nl/~jonker,~treur>}*

**This paper contributes trace semantics and a temporal trace language to provide a grounding and formalisation of the dynamics of mental states in relation to the dynamics of the interaction with the external world. The approach provides a foundation for the dynamical and interactionist perspective on cognitive phenomena.**

## 1. INTRODUCTION

In recent literature on cognitive science and philosophy of mind, perspectives on cognitive functioning are proposed, where dynamics and interaction with the environment are central; e.g. (Bickhard, 1993; Clark, 1997; Port and van Gelder, 1995; Clapin et al., 2000; Christensen and Hooker, 2000). For example, (Bickhard, 1993) emphasises the relation between the (mental) state of a system (or agent) and its past and future in the interaction with its environment:

‘When interaction is completed, the system will end in some one of its internal states - some of its possible final states. (...) The final state that the systems ends up in, then, serves to implicitly categorise together that class of environments that would yield that final state if interacted with. (...) The overall system, with its possible final states, therefore, functions as a *differentiator* of environments, with the final states implicitly defining the differentiation categories. (...) Representational content is constituted as indications of potential further interactions. (...) The claim is that such differentiated functional indications in the context of a goal-directed system constitute representation - emergent representation.’

This suggests that mental states need to be grounded in interaction histories on the one hand, and have to be related to future interactions on the other hand. Some of the questions addressed in this paper are the following. What exactly is an interaction history? How does this precisely relate to a mental state? What about future traces, if they depend as well on the environment’s dynamics? How do they relate to mental states? To answer these questions, the temporal aspect of the dynamics of mental states and the interaction with the environment is formalised in this paper on the basis of formally defined traces and an expressive temporal trace language.

First, as a basis, in Section 2 the notion of trace and the temporal trace language are defined. In Section 3 it is shown how internal state properties can be formally related to sets of interaction traces to obtain their representational content or semantics. Section 4 addresses how these sets of traces can be characterised by formulae in the temporal trace language. Criteria are identified that express when a temporal formula defines a class of interaction traces that can be related to a specific mental property. Such a temporal formula can be viewed as a temporal grounding or temporal representation of this mental property. Section 5 concludes the paper with a brief discussion.

## 2. DYNAMICS, TRACES AND TRACE LANGUAGE

In this section dynamics is formalised by states that change over time. A *state* for a state ontology  $\text{Ont}$  (i.e., names for state properties) is an assignment of truth-values from  $\{\text{true}, \text{false}\}$  to the set of ground atoms  $\text{At}(\text{Ont})$ . The *set of all possible states* for ontology  $\text{Ont}$  is denoted by  $\text{STATES}(\text{Ont})$ . In particular,  $\text{STATES}(\text{OvOnt})$  denotes the set of all possible *overall states* (of the agent, its body and external world together);  $\text{STATES}(\text{InOnt})$  is the set of all of the agent's possible *internal states*. Moreover,  $\text{STATES}(\text{InOnt})$  is the set of possible input states (e.g., sensor states), and  $\text{STATES}(\text{OutOnt})$  the set of output states (e.g., effector states);  $\text{STATES}(\text{InterfaceOnt})$  denotes the set of all *interface states* (i.e., input and output together) between agent and external world. The standard satisfaction relation between states and state properties is used:  $s \models p$  means that property  $p$  holds in state  $s$ . For a state  $s$  over ontology  $\text{Ont}$  with sub-ontology  $\text{Ont}'$ , a restriction of  $s$  to  $\text{Ont}'$  can be made, denoted by  $s|_{\text{Ont}'}$ ; this restriction is the member of  $\text{STATES}(\text{Ont}')$  defined by  $s|_{\text{Ont}'}(a) = S(a)$  if  $a \in \text{At}(\text{Ont}')$ . For example, if  $s$  is an overall state, i.e., a member of  $\text{STATES}(\text{OvOnt})$ , then the restriction of  $s$  to the internal atoms,  $s|_{\text{InOnt}}$  is an internal state, i.e., a member of  $\text{STATES}(\text{InOnt})$ .

To describe behaviour of the agent, explicit reference is made to time in a formal manner. A fixed *time frame*  $\tau$  is assumed which is linearly ordered. Depending on the application, it may be dense (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.

A *trace*  $\mathcal{X}$  over an ontology  $\text{Ont}$  and time frame  $\tau$  is a mapping  $\mathcal{X}: \tau \rightarrow \text{STATES}(\text{Ont})$ , i.e., a sequence of states  $\mathcal{X}_t (t \in \tau)$  in  $\text{STATES}(\text{Ont})$ . The set of all traces over ontology  $\text{Ont}$  is denoted by  $\text{TRACES}(\text{Ont})$ , i.e.,  $\text{TRACES}(\text{Ont}) = \text{STATES}(\text{Ont})^\tau$ . A *temporal world description*  $\mathcal{W}$  is a set of traces over the overall ontology, i.e.,  $\mathcal{W} \subseteq \text{TRACES}(\text{OvOnt})$ . This set represents all possible developments over time (respecting the world's laws) of the agent and part of the world considered. This set formalises the processes of the agent, its body and the world as one integrated whole, as advocated in (Clark, 1997): 'Putting Brain, Body and World Together Again'. Given a trace  $\mathcal{X}$  over the overall ontology  $\text{OvOnt}$ , the state of the input interface at time point  $t$ , i.e.,  $\mathcal{X}_t|_{\text{InOnt}}$ , is denoted by  $\text{state}(\mathcal{X}, t, \text{InOnt})$ . Analogously,  $\text{state}(\mathcal{X}, t, \text{OutOnt})$  denotes the state of the output interface of the agent, and  $\text{state}(\mathcal{X}, t, \text{InOnt})$  the internal state at time point  $t$ .

To focus on different aspects of the agent and time, traces can be restricted to a specific ontology  $\text{Ont}$  and time interval  $\text{Interval}$ . The ontology parameter  $\text{Ont}$  indicates which parts of the agent or world are considered. For example, when this parameter is  $\text{InOnt}$ , then only input information is present in the restriction. The time interval parameter  $\text{Interval}$  specifies the part of the time frame of interest. The *restriction*  $\mathcal{X}_{\text{Interval}}^{\text{Ont}}$  of a trace  $\mathcal{X}$  to time in  $\text{Interval}$  and information based on  $\text{Ont}$  is a mapping  $\mathcal{X}_{\text{Interval}}^{\text{Ont}}: \text{Interval} \rightarrow \text{STATES}(\text{Ont})$ , defined by:  $\mathcal{X}_{\text{Interval}}^{\text{Ont}}(t) = \mathcal{X}(t)|_{\text{Ont}}$  if  $t \in \text{Interval}$ . For example, the *interaction trace*  $\mathcal{X}_{\leq t}^{\text{InterfaceOnt}}$  denotes the restriction of  $\mathcal{X}$  to the past up to  $t$  and to interface atoms.

Comparable to the approach in situation calculus, the sorted predicate logic temporal trace language  $\text{TTL}$  is built on atoms referring to, e.g., traces, time and

state properties, such as  $\text{state}(\mathcal{W}, t, \text{OutOnt}) \models p$ . Here  $\models$  is a predicate symbol in the language, comparable to the  $\text{Holds-}$ predicate in situation calculus. Temporal formulae are built using the usual logical connectives and quantification (for example, over traces, time and state properties). The set  $\text{TFOR}(\text{Ont})$  is the set of all *temporal formulae* that only make use of state ontology  $\text{Ont}$ . A *past formula* for  $\mathcal{W}$  and  $t$  is a temporal formula  $\Psi(\mathcal{W}, t)$  such that each time variable different from  $t$  is restricted to the time interval before  $t$ . In other words, for every time quantifier for a variable  $s$  a restriction of the form  $s \leq t$ , or  $s < t$  is required within the formula. The set of past formulae over ontology  $\text{Ont}$  w.r.t. time point  $t$  is denoted by  $\text{PFOR}(\text{Ont}, t)$ . Note that for any past formula  $\Psi(\mathcal{W}, t)$  it holds:

$$\forall \mathcal{W}, \mathcal{X} \in \mathcal{W} \forall t [ \mathcal{W}_{\leq t} = \mathcal{X}_{\leq t} \Rightarrow [ \Psi(\mathcal{W}, t) \Leftrightarrow \Psi(\mathcal{X}, t) ] ].$$

Similarly,  $\text{FFOR}(\text{Ont}, t)$  denotes the set of future formulae over ontology  $\text{Ont}$  w.r.t. time point  $t$ : every time quantifier for a variable  $s$  is restricted by  $s \geq t$  or  $s > t$ .

### 3. INTERNAL STATE PROPERTIES AND INTERACTION DYNAMICS

As put forward in the introduction, according to the interactionist view, a possible internal state ‘... serves to implicitly categorise together that class of environments that would yield that final state if interacted with’, cf. (Bickhart, 1993). Using our framework introduced in Section 2 the set of interaction histories over ontology  $\text{Ont}$  leading to internal property  $p$ , is formally defined by

$$\text{PTRACES}(\text{Ont}, p) = \{ \mathcal{W}_{\leq t}^{\text{Ont}} \mid t \in T, \mathcal{W} \in \mathcal{W}, \text{state}(\mathcal{W}, t, \text{IntOnt}) \models p \}$$

Besides, the way in which internal properties themselves lead to particular possible types of future interactions is also crucial for their meaning (Bickhard, 1993). Therefore, for an internal state property  $p$ , and an ontology  $\text{Ont}$ , the set of *all future traces* for  $t$  over  $\text{Ont}$  *allowed by*  $p$  is defined by:

$$\text{FTRACES}(\text{Ont}, p) = \{ \mathcal{W}_{\geq t}^{\text{Ont}} \mid t \in T, \mathcal{W} \in \mathcal{W}, \text{state}(\mathcal{W}, t, \text{IntOnt}) \models p \}$$

Based on these formal definitions, the *representational content of an internal state property*  $p$  is defined as the pair of sets  $\text{PTRACES}(\text{InterfaceOnt}, p)$ ,  $\text{FTRACES}(\text{InterfaceOnt}, p)$ .

The concepts introduced are illustrated with a example of the internal state property  $\text{pain}$ . This property is assumed to have relationships to the input properties  $\text{injury}$  and  $\text{heat}$ . Each of  $\text{injury}$  or  $\text{heat}$  causes  $\text{pain}$  (1), and they are the only possible causes (2). The set of world traces  $\mathcal{W}$  for this example reflects this in the sense that for any trace, always after  $\text{injury}$  occurs at the input, the internal state property  $\text{pain}$  will occur further on in the trace, and the same holds for  $\text{heat}$  at the input (1). Moreover, if  $\text{pain}$  occurs in a trace, then earlier in the trace one of or both  $\text{heat}$  and  $\text{injury}$  occurred at the input (2). An example of an interaction history on the *input* of the agent leading to  $\text{pain}$ , i.e., an element of  $\text{PTRACES}(\text{InOnt}, \text{pain})$ , is the following (partially depicted) interaction trace:

$$t_0. \text{injury: false, heat: false} ; \quad t_1. \text{injury: true, heat: false} ; \quad t_2. \text{injury: true, heat: false}$$

Note that in such a trace a delay may occur between the occurrence of the sensory input and the occurrence of the internal state property  $\text{pain}$ . How much delay  $\geq 0$  is taken into account is easily expressible in the temporal approach introduced by taking the real numbers as time frame.

For the future perspective, the internal state property  $\text{pain}$  is assumed to have relationships to the output properties  $\text{move}$  and  $\text{ouch!}$ . The property  $\text{pain}$

unconditionally causes the cry  $\text{ouch!}$ , and causes the action  $\text{move}$  depending on whether or not the environment object that caused the pain stays close or returns (e.g., a wasp) (1). It is assumed that these outputs are only generated if the internal state property  $\text{pain}$  holds (2). The set of world traces  $\alpha_{\mathcal{W}}$  reflects this in the sense that always after a time point where  $\text{pain}$  occurs in the internal state, (a) the output property  $\text{ouch!}$  will occur further on in the trace, and (b) if later on at the input  $\text{present}$  occurs, i.e., the object causing the pain is still there or returned, then this is followed by  $\text{move}$  at the output later on in the trace (within a certain time  $d$ , which for simplicity will be left out). Note that this implies learned behaviour: e.g., all wasps encountered in future will trigger an avoidance reaction. An example of a (partially depicted) interaction future allowed by  $\text{pain}$ , i.e., in  $\text{FTRACES}(\text{InterfaceOnt}, \text{pain})$ , is as follows:

t0. present: false ; t1. present: true ; t2. move: true

#### 4. TEMPORAL REPRESENTATIONS OF MENTAL STATES

Until now the interaction histories and futures have been formally defined by semantic set-theoretic means. However, using the temporal trace language introduced in Section 2, sets of traces can be characterised by temporal formulae as well. In the pain example the set  $\text{PTRACES}(\text{InOnt}, \text{pain})$  is characterised by:

$$\mathcal{X}_{\leq t}^{\text{InOnt}} \in \text{PTRACES}(\text{InOnt}, \text{pain}) \Leftrightarrow \Psi_{\text{P}}(\mathcal{X}, t)$$

where  $\Psi_{\text{P}}(\mathcal{X}, t) \in \text{PFOR}(\text{InterfaceOnt}, t)$  is the past formula

$$\exists t_1 \leq t \text{ state}(\mathcal{X}, t_1, \text{InOnt}) \models \text{injury} \vee \exists t_2 \leq t \text{ state}(\mathcal{X}, t_2, \text{InOnt}) \models \text{heat}.$$

The formula  $\Psi_{\text{P}}(\mathcal{X}, t)$  can be considered as an *external temporal representation* of the internal state property  $\text{pain}$ . If a nonzero delay of at least  $d$  is taken into account,  $\Psi_{\text{P}}(\mathcal{X}, t)$  has to be replaced by  $\Psi_{\text{P}}(\mathcal{X}, t-d)$ , to make the equivalence hold. If a delay with some randomness between 0 and  $d$  is assumed, then  $\Psi_{\text{P}}(\mathcal{X}, t)$  can be replaced by  $\exists d' 0 \leq d' \leq d \ \Psi_{\text{P}}(\mathcal{X}, t-d')$ , to guarantee the implication  $\Rightarrow$ . However, the implication  $\Leftarrow$  then does not hold.

For the future, to characterise the set of traces  $\text{FTRACES}(\text{InterfaceOnt}, \text{pain})$  in the form

$$\mathcal{X}_{\geq t}^{\text{InterfaceOnt}} \in \text{FTRACES}(\text{InterfaceOnt}, \text{pain}) \Leftrightarrow \Psi_{\text{F}}(\mathcal{X}, t)$$

a candidate formula  $\Psi_{\text{F}}(\mathcal{X}, t) \in \text{FFOR}(\text{InterfaceOnt}, t)$  is:

$$\exists t_1 \geq t \text{ state}(\mathcal{X}, t_1, \text{OutOnt}) \models \text{ouch!} \ \& \ \forall t_2 \geq t [\text{state}(\mathcal{X}, t_2, \text{InOnt}) \models \text{present} \Rightarrow \exists t_3 \geq t_2 \text{ state}(\mathcal{X}, t_3, \text{OutOnt}) \models \text{move}]$$

Also here a zero or fixed delay is assumed. This guarantees the implication  $\Rightarrow$ . However, an additional problem here is caused by the conditional in  $\Psi_{\text{F}}(\mathcal{X}, t)$ . Traces may occur where never the condition on  $\text{present}$  comes to hold. Then the implication is trivially true. But it is not satisfactory on this basis to conclude by the implication  $\Leftarrow$  that there has been pain (although in the example, also the  $\text{ouch!}$  property may play a role, but this does not solve the principle of the problem).

The pain example illustrates that temporal formulae characterising the representational content of an internal notion (in the sense of the past and future traces sets), in a simple manner depends on two assumptions: (1) fixed delay (2) no conditionals. Moreover, a silent assumption was (3): within the state ontology an internal state property exists for the considered notion ( $\text{pain}$ ). For a mathematical modelling approach, assumption (1) is customary (although not

quite desirable), so one could live with that. From an interactionist perspective, assumption (2), however, is unacceptable, because it excludes the possibility of the agent to let its behaviour in the future interactively depend on conditions that may or may not occur in a specific future interaction trace. Assumption (3) is innocent in the study of internal state properties and their content. However, if the attribution of mental properties based on observed behaviour is addressed, then assumption (3) would be artificial.

To avoid the assumptions discussed, the following notions are introduced. Let  $\Psi_{P(\mathcal{O}\mathcal{C}, t)} \in \text{PFOR}(\text{InterfaceOnt}, t)$  be a past formula and  $\Psi_{F(\mathcal{O}\mathcal{C}, t)} \in \text{FFOR}(\text{InterfaceOnt}, t)$  a future formula over the interface ontology. Moreover, let  $\text{Ont}$  be a given ontology (e.g., the internal ontology), and  $\phi(\mathcal{O}\mathcal{C}, t) \in \text{TFOR}(\text{Ont})$  a temporal formula over  $\text{Ont}$ .

The past formula  $\Psi_{P(\mathcal{O}\mathcal{C}, t)}$  is a *sufficient past interaction grounding* for  $\phi(\mathcal{O}\mathcal{C}, t)$  if:

$$\forall \mathcal{O}\mathcal{C} \in \mathcal{O}\mathcal{W} \forall t [ \forall \mathcal{N} \in \mathcal{O}\mathcal{W} [ \mathcal{O}\mathcal{C}_{\geq t} = \mathcal{N}_{\geq t} \Rightarrow \exists t_1 \leq t \Psi_{P(\mathcal{N}, t_1)} ] \Rightarrow \exists t_2 \geq t \phi(\mathcal{O}\mathcal{C}, t_2) ]$$

The past formula  $\Psi_{P(\mathcal{O}\mathcal{C}, t)}$  is a *necessary past interaction grounding* for  $\phi(\mathcal{O}\mathcal{C}, t)$  if:

$$\forall \mathcal{O}\mathcal{C} \in \mathcal{O}\mathcal{W} \forall t [ \phi(\mathcal{O}\mathcal{C}, t) \Rightarrow \forall \mathcal{N} \in \mathcal{O}\mathcal{W} [ \mathcal{O}\mathcal{C}_{\geq t} = \mathcal{N}_{\geq t} \Rightarrow \exists t_1 \leq t \Psi_{P(\mathcal{N}, t_1)} ] ]$$

The future formula  $\Psi_{F(\mathcal{O}\mathcal{C}, t)}$  is a *sufficient future interaction grounding* for  $\phi(\mathcal{O}\mathcal{C}, t)$  if:

$$\forall \mathcal{O}\mathcal{C} \in \mathcal{O}\mathcal{W} \forall t [ \forall \mathcal{N} \in \mathcal{O}\mathcal{W} [ \mathcal{O}\mathcal{C}_{\leq t} = \mathcal{N}_{\leq t} \Rightarrow \exists t_1 \geq t \Psi_{F(\mathcal{N}, t_1)} ] \Rightarrow \exists t_2 \leq t \phi(\mathcal{O}\mathcal{C}, t_2) ]$$

The future formula  $\Psi_{F(\mathcal{O}\mathcal{C}, t)}$  is a *necessary future interaction grounding* for  $\phi(\mathcal{O}\mathcal{C}, t)$  if:

$$\forall \mathcal{O}\mathcal{C} \in \mathcal{O}\mathcal{W} \forall t [ \phi(\mathcal{O}\mathcal{C}, t) \Rightarrow \forall \mathcal{N} \in \mathcal{O}\mathcal{W} [ \mathcal{O}\mathcal{C}_{\leq t} = \mathcal{N}_{\leq t} \Rightarrow \exists t_1 \geq t \Psi_{F(\mathcal{N}, t_1)} ] ]$$

## 5. CONCLUSIONS AND FURTHER RESEARCH

In the discussion on representational content of mental states, often the argument is made that for most mental properties no satisfactory way can be found to relate them to the (physical) world state, and hence symbolic or logical means are of no use to describe cognitive phenomena (the symbol grounding problem). Alternatives put forward (cf. (Clapin et al., 2000)) include the dynamical systems approach, and the interactionist perspective; cf. (Port and van Gelder, 1995; Bickhard, 1993; Christensen and Hooker, 2000). In line with these, in this paper the dynamic and interactionist perspective is adopted. It is shown how, if an interactionist perspective is taken, logical means in the form of temporal languages and semantics can successfully be used to describe the dynamics of mental states and properties, in relation to the dynamics of the interaction with the external world. Using this temporal approach, mental states and properties get their semantics in a formal manner in the temporal traces describing past and future interaction with the external world, in accordance with what is proposed informally by, e.g., (Bickhard, 1993; Christensen and Hooker, 2000; Clark, 1997).

The major difference with the work as mentioned is that in our approach a formalisation is proposed. This throws a new light on the sometimes assumed symbolic versus dynamics controversy. It shows how symbolic means can be used to describe dynamics as well; dynamics as a variety of phenomena entails no commitment to either Dynamical Systems Theory (DST) or symbolic methods as means to describe it.

The approach presented here contributes on the one hand a solid foundation for perspectives on dynamics and interaction as occurring in the recent literature. On

the other hand, the use of the temporal trace language has a number of practical advantages as well. In the first place, it offers a well-defined language to formulate relevant dynamic relations in practical domains, with standard first order logic semantics. It has a high expressive power. For example, the possibility of explicit reference to *time points* and *time durations* enables modelling of the dynamics of real-time phenomena, such as sensory and neural activity patterns in relation to mental properties (cf. (Port and van Gelder, 1995)). Also difference and differential equations can be expressed (see the extended paper). These features go beyond the expressivity available in standard linear or branching time temporal logics.

An interesting challenge for the temporal and interactionist perspective presented here is found in work based on the dual-level hypothesis, expressing that cognitive processes can be modelled according to two levels: the conceptual level (e.g., based on a symbolic model) and the sub-conceptual level (e.g., based on a connectionist model); cf. (Sun, 2000). The dual-level hypothesis would suggest to obtain a more refined temporal description of the dynamics that takes into account three elements (and their dynamic interaction): conceptual level mental properties, sub-conceptual properties, and the environment, where the sub-conceptual properties in a sense mediate between the conceptual properties and the environment. As both symbolic models and DST-style models are expressible in our language, it might be expected that also combinations of such types of models (and their interaction) can be expressed. This is planned as one of the issues for further research.

#### ACKNOWLEDGEMENTS

The authors have benefit from discussions about the subject with Vera Stebletsova, Wieke de Vries and Wouter Wijngaards.

#### REFERENCES

- Bickhard, M.H. (1993). Representational Content in Humans and Machines. *Journal of Experimental and Theoretical Artificial Intelligence*, 5, 1993, pp. 285-333.
- Clapin, H., Staines, P., and Slezak, P. (2000). *Proc. of the Int. Conference on Representation in Mind: New Theories of Mental Representation*, 27-29th June 2000, University of Sydney. To be published by Elsevier.
- Clark, A. (1997). *Being There: Putting Brain, Body and World Together Again*. MIT Press, 1997.
- Christensen, W.D. and C.A. Hooker (2000). *Representation and the Meaning of Life*. In: (Clapin et al., 2000).
- Port, R.F., Gelder, T. van (eds.) (1995). *Mind as Motion: Explorations in the Dynamics of Cognition*. MIT Press, Cambridge, Mass.
- Sun, R. (2000). Symbol grounding: a new look at an old idea. *Philosophical Psychology*, 13, 2000, pp. 149-172.