

# A Two-Level BDI-Agent Model for Theory of Mind and its Use in Social Manipulation

Tibor Bosse, Zulfiqar A. Memon, Jan Treur

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

URL: <http://www.few.vu.nl/~{tbosse, zamemon, treur}> Email: {tbosse, zamemon, treur}@few.vu.nl

## Abstract

This paper introduces a formal BDI-based agent model for the concept of Theory of Mind. The model uses BDI-concepts to describe the reasoning process of an agent that reasons about the reasoning process of another agent, which is also based on BDI-concepts. A case study illustrates how the model can be used for social manipulation. This case study addresses the scenario of a manager that reasons about the task avoiding behaviour of his employee. For this scenario, a number of simulation experiments have been performed, and some of their results are discussed.

## 1. INTRODUCTION

To function efficiently in social life and within organisations, it is useful if agents can reason about the actual and potential behaviour of the agents around them. To this end, it is very helpful for these agents to have capabilities to predict in which circumstances other agents will show certain appropriate or inappropriate behaviours. If for a considered other agent, generation of actions is assumed to be based on a BDI-model, prediction of such actions will involve reasoning based on a Theory of Mind (Baron-Cohen, 1995; Bogdan, 1997; Malle, Moses, and Baldwin, 2001) involving beliefs, desires and intentions as a basis for the behaviour. Reasoning based on such a Theory of Mind can be exploited in two different manners. The first manner is just to predict the behaviour in advance, in order to be prepared that it will occur (*social anticipation*). For example, if an agent B has done things that are known as absolutely unacceptable for an organisation (or a relationship), then he or she may be able to predict and therefore be prepared on what will happen after a manager (or partner) agent A learns about it.

A second manner to exploit reasoning based on a Theory of Mind is to try to affect the occurrence of certain beliefs, desires and intentions at forehand, by manipulating the occurrence of circumstances that are likely to lead to them (*social manipulation*). For example, the agent B just mentioned can try to hide facts so that the manager (or partner) agent A will never learn about the issue. Such capabilities of anticipatory and manipulatory reasoning based on a Theory of Mind about the behaviour of colleague agents are considered quite important, not to say essential, to function smoothly in social life.

This type of reasoning has an information acquisition and analysis aspect, and a preparation and action aspect. To describe the latter aspect, for the agent using a Theory of Mind, a model for action preparation based on beliefs, desires and intentions can be used as well. For example, for agent B discussed above, the desire can be generated that agent A will not perform the action to fire (or break up with) him or her, and that agent A will in particular not generate the desire or intention to do so. Based on this desire, the refined desire can be generated that agent A will not learn about the

issue. Based on the latter desire, an intention and action can be generated to hide facts for agent A. Notice that agent B reasons on the basis of BDI-models at two different levels, one for B itself, and one as the basis for the Theory of Mind to reason about agent A. It is this two-level architecture that is worked out in this paper in a computational model.

The modelling approach used for this computational model is based on the modelling language LEADSTO (Bosse, Jonker, Meij, and Treur, 2007). In this language, direct temporal dependencies between two state properties in successive states are modelled by *executable dynamic properties*. The LEADSTO format is defined as follows. Let  $\alpha$  and  $\beta$  be state properties of the form ‘conjunction of ground atoms or negations of ground atoms’. In the LEADSTO language the notation  $\alpha \rightarrow_{e, f, g, h} \beta$ , means:

*If state property  $\alpha$  holds for a certain time interval with duration  $g$ , then after some delay (between  $e$  and  $f$ ) state property  $\beta$  will hold for a certain time interval of length  $h$ .*

Here, atomic state properties can have a qualitative, logical format, such as an expression *desire(d)*, expressing that desire  $d$  occurs, or a quantitative, numerical format such as an expression *has\_value(x, v)* which expresses that variable  $x$  has value  $v$ .

In Section 2, first the general BDI-model is explained. In Section 3, this BDI-model is illustrated by a case study about an employee that shows task avoiding behaviour. Next, Section 4 describes how the simple model can be extended to a BDI-model of an agent that reasons about another agent’s BDI-model and uses this for social manipulation. In Section 5, this two-level BDI-model is illustrated by a case study that elaborates upon the example addressed in Section 3. This case study addresses the scenario of a manager that reasons about the task avoiding behaviour of his employee, and how to prevent that behaviour. Based on this model, some simulation experiments and their results are presented in Section 6. Section 7 concludes the paper with a discussion.

## 2. THE BDI-MODEL

The BDI-model bases the preparation and performing of actions on beliefs, desires and intentions (e.g., Georgeff and Lansky, 1987; Rao and Georgeff, 1991; Jonker, Treur, and Wijngaards, 2003). This model shows a long tradition in the literature, going back to Aristotle’s analysis of how humans (and animals) can come to actions; cf. (Aristotle, 350a BC, 350b BC). He discusses how the occurrence of certain internal (mental) state properties within the living being entail or cause the occurrence of an action in the external world. Such internal state properties are sometimes called by him ‘things in the soul’, for example, sensation, reason and desire<sup>1</sup>. Here, sensation indicates the sensing of the environment by

<sup>1</sup> Now there are three things in the soul which control action and truth - sensation, reason, desire. (Aristotle, 350 BC, *Nicomachean Ethics*), Book VI, Part 2

the agent, which leads, (in modern terms) to internal representations, called beliefs. Reason indicates the (rational) choice of an action that is reasonable to fulfil the given desire. Based on this, Aristotle introduced the following pattern to explain action (called practical syllogism):

If A has a desire D  
and A has the belief that X is a (or: the best) means to achieve D  
then A will do X

The BDI-model incorporates such a pattern of reasoning to explain behaviour in a refined form. Instead of a process from desire to action in one step, as an intermediate stage first an intention is generated, and from the intention the action is generated. Thus the process is refined into a two-step process. See Figure 1 for the generic structure of the BDI-model in causal-graph-like style, as often used to visualise LEADSTO specifications. Here the box indicates the borders of the agent, the circles denote state properties, and the arrows indicate dynamic properties expressing that one state property leads to (or causes) another state property.

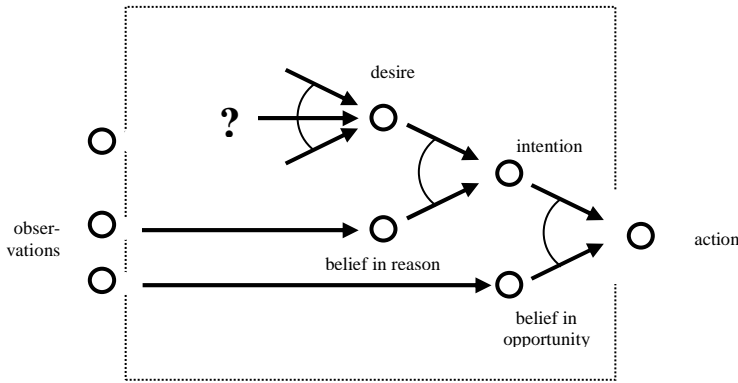


Figure 1 Structure of the general BDI-model

In this model, an action is performed when the subject has the intention to do this action and it has the belief that certain circumstances in the world are fulfilled such that the opportunity to do the action is there. Beliefs are created on the basis of observations. The intention to do a specific type of action is created if there is some desire D, and there is the belief that certain circumstances in the world state are there, that make it possible that performing this action will fulfil this desire (this is the kind of rationality criterion discussed above; e.g., what is called means-end analysis is covered by this). Whether or not a given action is adequate to fulfil a given desire depends on the current world state; therefore this belief may depend on other beliefs about the world state. Instantiated relations within the general BDI-model as depicted by arrows in graphical format in Figure 1 can be specified in formal LEADSTO format as follows:

$\text{desire}(D) \wedge \text{belief}(B1) \rightarrow \text{intention}(P)$   
 $\text{intention}(P) \wedge \text{belief}(B2) \rightarrow \text{performs}(P)$

with appropriate desire D, action P and beliefs B1, B2. Note that the beliefs used here both depend on observations, as shown in Figure 1. Furthermore,  $\wedge$  stands for the conjunction operator (and) between the atomic state properties (in the graphical format denoted by an arc connecting two (or more) arrows). Often, dynamic properties in LEADSTO are presented in *semi-formal* format, as follows:

At any point in time  
if desire D is present

and the belief B1 is present  
then the intention for action P will occur

At any point in time  
if the intention for action P is present  
and the belief B2 is present  
then the action P will be performed

As a generic template, including a reference to the agent X concerned, this can be expressed by:

For any desire D, world state property Z, and action Y such that  $\text{has\_reason\_for}(X, D, Z, Y)$  holds:

$\text{desire}(X, D) \wedge \text{belief}(X, Z) \rightarrow \text{intention}(X, Y)$

For any world state property Z and action Y such that  $\text{is\_opportunity\_for}(X, Z, Y)$  holds:

$\text{intention}(X, Y) \wedge \text{belief}(X, Z) \rightarrow \text{performs}(X, Y)$

Here  $\text{has\_reason\_for}(X, D, Z, Y)$  is a relation that can be used to specify which state property Z is considered a reason to choose a certain intention Y for desire D. Similarly  $\text{is\_opportunity\_for}(X, Z, Y)$  is a relation that can be used to specify which state property Z is considered an opportunity to actually perform an intended action Y.

Assuming that beliefs are available, what remains to be generated in this model are the desires. For desires, there is no generic way (known) in which they are to be generated in the standard model. Often, in applications, generation of desires depends on domain-specific knowledge.

### 3. A BDI-MODEL FOR TASK AVOIDANCE

To illustrate the BDI-model described above by a specific example, the following scenario is addressed (in the domain of an organisation); notice that here no Theory of Mind is involved.

#### Task Avoidance Case

A manager observes that a specific employee in the majority of cases functions quite cooperatively, but shows avoidance behaviour in other cases. In these latter cases, the employee starts trying to reject the task if he believes that his agenda already was full-booked for the short term, and he believes that capable colleagues are available with not full-booked agendas. Further observation by the manager reveals the pattern that the employee shows avoidance behaviour, in particular, in cases that a task is only asked shortly before its deadline, without the possibility to anticipate on the possibility of having the task allocated. The manager deliberates about this as follows:

*'If I know beforehand the possibility that a last-minute task will occur, I can tell him the possibility in advance, and in addition point out that I need his unique expertise for the task, in order to avoid the behaviour that he tries to avoid the task when it actually comes up.'*

Below, this example is formalised, using the BDI-model as introduced above. First, only the behaviour of the employee is addressed (in Section 5, the deliberation process of the manager is addressed as well). To this end, the example is made more precise as follows:

The *desire* to avoid a task is created after time t by the employee if the following holds:

- the employee has the belief at time t that a task is requested that has to be finished soon
- the employee has the belief at time t that he did not hear of the possibility of the task at any earlier time point

The *intention* to avoid a task is generated after time t if the following holds:

- the desire to avoid the task is available at time t
- the belief that capable not full-booked colleagues are available at time t

The *action* to avoid the task is generated after time  $t$  if the following holds:

- the intention to avoid the task is available at time  $t$
- the belief that the employee's own agenda is full-booked is available at time  $t$

Using the generic template discussed at the end of Section 2, via the relations

$\text{has\_reason\_for}(A, \text{lower\_workload}, \text{capable\_colleagues\_available}, \text{avoid\_task})$   
 $\text{is\_opportunity\_for}(A, \text{own\_agenda\_full}, \text{avoid\_task})$

the following model for agent A is obtained:

$\text{belief}(A, \text{task\_may\_come}) \wedge \text{belief}(\text{last\_minute\_request}) \rightarrow$   
 $\text{desire}(A, \text{lower\_workload})$

$\text{desire}(A, \text{lower\_workload}) \wedge \text{belief}(A, \text{capable\_colleagues\_available}) \rightarrow$   
 $\text{intention}(A, \text{avoid\_task})$

$\text{intention}(A, \text{avoid\_task}) \wedge \text{belief}(A, \text{own\_agenda\_full}) \rightarrow$   
 $\text{performs}(A, \text{avoid\_task})$

#### 4. THE TWO-LEVEL BDI-MODEL

As an instance of the *instrumentalist perspective* and opposed to explanations from a direct physical perspective (the physical stance), in (Dennett, 1987, 1991) the *intentional stance* (or folk-psychological stance) is put forward. In (Dennett, 1987), pp. 37-39, he explains the advantage of intentional stance explanations for mental phenomena over physical stance explanations.<sup>2</sup> According to the intentional stance, an agent is assumed to decide to act and communicate based on intentional notions such as beliefs about its environment and its desires and intentions. These decisions, and the intentional notions by which they can be explained and predicted, generally depend on circumstances in the environment, and, in particular, on the information on these circumstances just acquired by interaction (i.e., by observation and communication), but also on information acquired by interaction in the past. To be able to analyse the occurrence of intentional notions in the behaviour of an observed agent, the observable behavioural patterns over time form a basis; cf. (Dennett, 1991).

In the model presented in this paper, the instrumentalist perspective is taken as a point of departure for a Theory of Mind. More specifically, the model describes the reasoning process of an agent B that applies the intentional stance to another agent A by attributing beliefs, desires and intentions. Thus, for agent B a Theory of Mind is obtained using concepts for agent A's beliefs, desires and intentions. For example, in case a manager has an important last-minute task for his employee, but he knows that this employee often shows avoidance behaviour for last-minute tasks, he may analyse in more detail under which circumstances the employee may generate the desire and intention to avoid this task, and the related beliefs in reason and opportunity.

As a next step, the model is extended with BDI-concepts for agent B's own beliefs, desires and intentions as well. By doing this, agent B is able to not only *have* a theory about the mind of agent A, but also to *use* it within its own BDI-based reasoning processes. To this end, a number of meta-representations expressed by meta-predicates are introduced, e.g.:

$\text{belief}(B, \text{desire}(A, D))$

This expresses that agent B believes that agent A has desire D.

$\text{desire}(B, \text{not}(\text{intention}(A, X)))$

This expresses that agent B desires that agent A does not intend action X.

$\text{belief}(B, \text{depends\_on}(\text{performs}(A, X), \text{intention}(A, X)))$

This expresses that agent B believes that, whether A will perform action X depends on whether A intends to do X. Note that the third meta-statement has a more complex structure than the other two, since it represents a statement about a *dynamic property*, rather than a statement about a *state property*. These dependencies can be read from a graph such as depicted in Figures 1 and 2 (right hand side). For example, it is assumed that agent B knows part of this graph in his Theory of Mind, expressed by beliefs such as:

$\text{belief}(B, \text{depends\_on}(\text{performs}(A, X), \text{intention}(A, X)))$   
 $\text{belief}(B, \text{depends\_on}(\text{performs}(A, P), \text{belief}(A, B2)))$   
 $\text{belief}(B, \text{depends\_on}(\text{intention}(A, P), \text{desire}(A, D)))$   
 $\text{belief}(B, \text{depends\_on}(\text{intention}(A, P), \text{belief}(A, B1)))$   
 $\text{belief}(B, \text{depends\_on}(\text{desire}(A, D), \text{belief}(A, B3)))$   
 $\text{belief}(B, \text{depends\_on}(\text{belief}(A, X), \text{hears}(A, X)))$

Desire refinement in the BDI-model for an agent B attributing motivations to an agent A is formulated (in LEADSTO format) by:

$\text{desire}(B, X) \wedge \text{belief}(B, \text{depends\_on}(X, Y)) \rightarrow \text{desire}(B, Y)$

$\text{desire}(B, \text{not}(X)) \wedge \text{belief}(B, \text{depends\_on}(X, Y)) \rightarrow \text{desire}(B, \text{not}(Y))$

Moreover the following schemes for intention and action generation are included in the model. For any desire D, world state property Z, and action Y such that  $\text{has\_reason\_for}(B, D, Z, Y)$  holds:

$\text{desire}(B, D) \wedge \text{belief}(B, Z) \rightarrow \text{intention}(B, Y)$

For any world state property Z and action Y such that  $\text{is\_opportunity\_for}(B, Z, Y)$  holds:

$\text{intention}(B, Y) \wedge \text{belief}(B, Z) \rightarrow \text{performs}(B, Y)$

Moreover, some dynamic properties of the world are needed:

$\text{performs}(B, \text{tell}(A, C)) \rightarrow \text{holds\_in\_world}(\text{communication}(B, A, C))$

$\text{holds\_in\_world}(\text{communication}(B, A, C)) \rightarrow \text{hears}(A, C)$

For an overview of the complete two-level BDI-model, see Fig. 2.

#### 5. A TWO-LEVEL BDI-MODEL FOR REASONING ABOUT TASK AVOIDANCE

The above model can be used to describe how the manager agent (from the case described in Section 3) can reason and act in an anticipatory manner to avoid the employee's avoidance desire, intention and/or action to occur. The initial desire of B is that A does not perform the action to avoid the task:

$\text{desire}(B, \text{not}(\text{performs}(A, \text{avoid\_task})))$

Fulfilment of this desire can be obtained in the following three manners:

*Avoiding A's desire to occur*

This can be obtained when the employee hears in advance that possibly a last minute task may occur. This will make the second condition in A's desire generation as described in Section 3 fail.

*Avoiding A's intention to occur (given that the desire occurs)*

This can be obtained by refutation of the belief that plays the role of the reason to generate the intention in A's intention generation as described in Section 3, e.g., when the employee hears that colleagues do not have the required expertise.

*Avoiding A's action to occur (given that the intention occurs)*

This can be obtained by refutation of the belief that plays the role of opportunity in A's desire action as described in Section 3, e.g., when his agenda is not full-booked.

For convenience, the model does not make a selection but addresses all three options to prevent the avoidance action. This means that B generates desires for:

- A hears about the possibility of a last-minute task in advance  
 $\text{hears}(A, \text{task\_may\_come})$
- A hears that no colleagues that are capable of performing the task are available  
 $\text{hears}(A, \text{not}(\text{capable\_colleagues\_available}))$
- A hears that his agenda is not full-booked  
 $\text{hears}(A, \text{not}(\text{own\_agenda\_full}))$

<sup>2</sup> 'Predicting that someone will duck if you throw a brick at him is easy from the folk-psychological stance; it is and will always be intractable if you have to trace the photons from brick to eyeball, the neurotransmitters from optic nerve to motor nerve, and so forth.' (Dennett, 1991), p. 42

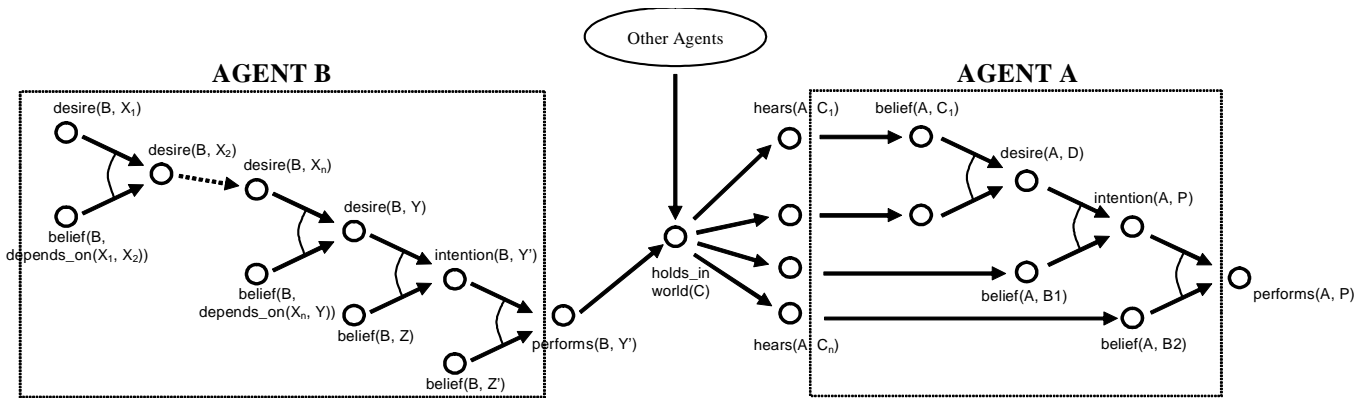


Figure 2 Structure of the Two-Level BDI-model

To fulfil these desires, intentions are to be generated by B to perform actions such as:

- B tells A about the possibility of a last-minute task in advance  
performs(B, tell(A, task\_will\_come))
- B tells A that no colleagues that are capable of performing the task are available  
performs(B, tell(A, not(capable\_colleagues\_available)))
- B tells A that some of the (perhaps less interesting) tasks were taken from A's agenda and were re-allocated to a colleague  
performs(B, tell(A, not(own\_agenda\_full)))

Reason for B to choose for these actions is

- the belief of B that telling something will lead to the person hearing it  
belief(B, adequate\_communication(B, A))

Moreover, these intentions of B can lead to the corresponding actions when the following belief of B in opportunity is there:

- the belief that A is available for B to talk to  
belief(B, available\_for(A, B))

In addition to the generic BDI-model shown in Section 4, the following specific relations are used to model the case study:

has\_reason\_for(B, hears(A, C), adequate\_communication, tell(A, C))

is\_opportunity\_for(B, available\_for(A, B), tell(A, C))

Note that the last minute request itself is an event that not necessarily comes from agent B; it can come from any agent, for example a Director agent. It is modelled as an event in LEADSTO.

## 6. SIMULATION EXPERIMENTS

In a number of simulation experiments, the two-level BDI-model has been applied to the case study as described in Section 5. To this end, the LEADSTO software environment (Bosse, Jonker, Meij, and Treur, 2007) has been used. In Figure 3, an example of a resulting simulation trace is shown. In this figure, time is on the horizontal axis; the state properties are on the vertical axis. A box on top of a line indicates that a state property is true. Note that, due to space limitations, only a selection of the relevant atoms is shown.

Figure 3 shows that the manager initially desires that the employee does not perform the action to avoid the task (desire(manager, not(performs(employee, avoid\_task))). Based on this, he eventually generates number of more detailed desires about what the employee should hear (see, for example, the state property desire(manager, not(hears(employee, capable\_colleagues\_available))) at time point 3). Next, the manager uses these desires to generate some intentions to fulfil these desires (e.g., the state property intention(manager, tell(employee, not(capable\_colleagues\_available))) at time point 4). Eventually, these intentions are performed, and the employee receives some new inputs (e.g., the state property hears(employee, not(capable\_colleagues\_available)) at time point 7). As a result, when the employee receives a last minute request at time

point 11 (hears(employee, last\_minute\_request)), he does not generate the action to avoid the task.

Note that in the scenario sketched in Figure 3, the manager takes all possible actions (within the given conceptualisation) to fulfil its desires. This is a rather extreme case, since according to the employee's BDI-model, modifying only one of its input will be sufficient to make sure that (s)he does not avoid the task. Other traces can be generated in which the manager takes less actions to fulfil its desires.



Figure 3 Example Simulation Trace

## 7. DISCUSSION

In order to function efficiently in social life, it is very helpful for an agent to have capabilities to predict in which circumstances the agents in its environment will show certain behaviours. To this end, such an agent will have to perform reasoning based on a Theory of Mind (Baron-Cohen, 1995). This paper presents a model for reasoning based on a Theory of Mind, which makes use of BDI-concepts at two different levels. First, the model uses BDI-concepts *within* the Theory of Mind (i.e., it makes use of beliefs, desires and intentions to describe the reasoning process of another agent). Second, it uses BDI-concepts for reasoning *about* the Theory of Mind (i.e., it makes use of beliefs, desires and intentions to describe an agent's meta-reasoning about the reasoning process of another agent). At this second level, meta-statements are involved, such as 'B believes that A desires d' or 'B desires that A does not intend a'. These meta-statements are about the states occurring within the other agent. In addition, meta-statements are involved about the dynamics occurring within the other agents. An example of such a

(more complex) meta-statement is 'B believes that, if A performs a, then earlier he or she intended a'.

The two-level BDI-based model can be exploited both for *social anticipation* (i.e., in order to be prepared for the behaviour of another agent) and for *social manipulation* (i.e., in order to affect the behaviour of another agent at forehand). The model has been formalised using the high-level modelling language LEADSTO, which describes dynamics in terms of direct temporal dependencies between state properties in successive states. Based on the formal model, a number of simulation experiments have been performed within a specific case study, addressing the scenario of a manager that reasons about the task avoiding behaviour of his employee.

The case study illustrates how the two-level model can be used for social manipulation. For this purpose, the crucial steps are to find out which situations would lead to undesired behaviour of another agent, and to prevent these situations from occurring (or, similarly, to establish situations that would lead to desired behaviour). In addition, the model can be used for social anticipation. In that case, the main steps are to predict the behaviour that another agent will show given the current situation, and to prepare for this. Also this second type of reasoning based on a Theory of Mind is essential to function smoothly in social life. Applying the model to social anticipation will be the subject of future research.

## REFERENCES

- Aristotle (350a BC). *Nicomachean Ethics* (translated by W.D. Ross)  
Aristotle (350b BC). *De Motu Animalium* On the Motion of Animals (translated by A. S. L. Farquharson)

- Baron-Cohen, S. (1995). *Mindblindness*. MIT Press.  
Bogdan, R.J. (1997). *Interpreting Minds*. MIT Press.  
Bosse, T., Jonker, C.M., Meij, L. van der, and Treur, J. (2007). A Language and Environment for Analysis of Dynamics by Simulation. *International Journal of Artificial Intelligence Tools*. To appear, 2007. Shorter, earlier version in: Eymann, T., Kluegl, F., Lamersdorf, W., Klusch, M., and Huhns, M.N. (eds.), *Proc. of the Third German Conference on Multi-Agent System Technologies, MATES'05*. Lecture Notes in Artificial Intelligence, vol. 3550. Springer Verlag, 2005, pp. 165-178.  
Dennett, D.C. (1987). *The Intentional Stance*. MIT Press. Cambridge Mass.  
Dennett, D.C. (1991). Real Patterns. *The Journal of Philosophy*, vol. 88, pp. 27-51.  
Georgeff, M. P., and Lansky, A. L. (1987). Reactive Reasoning and Planning. In: *Proceedings of the Sixth National Conference on Artificial Intelligence, AAAI'87*. Menlo Park, California. American Association for Artificial Intelligence, 1987, pp. 677-682.  
Jonker, C.M., Treur, J., and Wijngaards, W.C.A., (2003). A Temporal Modelling Environment for Internally Grounded Beliefs, Desires and Intentions. *Cognitive Systems Research Journal*, vol. 4(3), 2003, pp. 191-210.  
Malle, B.F., Moses, L.J., Baldwin, D.A. (2001). *Intentions and Intentionality: Foundations of Social Cognition*. MIT Press.  
Rao, A.S. and Georgeff, M.P. (1991). Modelling Rational Agents within a BDI-architecture. In: Allen, J., Fikes, R. and Sandewall, E. (eds.), *Proceedings of the Second International Conference on Principles of Knowledge Representation and Reasoning, (KR'91)*. Morgan Kaufmann, pp. 473-484.