

Intelligent Agents: The Key Concepts

Michael Wooldridge

Department of Computer Science
University of Liverpool
Liverpool L69 7ZF, UK
mjl@csc.liv.ac.uk

Abstract. This chapter aims to introduce the reader to the basic issues surrounding the design and implementation of intelligent agents. It begins by motivating the idea of an agent, presents a definition of agents and intelligent agents. The article then goes on to discuss four major approaches to building agents. First, logic based architectures are reviewed, in which decision-making is viewed as logical deduction. Second, reactive architectures are discussed, in which symbolic representations and models are eschewed in favour of a closer relationship between agent perception and action. Third, we discuss belief-desire-intention architectures, in which decision making is viewed as practical reasoning from beliefs about how the world is and will be to the options available to an agent, and finally to intentions and actions. Fourth, we review layered agent architectures, in which decision making is partitioned into a number of different decision making layers, each dealing with the agent's environment at a different level of abstraction.

1 Introduction

Computers are not very good at knowing what to do: every action a computer performs must be explicitly anticipated, planned for, and coded by a programmer. If a computer program ever encounters a situation that its designer did not anticipate, then the result is not usually pretty – a system crash at best, multiple loss of life at worst. This mundane fact is at the heart of our relationship with computers. It is so self-evident to the computer literate that it is rarely mentioned. And yet it comes as a complete surprise to those encountering computers for the first time.

For the most part, we are happy to accept computers as obedient, literal, unimaginative servants. For many applications (such as payroll processing), it is entirely acceptable. However, for an increasingly large number of applications, we require systems that can *decide for themselves* what they need to do in order to satisfy their design objectives. Such computer systems are known as *agents*. Agents that must operate robustly in rapidly changing, unpredictable, or open environments, where there is a significant possibility that actions can *fail* are known as *intelligent agents*, or sometimes *autonomous agents*. Here are examples of recent application areas for intelligent agents:

- When a space probe makes its long flight from Earth to the outer planets, a ground crew is usually required to continually track its progress, and decide how to deal with unexpected eventualities. This is costly and, if decisions are required *quickly*, it is simply not practicable. For these reasons, organisations like NASA are seriously investigating the possibility of making probes more autonomous – giving them richer decision making capabilities and responsibilities.
- Searching the Internet for the answer to a specific query can be a long and tedious process. So, why not allow a computer program – an agent – do searches for us? The agent would typically be given a query that would require synthesising pieces of information from various different Internet information sources. Failure would occur when a particular resource was unavailable, (perhaps due to network failure), or where results could not be obtained.

This chapter is about intelligent agents. Specifically, it aims to give you a thorough introduction to the main issues associated with the design and implementation of intelligent agents. After reading it, you will understand:

- why agents are believed to be an important new way of conceptualising and implementing certain types of software application;
- what intelligent agents are (and are not);
- the main approaches that have been advocated for designing and implementing intelligent agents, the issues surrounding these approaches, their relative merits, and the challenges that face the agent implementor.

The chapter is structured as follows. First, Sect. 2 describes what is meant by the term *agent*. Section 3, presents some *abstract architectures* for agents. That is, some general models and properties of agents are discussed without regard to how they might be implemented. Section 4, discusses *concrete* architectures for agents. The various major design routes that one can follow in implementing an agent system are outlined in this section. In particular, *logic-based* architectures, *reactive* architectures, *belief-desire-intention* architectures, and finally, *layered* architectures for intelligent agents are described in detail.

Comments on Notation. This chapter makes use of simple mathematical notation in order to make ideas precise. The formalism used is that of discrete maths: a basic grounding in sets and first-order logic should be quite sufficient to make sense of the various definitions presented. In addition: if S is an arbitrary set, then $\wp(S)$ is the powerset of S , and S^* is the set of sequences of elements of S ; the symbol \neg is used for logical negation (so $\neg p$ is read “not p ”); \wedge is used for conjunction (so $p \wedge q$ is read “ p and q ”); \vee is used for disjunction (so $p \vee q$ is read “ p or q ”); and finally, \Rightarrow is used for material implication (so $p \Rightarrow q$ is read “ p implies q ”).

2 What Are Agents?

An obvious way to open this chapter would be by presenting a definition of the term *agent*. After all, this is a book about multi-agent systems – surely we must all agree on what an agent is? Surprisingly, there is no such agreement: there is no universally accepted definition of the term agent, and indeed there is a good deal of ongoing debate and controversy on this very subject. Essentially, while there is a general consensus that *autonomy* is central to the notion of agency, there is little agreement beyond this. Part of the difficulty is that various attributes associated with agency are of differing importance for different domains. Thus, for some applications, the ability of agents to *learn* from their experiences is of paramount importance; for other applications, learning is not only unimportant, it is undesirable.

Nevertheless, some sort of definition is important – otherwise, there is a danger that the term will lose all meaning (cf. “user friendly”). The definition presented here is adapted from [55]: An *agent* is a computer system that is *situated* in some *environment*, and that is capable of *autonomous action* in this environment in order to meet its design objectives.

There are several points to note about this definition. First, the definition refers to “agents” and not “intelligent agents”. The distinction is deliberate: it is discussed in more detail below. Second, the definition does not say anything about what *type* of environment an agent occupies. Again, this is deliberate: agents can occupy many different types of environment, as we shall see below. Third, we have not defined *autonomy*. Like agency itself, autonomy is a somewhat tricky concept to tie down precisely. In this chapter, it is used to mean that agents are able to act without the intervention of humans or other systems: they have control both over their own internal state, and over their behavior.

Figure 1 gives an abstract, top-level view of an agent. In this diagram, we can see the action output generated by the agent in order to affect its environment. In most domains of reasonable complexity, an agent will not have *complete* control over its environment. It will have at best *partial* control, in that it can *influence* it. From the point of view of the agent, this means that the same action performed twice in apparently identical circumstances might appear to have entirely different effects, and in particular, it may *fail* to have the desired effect. Thus agents in all but the most trivial of environments must be prepared for the possibility of *failure*. We can sum this situation up formally by saying that environments are *non-deterministic*.

Normally, an agent will have a repertoire of actions available to it. This set of possible actions represents the agents *effectoric capability*: its ability to modify its environments. Note that not all actions can be performed in all situations. For example, an action “lift table” is only applicable in situations where the weight of the table is sufficiently small that the agent *can* lift it. Similarly, the action “purchase a Ferrari” will fail if insufficient funds are available to do so. Actions therefore have *pre-conditions* associated with them, which define the possible situations in which they can be applied.

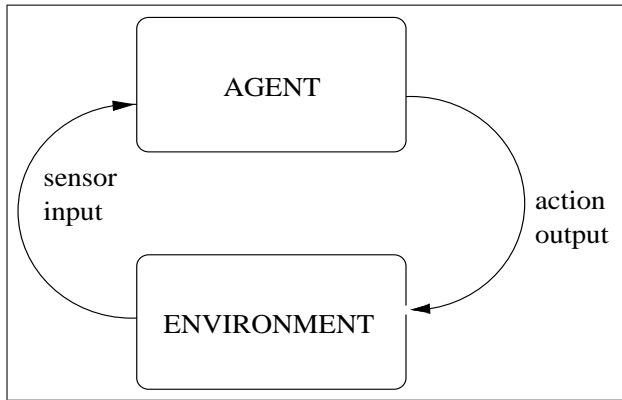


Fig. 1. An agent in its environment. The agent takes sensory input from the environment, and produces as output actions that affect it. The interaction is usually an ongoing, non-terminating one.

The key problem facing an agent is that of deciding *which* of its actions it should perform in order to best satisfy its design objectives. *Agent architectures*, of which we shall see several examples later in this chapter, are really software architectures for decision making systems that are embedded in an environment. The complexity of the decision-making process can be affected by a number of different environmental properties. Russell and Norvig suggest the following classification of environment properties [49, p46]:

- *Accessible vs inaccessible.*

An accessible environment is one in which the agent can obtain complete, accurate, up-to-date information about the environment’s state. Most moderately complex environments (including, for example, the everyday physical world and the Internet) are inaccessible. The more accessible an environment is, the simpler it is to build agents to operate in it.

- *Deterministic vs non-deterministic.*

As we have already mentioned, a deterministic environment is one in which any action has a single guaranteed effect – there is no uncertainty about the state that will result from performing an action. The physical world can to all intents and purposes be regarded as non-deterministic. Non-deterministic environments present greater problems for the agent designer.

- *Episodic vs non-episodic.*

In an episodic environment, the performance of an agent is dependent on a number of discrete episodes, with no link between the performance of an agent in different scenarios. An example of an episodic environment would be a mail sorting system [50]. Episodic environments are simpler from the agent developer’s perspective because the agent can decide what action to perform based only on the current episode – it need not reason about the interactions between this and future episodes.

– *Static vs dynamic.*

A static environment is one that can be assumed to remain unchanged except by the performance of actions by the agent. A dynamic environment is one that has other processes operating on it, and which hence changes in ways beyond the agent’s control. The physical world is a highly dynamic environment.

– *Discrete vs continuous.*

An environment is discrete if there are a fixed, finite number of actions and percepts in it. Russell and Norvig give a chess game as an example of a discrete environment, and taxi driving as an example of a continuous one.

As Russell and Norvig observe [49, p46], if an environment is sufficiently complex, then the fact that it is *actually* deterministic is not much help: to all intents and purposes, it may as well be non-deterministic. The most complex general class of environments are those that are inaccessible, non-deterministic, non-episodic, dynamic, and continuous.

2.1 Examples of Agents

At this point, it is worth pausing to consider some examples of agents (though not, as yet, intelligent agents):

- Any *control* system can be viewed as an agent. A simple (and overused) example of such a system is a thermostat. Thermostats have a sensor for detecting room temperature. This sensor is directly embedded within the environment (i.e., the room), and it produces as output one of two signals: one that indicates that the temperature is too low, another which indicates that the temperature is OK. The actions available to the thermostat are “heating on” or “heating off”. The action “heating on” will generally have the effect of raising the room temperature, but this cannot be a *guaranteed* effect – if the door to the room is open, for example, switching on the heater may have no effect. The (extremely simple) decision making component of the thermostat implements (usually in electro-mechanical hardware) the following rules:

too cold \rightarrow heating on
 temperature OK \rightarrow heating off

More complex environment control systems, of course, have considerably richer decision structures. Examples include autonomous space probes, fly-by-wire aircraft, nuclear reactor control systems, and so on.

- Most software daemons, (such as background processes in the UNIX operating system), which monitor a software environment and perform actions to modify it, can be viewed as agents. An example is the X Windows program `xbiff`. This utility continually monitors a user’s incoming email, and indicates via a GUI icon whether or not they have unread messages. Whereas our thermostat agent in the previous example inhabited a *physical* environment

– the physical world – the `xbiff` program inhabits a *software* environment. It obtains information about this environment by carrying out software functions (by executing system programs such as `ls`, for example), and the actions it performs are software actions (changing an icon on the screen, or executing a program). The decision making component is just as simple as our thermostat example.

To summarize, agents are simply computer systems that are capable of autonomous action in some environment in order to meet their design objectives. An agent will typically sense its environment (by physical sensors in the case of agents situated in part of the real world, or by software sensors in the case of software agents), and will have available a repertoire of actions that can be executed to modify the environment, which may appear to respond non-deterministically to the execution of these actions.

2.2 Intelligent Agents

We are not used to thinking of thermostats or UNIX daemons as agents, and certainly not as *intelligent* agents. So, when do we consider an agent to be intelligent? The question, like the question *what is intelligence?* itself, is not an easy one to answer. But for the purposes of this chapter, an intelligent agent is one that is capable of *flexible* autonomous action in order to meet its design objectives, where flexibility means three things [55]:

- *reactivity*: intelligent agents are able to perceive their environment, and respond in a timely fashion to changes that occur in it in order to satisfy their design objectives;
- *pro-activeness*: intelligent agents are able to exhibit goal-directed behavior by *taking the initiative* in order to satisfy their design objectives;
- *social ability*: intelligent agents are capable of interacting with other agents (and possibly humans) in order to satisfy their design objectives.

These properties are more demanding than they might at first appear. To see why, let us consider them in turn. First, consider *pro-activeness*: goal directed behavior. It is not hard to build a system that exhibits goal directed behavior – we do it every time we write a procedure in PASCAL, a function in C, or a method in JAVA. When we write such a procedure, we describe it in terms of the *assumptions* on which it relies (formally, its *pre-condition*) and the *effect* it has if the assumptions are valid (its *post-condition*). The effects of the procedure are its *goal*: what the author of the software intends the procedure to achieve. If the pre-condition holds when the procedure is invoked, then we expect that the procedure will execute *correctly*: that it will terminate, and that upon termination, the post-condition will be true, i.e., the goal will be achieved. This is goal directed behavior: the procedure is simply a plan or recipe for achieving the goal. This programming model is fine for many environments. For example, it works well when we consider *functional systems* – those that simply take some input x , and

produce as output some some function $f(x)$ of this input. Compilers are a classic example of functional systems.

But for non-functional systems, this simple model of goal directed programming is not acceptable, as it makes some important limiting assumptions. In particular, it assumes that the environment *does not change* while the procedure is executing. If the environment does change, and in particular, if the assumptions (pre-condition) underlying the procedure become false while the procedure is executing, then the behavior of the procedure may not be defined – often, it will simply crash. Also, it is assumed that the goal, that is, the reason for executing the procedure, remains valid at least until the procedure terminates. If the goal does *not* remain valid, then there is simply no reason to continue executing the procedure.

In many environments, neither of these assumptions are valid. In particular, in domains that are *too complex* for an agent to observe completely, that are *multi-agent* (i.e., they are populated with more than one agent that can change the environment), or where there is *uncertainty* in the environment, these assumptions are not reasonable. In such environments, blindly executing a procedure without regard to whether the assumptions underpinning the procedure are valid is a poor strategy. In such dynamic environments, an agent must be *reactive*, in just the way that we described above. That is, it must be responsive to events that occur in its environment, where these events affect either the agent's goals or the assumptions which underpin the procedures that the agent is executing in order to achieve its goals.

As we have seen, building purely goal directed systems is not hard. As we shall see later in this chapter, building *purely reactive* systems – ones that *continually* respond to their environment – is also not difficult. However, what turns out to be hard is building a system that achieves an effective *balance* between goal-directed and reactive behavior. We want agents that will attempt to achieve their goals systematically, perhaps by making use of complex procedure-like patterns of action. But we don't want our agents to continue blindly executing these procedures in an attempt to achieve a goal either when it is clear that the procedure will not work, or when the goal is for some reason no longer valid. In such circumstances, we want our agent to be able to react to the new situation, in time for the reaction to be of some use. However, we do not want our agent to be *continually* reacting, and hence never focussing on a goal long enough to actually achieve it.

On reflection, it should come as little surprise that achieving a good balance between goal directed and reactive behavior is hard. After all, it is comparatively rare to find humans that do this very well. How many of us have had a manager who stayed blindly focussed on some project long after the relevance of the project was passed, or it was clear that the project plan was doomed to failure? Similarly, how many have encountered managers who seem unable to stay focussed at all, who flit from one project to another without ever managing to pursue a goal long enough to achieve *anything*? This problem – of effectively integrating goal-directed and reactive behavior – is one of the key problems fac-

ing the agent designer. As we shall see, a great many proposals have been made for how to build agents that can do this – but the problem is essentially still open.

Finally, let us say something about *social ability*, the final component of flexible autonomous action as defined here. In one sense, social ability is trivial: every day, millions of computers across the world routinely exchange information with both humans and other computers. But the ability to exchange bit streams is not really social ability. Consider that in the human world, comparatively few of our meaningful goals can be achieved without the *cooperation* of other people, who cannot be assumed to *share* our goals – in other words, they are themselves autonomous, with their own agenda to pursue. To achieve our goals in such situations, we must *negotiate* and *cooperate* with others. We may be required to understand and reason about the goals of others, and to perform actions (such as paying them money) that we would not otherwise choose to perform, in order to get them to cooperate with us, and achieve our goals. This type of social ability is much more complex, and much less well understood, than simply the ability to exchange binary information. Social ability in general (and topics such as negotiation and cooperation in particular) are dealt with elsewhere in this book, and will not therefore be considered here. In this chapter, we will be concerned with the decision making of *individual* intelligent agents in environments which may be dynamic, unpredictable, and uncertain, but do not contain other agents.

Sources and Further Reading. A view of artificial intelligence as the process of agent design is presented in [49], and in particular, Chapter 2 of [49] presents much useful material. The definition of agents presented here is based on [55], which also contains an extensive review of agent architectures and programming languages. In addition, [55] contains a detailed survey of *agent theories* – formalisms for reasoning about intelligent, rational agents, which is outside the scope of this chapter. This question of “what is an agent” is one that continues to generate some debate; a collection of answers may be found in [39]. The relationship between agents and objects has not been widely discussed in the literature, but see [20]. Other readable introductions to the idea of intelligent agents include [26] and [11].

3 Abstract Architectures for Intelligent Agents

We can easily formalize the abstract view of agents presented so far. First, we will assume that the state of the agent’s environment can be characterized as a set $S = \{s_1, s_2, \dots\}$ of *environment states*. At any given instant, the environment is assumed to be in one of these states. The effectoric capability of an agent is assumed to be represented by a set $A = \{a_1, a_2, \dots\}$ of *actions*. Then abstractly, an agent can be viewed as a function

$$\text{action} : S^* \rightarrow A$$

which maps sequences of environment states to actions. We will refer to an agent modelled by a function of this form as a *standard agent*. The intuition is that an agent decides what action to perform on the basis of its history – its experiences to date. These experiences are represented as a sequence of environment states – those that the agent has thus far encountered.

The (non-deterministic) behavior of an an environment can be modelled as a function

$$env : S \times A \rightarrow \wp(S)$$

which takes the current state of the environment $s \in S$ and an action $a \in A$ (performed by the agent), and maps them to a set of environment states $env(s, a)$ – those that could result from performing action a in state s . If all the sets in the range of env are all singletons, (i.e., if the result of performing any action in any state is a set containing a single member), then the environment is *deterministic*, and its behavior can be accurately predicted.

We can represent the interaction of agent and environment as a *history*. A history h is a sequence:

$$h : s_0 \xrightarrow{a_0} s_1 \xrightarrow{a_1} s_2 \xrightarrow{a_2} s_3 \xrightarrow{a_3} \dots \xrightarrow{a_{u-1}} s_u \xrightarrow{a_u} \dots$$

where s_0 is the initial state of the environment (i.e., its state when the agent starts executing), a_u is the u 'th action that the agent chose to perform, and s_u is the u 'th environment state (which is one of the possible results of executing action a_{u-1} in state s_{u-1}). If $action : S^* \rightarrow A$ is an agent, $env : S \times A \rightarrow \wp(S)$ is an environment, and s_0 is the initial state of the environment, then the sequence

$$h : s_0 \xrightarrow{a_0} s_1 \xrightarrow{a_1} s_2 \xrightarrow{a_2} s_3 \xrightarrow{a_3} \dots \xrightarrow{a_{u-1}} s_u \xrightarrow{a_u} \dots$$

will represent a possible history of the agent in the environment iff the following two conditions hold:

$$\forall u \in \mathbb{N}, a_u = action((s_0, s_1, \dots, s_u))$$

and

$$\forall u \in \mathbb{N} \text{ such that } u > 0, s_u \in env(s_{u-1}, a_{u-1}).$$

The *characteristic behavior* of an agent $action : S^* \rightarrow A$ in an environment $env : S \times A \rightarrow \wp(S)$ is the set of all the histories that satisfy these properties. If some property ϕ holds of all these histories, this property can be regarded as an invariant property of the agent in the environment. For example, if our agent is a nuclear reactor controller, (i.e., the environment is a nuclear reactor), and in all possible histories of the controller/reactor, the reactor does not blow up, then this can be regarded as a (desirable) invariant property. We will denote by $hist(agent, environment)$ the set of all histories of *agent* in *environment*. Two agents ag_1 and ag_2 are said to be *behaviorally equivalent* with respect to environment env iff $hist(ag_1, env) = hist(ag_2, env)$, and simply behaviorally equivalent iff they are behaviorally equivalent with respect to all environments.

In general, we are interested in agents whose interaction with their environment *does not end*, i.e., they are *non-terminating*. In such cases, the histories that we consider will be infinite.

3.1 Purely Reactive Agents

Certain types of agents decide what to do without reference to their history. They base their decision making entirely on the present, with no reference at all to the past. We will call such agents *purely reactive*, since they simply respond directly to their environment. Formally, the behavior of a purely reactive agent can be represented by a function

$$action : S \rightarrow A.$$

It should be easy to see that for every purely reactive agent, there is an equivalent standard agent; the reverse, however, is not generally the case.

Our thermostat agent is an example of a purely reactive agent. Assume, without loss of generality, that the thermostat's environment can be in one of two states – either too cold, or temperature OK. Then the thermostat's action function is simply

$$action(s) = \begin{cases} \text{heater off if } s = \text{temperature OK} \\ \text{heater on otherwise.} \end{cases}$$

3.2 Perception

Viewing agents at this abstract level makes for a pleasantly simple analysis. However, it does not help us to construct them, since it gives us no clues about how to design the decision function *action*. For this reason, we will now begin to *refine* our abstract model of agents, by breaking it down into sub-systems in exactly the way that one does in standard software engineering. As we refine our view of agents, we find ourselves making *design choices* that mostly relate to the subsystems that go to make up an agent – what data and control structures will be present. An *agent architecture* is essentially a map of the internals of an agent – its data structures, the operations that may be performed on these data structures, and the control flow between these data structures. Later in this chapter, we will discuss a number of different types of agent architecture, with very different views on the data structures and algorithms that will be present within an agent. In the remainder of this section, however, we will survey some fairly high-level design decisions. The first of these is the separation of an agent's decision function into *perception* and *action* subsystems: see Fig. 2.

The idea is that the function *see* captures the agent's ability to observe its environment, whereas the *action* function represents the agent's decision making process. The *see* function might be implemented in hardware in the case of an agent situated in the physical world: for example, it might be a video camera or an infra-red sensor on a mobile robot. For a software agent, the sensors might be system commands that obtain information about the software environment, such as `ls`, `finger`, or suchlike. The *output* of the *see* function is a *percept* – a perceptual input. Let P be a (non-empty) set of percepts. Then *see* is a function

$$see : S \rightarrow P$$

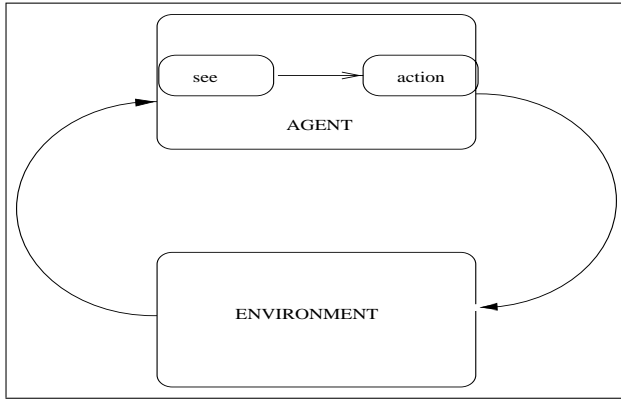


Fig. 2. Perception and action subsystems.

which maps environment states to percepts, and *action* is now a function

$$action : P^* \rightarrow A$$

which maps sequences of percepts to actions.

These simple definitions allow us to explore some interesting properties of agents and perception. Suppose that we have two environment states, $s_1 \in S$ and $s_2 \in S$, such that $s_1 \neq s_2$, but $see(s_1) = see(s_2)$. Then two *different* environment states are mapped to the *same* percept, and hence the agent would receive the same perceptual information from different environment states. As far as the agent is concerned, therefore, s_1 and s_2 are *indistinguishable*. To make this example concrete, let us return to the thermostat example. Let x represent the statement

“the room temperature is OK”

and let y represent the statement

“John Major is Prime Minister”.

If these are the only two facts about our environment that we are concerned with, then the set S of environment states contains exactly four elements:

$$S = \underbrace{\{\{-x, \neg y\}\}}_{s_1}, \underbrace{\{\{-x, y\}\}}_{s_2}, \underbrace{\{\{x, \neg y\}\}}_{s_3}, \underbrace{\{\{x, y\}\}}_{s_4}$$

Thus in state s_1 , the room temperature is not OK, and John Major is not Prime Minister; in state s_2 , the room temperature is not OK, and John Major *is* Prime Minister. Now, our thermostat is sensitive *only* to temperatures in the room. This room temperature is not causally related to whether or not John Major is Prime Minister. Thus the states where John Major is and is not

Prime Minister are literally *indistinguishable* to the thermostat. Formally, the *see* function for the thermostat would have two percepts in its range, p_1 and p_2 , indicating that the temperature is too cold or OK respectively. The *see* function for the thermostat would behave as follows:

$$see(s) = \begin{cases} p_1 & \text{if } s = s_1 \text{ or } s = s_2 \\ p_2 & \text{if } s = s_3 \text{ or } s = s_4. \end{cases}$$

Given two environment states $s \in S$ and $s' \in S$, let us write $s \equiv s'$ if $see(s) = see(s')$. It is not hard to see that \equiv is an *equivalence relation* over environment states, which partitions S into mutually indistinguishable sets of states. Intuitively, the coarser these equivalence classes are, the less effective is the agent's perception. If $|\equiv| = |S|$, (i.e., the number of distinct percepts is equal to the number of different environment states), then the agent can distinguish *every* state – the agent has perfect perception in the environment; it is *omniscient*. At the other extreme, if $|\equiv| = 1$, then the agent's perceptual ability is non-existent – it cannot distinguish between *any* different states. In this case, as far as the agent is concerned, all environment states are identical.

3.3 Agents with State

We have so far been modelling an agent's decision function *action* as from *sequences* of environment states or percepts to actions. This allows us to represent agents whose decision making is influenced by history. However, this is a somewhat unintuitive representation, and we shall now replace it by an equivalent, but somewhat more natural scheme. The idea is that we now consider agents that *maintain state* – see Fig. 3.

These agents have some internal data structure, which is typically used to record information about the environment state and history. Let I be the set

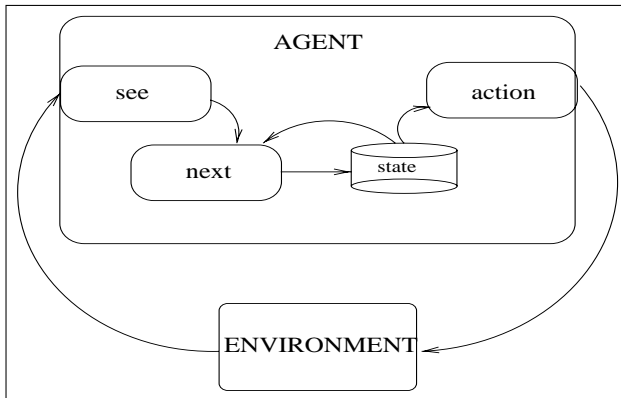


Fig. 3. Agents that maintain state.

of all internal states of the agent. An agent's decision making process is then based, at least in part, on this information. The perception function *see* for a state-based agent is unchanged, mapping environment states to percepts as before:

$$see : S \rightarrow P$$

The action-selection function *action* is now defined as a mapping

$$action : I \rightarrow A$$

from internal states to actions. An additional function *next* is introduced, which maps an internal state and percept to an internal state:

$$next : I \times P \rightarrow I$$

The behavior of a state-based agent can be summarized as follows. The agent starts in some initial internal state i_0 . It then observes its environment state s , and generates a percept $see(s)$. The internal state of the agent is then updated via the *next* function, becoming set to $next(i_0, see(s))$. The action selected by the agent is then $action(next(i_0, see(s)))$. This action is then performed, and the agent enters another cycle, perceiving the world via *see*, updating its state via *next*, and choosing an action to perform via *action*.

It is worth observing that state-based agents as defined here are in fact no more powerful than the standard agents we introduced earlier. In fact, they are *identical* in their expressive power – every state-based agent can be transformed into a standard agent that is behaviorally equivalent.

Sources and Further Reading. The abstract model of agents presented here is based on that given in [21, Chapter 13], and also makes use of some ideas from [51,50]. The properties of perception as discussed in this section lead to *knowledge theory*, a formal analysis of the information implicit within the state of computer processes, which has had a profound effect in theoretical computer science. The definitive reference is [12], and an introductory survey is [25].

4 Concrete Architectures for Intelligent Agents

Thus far, we have considered agents only in the abstract. So while we have examined the properties of agents that do and do not maintain state, we have not stopped to consider what this state might look like. Similarly, we have modelled an agent's decision making as an abstract function *action*, which somehow manages to indicate which action to perform – but we have not discussed how this function might be implemented. In this section, we will rectify this omission. We will consider four classes of agents:

- *logic based agents* – in which decision making is realized through logical deduction;

- *reactive agents* – in which decision making is implemented in some form of direct mapping from situation to action;
- *belief-desire-intention agents* – in which decision making depends upon the manipulation of data structures representing the beliefs, desires, and intentions of the agent; and finally,
- *layered architectures* – in which decision making is realized via various software layers, each of which is more-or-less explicitly reasoning about the environment at different levels of abstraction.

In each of these cases, we are moving away from the abstract view of agents, and beginning to make quite specific commitments about the internal structure and operation of agents. Each section explains the nature of these commitments, the assumptions upon which the architectures depend, and the relative advantages and disadvantages of each.

4.1 Logic-Based Architectures

The “traditional” approach to building artificially intelligent systems, (known as *symbolic AI*) suggests that intelligent behavior can be generated in a system by giving that system a *symbolic* representation of its environment and its desired behavior, and syntactically manipulating this representation. In this section, we focus on the apotheosis of this tradition, in which these symbolic representations are *logical formulae*, and the syntactic manipulation corresponds to *logical deduction*, or *theorem proving*.

The idea of agents as theorem provers is seductive. Suppose we have some theory of agency – some theory that explains how an intelligent agent should behave. This theory might explain, for example, how an agent generates goals so as to satisfy its design objective, how it interleaves goal-directed and reactive behavior in order to achieve these goals, and so on. Then this theory ϕ can be considered as a *specification* for how an agent should behave. The traditional approach to implementing a system that will satisfy this specification would involve *refining* the specification through a series of progressively more concrete stages, until finally an implementation was reached. In the view of agents as theorem provers, however, no such refinement takes place. Instead, ϕ is viewed as an *executable specification*: it is *directly executed* in order to produce the agent’s behavior.

To see how such an idea might work, we shall develop a simple model of logic-based agents, which we shall call *deliberate* agents. In such agents, the internal state is assumed to be a database of formulae of classical first-order predicate logic. For example, the agent’s database might contain formulae such as:

Open(valve221)
Temperature(reactor4726, 321)
Pressure(tank776, 28)

It is not difficult to see how formulae such as these can be used to represent the properties of some environment. The database is the *information* that the agent

has about its environment. An agent's database plays a somewhat analogous role to that of *belief* in humans. Thus a person might have a belief that valve 221 is open – the agent might have the predicate $Open(valve221)$ in its database. Of course, just like humans, agents can be wrong. Thus I might believe that valve 221 is open when it is in fact closed; the fact that an agent has $Open(valve221)$ in its database does not mean that valve 221 (or indeed any valve) is open. The agent's sensors may be faulty, its reasoning may be faulty, the information may be out of date, or the interpretation of the formula $Open(valve221)$ intended by the agent's designer may be something entirely different.

Let L be the set of sentences of classical first-order logic, and let $D = \wp(L)$ be the set of L databases, i.e., the set of sets of L -formulae. The internal state of an agent is then an element of D . We write Δ, Δ_1, \dots for members of D . The internal state of an agent is then simply a member of the set D . An agent's decision making process is modelled through a set of *deduction rules*, ρ . These are simply rules of inference for the logic. We write $\Delta \vdash_\rho \phi$ if the formula ϕ can be proved from the database Δ using only the deduction rules ρ . An agent's perception function *see* remains unchanged:

$$see : S \rightarrow P.$$

Similarly, our *next* function has the form

$$next : D \times P \rightarrow D$$

It thus maps a database and a percept to a new database. However, an agent's action selection function, which has the signature

$$action : D \rightarrow A$$

is defined in terms of its deduction rules. The pseudo-code definition of this function is as follows.

1. function $action(\Delta : D) : A$
2. begin
3. for each $a \in A$ do
4. if $\Delta \vdash_\rho Do(a)$ then
5. return a
6. end-if
7. end-for
8. for each $a \in A$ do
9. if $\Delta \not\vdash_\rho \neg Do(a)$ then
10. return a
11. end-if
12. end-for
13. return *null*
14. end function $action$

The idea is that the agent programmer will encode the deduction rules ρ and database Δ in such a way that if a formula $Do(a)$ can be derived, where a is a term that denotes an action, then a is the best action to perform. Thus, in the first part of the function (lines (3)–(7)), the agent takes each of its possible actions a in turn, and attempts to prove the form the formula $Do(a)$ from its database (passed as a parameter to the function) using its deduction rules ρ . If the agent succeeds in proving $Do(a)$, then a is returned as the action to be performed.

What happens if the agent fails to prove $Do(a)$, for all actions $a \in A$? In this case, it attempts to find an action that is *consistent* with the rules and database, i.e., one that is not explicitly forbidden. In lines (8)–(12), therefore, the agent attempts to find an action $a \in A$ such that $\neg Do(a)$ cannot be derived from its database using its deduction rules. If it can find such an action, then this is returned as the action to be performed. If, however, the agent fails to find an action that is at least consistent, then it returns a special action *null* (or *noop*), indicating that no action has been selected.

In this way, the agent’s behavior is determined by the agent’s deduction rules (its “program”) and its current database (representing the information the agent has about its environment).

To illustrate these ideas, let us consider a small example (based on the vacuum cleaning world example of [49, p51]). The idea is that we have a small robotic agent that will clean up a house. The robot is equipped with a sensor that will tell it whether it is over any dirt, and a vacuum cleaner that can be used to suck up dirt. In addition, the robot always has a definite orientation (one of *north*, *south*, *east*, or *west*). In addition to being able to suck up dirt, the agent can move forward one “step” or turn right 90° . The agent moves around a room, which is divided grid-like into a number of equally sized squares (conveniently corresponding to the unit of movement of the agent). We will assume that our agent does nothing but clean – it never leaves the room, and further, we will assume in the interests of simplicity that the room is a 3×3 grid, and the agent always starts in grid square $(0, 0)$ facing north.

To summarize, our agent can receive a percept *dirt* (signifying that there is dirt beneath it), or *null* (indicating no special information). It can perform any one of three possible actions: *forward*, *suck*, or *turn*. The goal is to traverse the room continually searching for and removing dirt. See Fig. 4 for an illustration of the vacuum world.

First, note that we make use of three simple *domain predicates* in this exercise:

$In(x, y)$ agent is at (x, y)
 $Dirt(x, y)$ there is dirt at (x, y)
 $Facing(d)$ the agent is facing direction d

Now we specify our *next* function. This function must look at the perceptual information obtained from the environment (either *dirt* or *null*), and generate a new database which includes this information. But in addition, it must *remove* old or irrelevant information, and also, it must try to figure out the new location

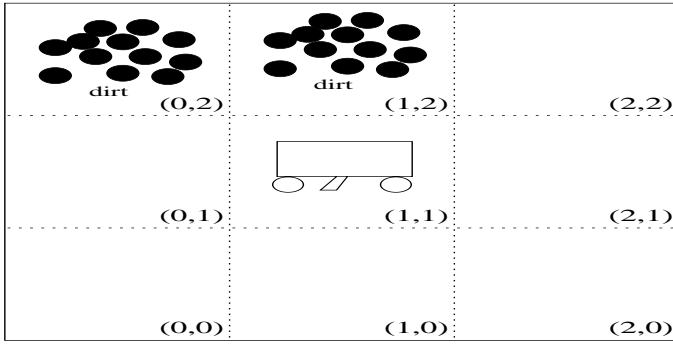


Fig. 4. Vacuum world

and orientation of the agent. We will therefore specify the *next* function in several parts. First, let us write $old(\Delta)$ to denote the set of “old” information in a database, which we want the update function *next* to remove:

$$old(\Delta) = \{P(t_1, \dots, t_n) \mid P \in \{In, Dirt, Facing\} \text{ and } P(t_1, \dots, t_n) \in \Delta\}$$

Next, we require a function *new*, which gives the set of new predicates to add to the database. This function has the signature

$$new : D \times P \rightarrow D$$

The definition of this function is not difficult, but it is rather lengthy, and so we will leave it as an exercise. (It must generate the predicates $In(\dots)$, describing the new position of the agent, $Facing(\dots)$ describing the orientation of the agent, and $Dirt(\dots)$ if dirt has been detected at the new position.) Given the *new* and *old* functions, the *next* function is defined as follows:

$$next(\Delta, p) = (\Delta \setminus old(\Delta)) \cup new(\Delta, p)$$

Now we can move on to the rules that govern our agent’s behavior. The deduction rules have the form

$$\phi(\dots) \longrightarrow \psi(\dots)$$

where ϕ and ψ are predicates over some arbitrary list of constants and variables. The idea being that if ϕ matches against the agent’s database, then ψ can be concluded, with any variables in ψ instantiated.

The first rule deals with the basic cleaning action of the agent: this rule will take priority over all other possible behaviors of the agent (such as navigation).

$$In(x, y) \wedge Dirt(x, y) \longrightarrow Do(suck) \tag{1}$$

Hence if the agent is at location (x, y) and it perceives dirt, then the prescribed action will be to suck up dirt. Otherwise, the basic action of the agent will be to

traverse the world. Taking advantage of the simplicity of our environment, we will hardwire the basic navigation algorithm, so that the robot will always move from $(0, 0)$ to $(0, 1)$ to $(0, 2)$ and then to $(1, 2)$, $(1, 1)$ and so on. Once the agent reaches $(2, 2)$, it must head back to $(0, 0)$. The rules dealing with the traversal up to $(0, 2)$ are very simple.

$$\text{In}(0, 0) \wedge \text{Facing}(\text{north}) \wedge \neg \text{Dirt}(0, 0) \longrightarrow \text{Do}(\text{forward}) \quad (2)$$

$$\text{In}(0, 1) \wedge \text{Facing}(\text{north}) \wedge \neg \text{Dirt}(0, 1) \longrightarrow \text{Do}(\text{forward}) \quad (3)$$

$$\text{In}(0, 2) \wedge \text{Facing}(\text{north}) \wedge \neg \text{Dirt}(0, 2) \longrightarrow \text{Do}(\text{turn}) \quad (4)$$

$$\text{In}(0, 2) \wedge \text{Facing}(\text{east}) \longrightarrow \text{Do}(\text{forward}) \quad (5)$$

Notice that in each rule, we must explicitly check whether the antecedent of rule (1) fires. This is to ensure that we only ever prescribe one action via the $\text{Do}(\dots)$ predicate. Similar rules can easily be generated that will get the agent to $(2, 2)$, and once at $(2, 2)$ back to $(0, 0)$. It is not difficult to see that these rules, together with the *next* function, will generate the required behavior of our agent.

At this point, it is worth stepping back and examining the pragmatics of the logic-based approach to building agents. Probably the most important point to make is that a literal, naive attempt to build agents in this way would be more or less entirely impractical. To see why, suppose we have designed our agent's rule set ρ such that for any database Δ , if we can prove $\text{Do}(a)$ then a is an *optimal* action – that is, a is the best action that could be performed when the environment is as described in Δ . Then imagine we start running our agent. At time t_1 , the agent has generated some database Δ_1 , and begins to apply its rules ρ in order to find which action to perform. Some time later, at time t_2 , it manages to establish $\Delta_1 \vdash_{\rho} \text{Do}(a)$ for some $a \in A$, and so a is the optimal action that the agent could perform at time t_1 . But if the environment has *changed* between t_1 and t_2 , then there is no guarantee that a will *still* be optimal. It could be far from optimal, particularly if much time has elapsed between t_1 and t_2 . If $t_2 - t_1$ is infinitesimal – that is, if decision making is effectively instantaneous – then we could safely disregard this problem. But in fact, we know that reasoning of the kind our logic-based agents use will be anything *but* instantaneous. (If our agent uses classical first-order predicate logic to represent the environment, and its rules are sound and complete, then there is no guarantee that the decision making procedure will even *terminate*.) An agent is said to enjoy the property of *calculative rationality* if and only if its decision making apparatus will suggest an action that was optimal *when the decision making process began*. Calculative rationality is clearly not acceptable in environments that change faster than the agent can make decisions – we shall return to this point later.

One might argue that this problem is an artifact of the pure logic-based approach adopted here. There is an element of truth in this. By moving away from strictly logical representation languages and complete sets of deduction rules, one can build agents that enjoy respectable performance. But one also loses what is arguably the greatest advantage that the logical approach brings: a simple, elegant logical semantics.

There are several other problems associated with the logical approach to agency. First, the *see* function of an agent, (its perception component), maps its environment to a percept. In the case of a logic-based agent, this percept is likely to be symbolic – typically, a set of formulae in the agent’s representation language. But for many environments, it is not obvious how the mapping from environment to symbolic percept might be realized. For example, the problem of transforming an image to a set of declarative statements representing that image has been the object of study in AI for decades, and is still essentially open. Another problem is that actually *representing* properties of dynamic, real-world environments is extremely hard. As an example, representing and reasoning about *temporal information* – how a situation changes over time – turns out to be extraordinarily difficult. Finally, as the simple vacuum world example illustrates, representing even rather simple *procedural* knowledge (i.e., knowledge about “what to do”) in traditional logic can be rather unintuitive and cumbersome.

To summarize, in logic-based approaches to building agents, decision making is viewed as deduction. An agent’s “program” – that is, its decision making strategy – is encoded as a logical theory, and the process of selecting an action reduces to a problem of proof. Logic-based approaches are elegant, and have a clean (logical) semantics – wherein lies much of their long-lived appeal. But logic-based approaches have many disadvantages. In particular, the inherent computational complexity of theorem proving makes it questionable whether agents as theorem provers can operate effectively in time-constrained environments. Decision making in such agents is predicated on the assumption of calculative rationality – the assumption that the world will not change in any significant way while the agent is deciding what to do, and that an action which is rational when decision making begins will be rational when it concludes. The problems associated with representing and reasoning about complex, dynamic, possibly physical environments are also essentially unsolved.

Sources and Further Reading. My presentation of logic based agents is based largely on the discussion of *deliberate agents* presented in [21, Chapter 13], which represents the logic-centric view of AI and agents very well. The discussion is also partly based on [30]. A number of more-or-less “pure” logical approaches to agent programming have been developed. Well-known examples include the CONGOLOG system of Lespérance and colleagues [31] (which is based on the *situation calculus* [36]) and the METATEM and Concurrent METATEM programming languages developed by Fisher and colleagues [3,19] (in which agents are programmed by giving them *temporal logic* specifications of the behavior they should exhibit). Note that these architectures (and the discussion above) assume that if one adopts a logical approach to agent-building, then this means agents are essentially theorem provers, employing explicit symbolic reasoning (theorem proving) in order to make decisions. But just because we find logic a useful tool for conceptualising or specifying agents, this does not mean that we must view decision-making as logical manipulation. An alternative is to *compile* the logical

specification of an agent into a form more amenable to efficient decision making. The difference is rather like the distinction between interpreted and compiled programming languages. The best-known example of this work is the *situated automata* paradigm of Leslie Kaelbling and Stanley Rosenschein [48]. A review of the role of logic in intelligent agents may be found in [54]. Finally, for a detailed discussion of calculative rationality and the way that it has affected thinking in AI, see [50].

4.2 Reactive Architectures

The seemingly intractable problems with symbolic/logical approaches to building agents led some researchers to question, and ultimately reject, the assumptions upon which such approaches are based. These researchers have argued that minor changes to the symbolic approach, such as weakening the logical representation language, will not be sufficient to build agents that can operate in time-constrained environments: nothing less than a whole new approach is required. In the mid-to-late 1980s, these researchers began to investigate alternatives to the symbolic AI paradigm. It is difficult to neatly characterize these different approaches, since their advocates are united mainly by a rejection of symbolic AI, rather than by a common manifesto. However, certain themes do recur:

- the rejection of symbolic representations, and of decision making based on syntactic manipulation of such representations;
- the idea that intelligent, rational behavior is seen as innately linked to the *environment* an agent occupies – intelligent behavior is not disembodied, but is a product of the *interaction* the agent maintains with its environment;
- the idea that intelligent behavior *emerges* from the interaction of various simpler behaviors.

Alternative approaches to agency are sometime referred to as *behavioral* (since a common theme is that of developing and combining individual behaviors), *situated* (since a common theme is that of agents actually situated in some environment, rather than being disembodied from it), and finally – the term used in this chapter – *reactive* (because such systems are often perceived as simply reacting to an environment, without reasoning about it). This section presents a survey of the *subsumption architecture*, which is arguably the best-known reactive agent architecture. It was developed by Rodney Brooks – one of the most vocal and influential critics of the symbolic approach to agency to have emerged in recent years.

There are two defining characteristics of the subsumption architecture. The first is that an agent’s decision-making is realized through a set of *task accomplishing behaviors*; each behavior may be thought of as an individual *action* function, as we defined above, which continually takes perceptual input and maps it to an action to perform. Each of these behavior modules is intended to achieve some particular task. In Brooks’ implementation, the behavior modules

are finite state machines. An important point to note is that these task accomplishing modules are assumed to include *no* complex symbolic representations, and are assumed to do *no* symbolic reasoning at all. In many implementations, these behaviors are implemented as rules of the form

$$\text{situation} \longrightarrow \text{action}$$

which simply map perceptual input directly to actions.

The second defining characteristic of the subsumption architecture is that many behaviors can “fire” simultaneously. There must obviously be a mechanism to choose between the different actions selected by these multiple actions. Brooks proposed arranging the modules into a *subsumption hierarchy*, with the behaviors arranged into *layers*. Lower layers in the hierarchy are able to *inhibit* higher layers: the lower a layer is, the higher is its priority. The idea is that higher layers represent more abstract behaviors. For example, one might desire a behavior in a mobile robot for the behavior “avoid obstacles”. It makes sense to give obstacle avoidance a high priority – hence this behavior will typically be encoded in a *low-level* layer, which has *high* priority. To illustrate the subsumption architecture in more detail, we will now present a simple formal model of it, and illustrate how it works by means of a short example. We then discuss its relative advantages and shortcomings, and point at other similar reactive architectures.

The *see* function, which represents the agent’s perceptual ability, is assumed to remain unchanged. However, in implemented subsumption architecture systems, there is assumed to be quite tight coupling between perception and action – raw sensor input is not processed or transformed much, and there is certainly no attempt to transform images to symbolic representations.

The decision function *action* is realized through a set of behaviors, together with an *inhibition* relation holding between these behaviors. A behavior is a pair (c, a) , where $c \subseteq P$ is a set of percepts called the *condition*, and $a \in A$ is an action. A behavior (c, a) will *fire* when the environment is in state $s \in S$ iff $\text{see}(s) \in c$. Let $\text{Beh} = \{(c, a) \mid c \subseteq P \text{ and } a \in A\}$ be the set of all such rules.

Associated with an agent’s set of behavior rules $R \subseteq \text{Beh}$ is a binary *inhibition relation* on the set of behaviors: $\prec \subseteq R \times R$. This relation is assumed to be a total ordering on R (i.e., it is transitive, irreflexive, and antisymmetric). We write $b_1 \prec b_2$ if $(b_1, b_2) \in \prec$, and read this as “ b_1 inhibits b_2 ”, that is, b_1 is lower in the hierarchy than b_2 , and will hence get priority over b_2 . The action function is then defined as follows:

1. function $\text{action}(p : P) : A$
2. var $\text{fired} : \wp(R)$
3. var $\text{selected} : A$
4. begin
5. $\text{fired} := \{(c, a) \mid (c, a) \in R \text{ and } p \in c\}$
6. for each $(c, a) \in \text{fired}$ do
7. if $\neg(\exists(c', a') \in \text{fired} \text{ such that } (c', a') \prec (c, a))$ then
8. return a

9. end-if
10. end-for
11. return *null*
12. end function *action*

Thus action selection begins by first computing the set *fired* of all behaviors that fire (5). Then, each behavior (c, a) that fires is checked, to determine whether there is some other higher priority behavior that fires. If not, then the action part of the behavior, a , is returned as the selected action (8). If no behavior fires, then the distinguished action *null* will be returned, indicating that no action has been chosen.

Given that one of our main concerns with logic-based decision making was its theoretical complexity, it is worth pausing to examine how well our simple behavior-based system performs. The overall time complexity of the subsumption action function is no worse than $O(n^2)$, where n is the larger of the number of behaviors or number of percepts. Thus, even with the naive algorithm above, decision making is tractable. In practice, we can do *considerably* better than this: the decision making logic can be encoded into hardware, giving *constant* decision time. For modern hardware, this means that an agent can be guaranteed to select an action within nano-seconds. Perhaps more than anything else, this computational simplicity is the strength of the subsumption architecture.

To illustrate how the subsumption architecture in more detail, we will show how subsumption architecture agents were built for the following scenario (this example is adapted from [53]):

The objective is to explore a distant planet, more concretely, to collect samples of a particular type of precious rock. The location of the rock samples is not known in advance, but they are typically clustered in certain spots. A number of autonomous vehicles are available that can drive around the planet collecting samples and later reenter the a mothership spacecraft to go back to earth. There is no detailed map of the planet available, although it is known that the terrain is full of obstacles – hills, valleys, etc. – which prevent the vehicles from exchanging any communication.

The problem we are faced with is that of building an agent control architecture for each vehicle, so that they will cooperate to collect rock samples from the planet surface as efficiently as possible. Luc Steels argues that logic-based agents, of the type we described above, are “entirely unrealistic” for this problem [53]. Instead, he proposes a solution using the subsumption architecture.

The solution makes use of two mechanisms introduced by Steels: The first is a *gradient field*. In order that agents can know in which direction the mothership lies, the mothership generates a radio signal. Now this signal will obviously weaken as distance to the source increases – to find the direction of the mothership, an agent need therefore only travel “up the gradient” of signal strength. The signal need not carry any information – it need only exist.

The second mechanism enables agents to communicate with one another. The characteristics of the terrain prevent direct communication (such as message passing), so Steels adopted an *indirect* communication method. The idea is that agents will carry “radioactive crumbs”, which can be dropped, picked up, and detected by passing robots. Thus if an agent drops some of these crumbs in a particular location, then later, another agent happening upon this location will be able to detect them. This simple mechanism enables a quite sophisticated form of cooperation.

The behavior of an individual agent is then built up from a number of behaviors, as we indicated above. First, we will see how agents can be programmed to *individually* collect samples. We will then see how agents can be programmed to generate a *cooperative* solution.

For individual (non-cooperative) agents, the lowest-level behavior, (and hence the behavior with the highest “priority”) is obstacle avoidance. This behavior can be represented in the rule:

$$\textit{if detect an obstacle then change direction.} \quad (6)$$

The second behavior ensures that any samples carried by agents are dropped back at the mother-ship.

$$\textit{if carrying samples and at the base then drop samples} \quad (7)$$

$$\textit{if carrying samples and not at the base then travel up gradient.} \quad (8)$$

Behavior (8) ensures that agents carrying samples will return to the mother-ship (by heading towards the origin of the gradient field). The next behavior ensures that agents will collect samples they find.

$$\textit{if detect a sample then pick sample up.} \quad (9)$$

The final behavior ensures that an agent with “nothing better to do” will explore randomly.

$$\textit{if true then move randomly.} \quad (10)$$

The pre-condition of this rule is thus assumed to always fire. These behaviors are arranged into the following hierarchy:

$$(6) \prec (7) \prec (8) \prec (9) \prec (10)$$

The subsumption hierarchy for this example ensures that, for example, an agent will *always* turn if any obstacles are detected; if the agent is at the mother-ship and is carrying samples, then it will *always* drop them if it is not in any immediate danger of crashing, and so on. The “top level” behavior – a random walk – will only ever be carried out if the agent has nothing more urgent to do. It is not difficult to see how this simple set of behaviors will solve the problem: agents will search for samples (ultimately by searching randomly), and when they find them, will return them to the mother-ship.

If the samples are distributed across the terrain entirely at random, then equipping a large number of robots with these very simple behaviors will work extremely well. But we know from the problem specification, above, that this is not the case: the samples tend to be located in clusters. In this case, it makes sense to have agents *cooperate* with one-another in order to find the samples. Thus when one agent finds a large sample, it would be helpful for it to communicate this to the other agents, so they can help it collect the rocks. Unfortunately, we also know from the problem specification that *direct* communication is impossible. Steels developed a simple solution to this problem, partly inspired by the foraging behavior of ants. The idea revolves around an agent creating a “trail” of radioactive crumbs whenever it finds a rock sample. The trail will be created when the agent returns the rock samples to the mother ship. If at some later point, another agent comes across this trail, then it need only follow it down the gradient field to locate the source of the rock samples. Some small refinements improve the efficiency of this ingenious scheme still further. First, as an agent follows a trail to the rock sample source, it picks up some of the crumbs it finds, hence making the trail fainter. Secondly, the trail is *only* laid by agents returning to the mothership. Hence if an agent follows the trail out to the source of the nominal rock sample only to find that it contains no samples, it will reduce the trail on the way out, and will not return with samples to reinforce it. After a few agents have followed the trail to find no sample at the end of it, the trail will in fact have been removed.

The modified behaviors for this example are as follows. Obstacle avoidance, (6), remains unchanged. However, the two rules determining what to do if carrying a sample are modified as follows.

if carrying samples *and* at the base *then* drop samples (11)

if carrying samples *and not* at the base
then drop 2 crumbs *and* travel up gradient. (12)

The behavior (12) requires an agent to drop crumbs when returning to base with a sample, thus either reinforcing or creating a trail. The “pick up sample” behavior, (9), remains unchanged. However, an additional behavior is required for dealing with crumbs.

if sense crumbs *then* pick up 1 crumb *and* travel down gradient (13)

Finally, the random movement behavior, (10), remains unchanged. These behaviors are then arranged into the following subsumption hierarchy.

(6) < (11) < (12) < (9) < (13) < (10)

Steels shows how this simple adjustment achieves near-optimal performance in many situations. Moreover, the solution is *cheap* (the computing power required by each agent is minimal) and *robust* (the loss of a single agent will not affect the overall system significantly).

In summary, there are obvious advantages to reactive approaches such as that Brooks' subsumption architecture: simplicity, economy, computational tractability, robustness against failure, and elegance all make such architectures appealing. But there are some fundamental, unsolved problems, not just with the subsumption architecture, but with other purely reactive architectures:

- If agents do not employ models of their environment, then they must have sufficient information available in their *local* environment for them to determine an acceptable action.
- Since purely reactive agents make decisions based on *local* information, (i.e., information about the agents *current* state), it is difficult to see how such decision making could take into account *non-local* information – it must inherently take a “short term” view.
- It is difficult to see how purely reactive agents can be designed that *learn* from experience, and improve their performance over time.
- A major selling point of purely reactive systems is that overall behavior *emerges* from the interaction of the component behaviors when the agent is placed in its environment. But the very term “emerges” suggests that the relationship between individual behaviors, environment, and overall behavior is not understandable. This necessarily makes it very hard to *engineer* agents to fulfill specific tasks. Ultimately, there is no principled *methodology* for building such agents: one must use a laborious process of experimentation, trial, and error to engineer an agent.
- While effective agents can be generated with small numbers of behaviors (typically less than ten layers), it is *much* harder to build agents that contain many layers. The dynamics of the interactions between the different behaviors become too complex to understand.

Various solutions to these problems have been proposed. One of the most popular of these is the idea of *evolving* agents to perform certain tasks. This area of work has largely broken away from the mainstream AI tradition in which work on, for example, logic-based agents is carried out, and is documented primarily in the *artificial life* (alife) literature.

Sources and Further Reading. Brooks' original paper on the subsumption architecture – the one that started all the fuss – was published as [7]. The description and discussion here is partly based on [13]. This original paper seems to be somewhat less radical than many of his later ones, which include [8,10,9]. The version of the subsumption architecture used in this chapter is actually a simplification of that presented by Brooks. The subsumption architecture is probably the best-known reactive architecture around – but there are many others. The collection of papers edited by Pattie Maes [33] contains papers that describe many of these, as does the collection by Agre and Rosenschein [2]. Other approaches include:

- the *agent network architecture* developed by Pattie Maes [32,34,35];

- Nilsson’s *teleo reactive programs* [40];
- Rosenchein and Kaelbling’s *situated automata* approach, which is particularly interesting in that it shows how agents can be *specified* in an abstract, logical framework, and *compiled* into equivalent, but computationally very simple machines [47,28,27,48];
- Agre and Chapman’s PENGU system [1];
- Schoppers’ *universal plans* – which are essentially decision trees that can be used to efficiently determine an appropriate action in any situation [52];
- Firby’s *reactive action packages* [17].

Kaelbling [26] gives a good discussion of the issues associated with developing resource-bounded rational agents, and proposes an agent architecture somewhat similar to that developed by Brooks.

4.3 Belief-Desire-Intention Architectures

In this section, we shall discuss *belief-desire-intention* (BDI) architectures. These architectures have their roots in the philosophical tradition of understanding *practical reasoning* – the process of deciding, moment by moment, which action to perform in the furtherance of our goals.

Practical reasoning involves two important processes: deciding *what* goals we want to achieve, and *how* we are going to achieve these goals. The former process is known as *deliberation*, the latter as *means-ends* reasoning. To gain an understanding of the BDI model, it is worth considering a simple example of practical reasoning. When you leave university with a first degree, you are faced with a decision to make – about what to do with your life. The decision process typically begins by trying to understand what the *options* available to you are. For example, if you gain a good first degree, then one option is that of becoming an academic. (If you fail to obtain a good degree, this option is not available to you.) Another option is entering industry. After generating this set of alternatives, you must *choose between them*, and *commit* to some. These chosen options become *intentions*, which then determine the agent’s actions. Intentions then feed back into the agent’s future practical reasoning. For example, if I decide I want to be an academic, then I should commit to this objective, and devote time and effort to bringing it about.

Intentions play a crucial role in the practical reasoning process. Perhaps the most obvious property of intentions is that they tend to lead to action. If I truly have an intention to become an academic, then you would expect me to *act* on that intention – to try to achieve it. For example, you might expect me to apply to various PhD programs. You would expect to make a *reasonable attempt* to achieve the intention. Thus you would expect me to carry out some course of action that I believed would best satisfy the intention. Moreover, if a course of action fails to achieve the intention, then you would expect me to *try again* – you would not expect me to simply give up. For example, if my first application for a PhD programme is rejected, then you might expect me to apply to alternative universities.

In addition, once I have adopted an intention, then the very fact of having this intention will constrain my future practical reasoning. For example, while I hold some particular intention, I will not entertain options that are inconsistent with that intention. Intending to become an academic, for example, would preclude the option of partying every night: the two are mutually exclusive.

Next, intentions *persist*. If I adopt an intention to become an academic, then I should *persist* with this intention and attempt to achieve it. For if I immediately drop my intentions without devoting resources to achieving them, then I will never achieve anything. However, I should not persist with my intention for too long – if it becomes clear to me that I will *never* become an academic, then it is only rational to drop my intention to do so. Similarly, if the reason for having an intention goes away, then it is rational of me to drop the intention. For example, if I adopted the intention to become an academic because I believed it would be an easy life, but then discover that I would be expected to actually *teach*, then the justification for the intention is no longer present, and I should drop the intention.

Finally, intentions are closely related to beliefs about the future. For example, if I intend to become an academic, then I should believe that I will indeed become an academic. For if I truly believe that I will never be an academic, it would be non-sensical of me to have an intention to become one. Thus if I intend to become an academic, I should at least believe that there is a good chance I will indeed become one.

From this discussion, we can see that intentions play a number of important roles in practical reasoning:

– *Intentions drive means-ends reasoning.*

If I have formed an intention to become an academic, then I will attempt to achieve the intention, which involves, amongst other things, deciding *how* to achieve it, for example, by applying for a PhD programme. Moreover, if one particular course of action fails to achieve an intention, then I will typically attempt others. Thus if I fail to gain a PhD place at one university, I might try another university.

– *Intentions constrain future deliberation.*

If I intend to become an academic, then I will not entertain options that are inconsistent with this intention. For example, a rational agent would not consider being rich as an option while simultaneously intending to be an academic. (While the two are not actually mutually exclusive, the probability of simultaneously achieving both is infinitesimal.)

– *Intentions persist.*

I will not usually give up on my intentions without good reason – they will persist, typically until either I believe I have successfully achieved them, I believe I cannot achieve them, or else because the purpose for the intention is no longer present.

– *Intentions influence beliefs upon which future practical reasoning is based.*

If I adopt the intention to become an academic, then I can plan for the future on the assumption that I *will* be an academic. For if I intend to be

an academic while simultaneously believing that I will never be one, then I am being irrational.

A key problem in the design of practical reasoning agents is that of achieving a good *balance* between these different concerns. Specifically, it seems clear that an agent should at times drop some intentions (because it comes to believe that either they will never be achieved, they are achieved, or else because the reason for having the intention is no longer present). It follows that, from time to time, it is worth an agent stopping to *reconsider* its intentions. But reconsideration has a cost – in terms of both time and computational resources. But this presents us with a dilemma:

- an agent that does not stop to reconsider sufficiently often will continue attempting to achieve its intentions even after it is clear that they cannot be achieved, or that there is no longer any reason for achieving them;
- an agent that *constantly* reconsiders its attentions may spend insufficient time actually working to achieve them, and hence runs the risk of never actually achieving them.

This dilemma is essentially the problem of balancing pro-active (goal directed) and reactive (event driven) behavior, that we introduced in section 2.2.

There is clearly a tradeoff to be struck between the degree of commitment and reconsideration at work here. The nature of this tradeoff was examined by David Kinny and Michael Georgeff, in a number of experiments carried out with a BDI agent framework called dMARS [29]. They investigate how *bold* agents (those that never stop to reconsider) and *cautious* agents (those that are constantly stopping to reconsider) perform in a variety of different environments. The most important parameter in these experiments was the *rate of world change*, γ . The key results of Kinny and Georgeff were as follows.

- If γ is low, (i.e., the environment does not change quickly), then bold agents do well compared to cautious ones, because cautious ones waste time reconsidering their commitments while bold agents are busy working towards – and achieving – their goals.
- If γ is high, (i.e., the environment changes frequently), then cautious agents tend to outperform bold agents, because they are able to recognize when intentions are doomed, and also to take advantage of serendipitous situations and new opportunities.

The lesson is that different types of environment require different types of decision strategies. In static, unchanging environment, purely pro-active, goal directed behavior is adequate. But in more dynamic environments, the ability to react to changes by modifying intentions becomes more important.

The process of practical reasoning in a BDI agent is summarized in Fig. 5. As this figure illustrates, there are seven main components to a BDI agent:

- a set of current *beliefs*, representing information the agent has about its current environment;

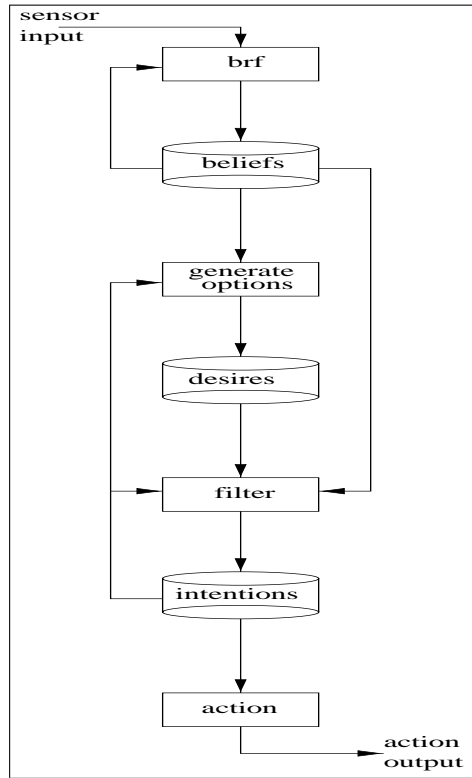


Fig. 5. Schematic diagram of a generic belief-desire-intention architecture.

- a *belief revision function*, (*brf*), which takes a perceptual input and the agent’s current beliefs, and on the basis of these, determines a new set of beliefs;
- an *option generation function*, (*options*), which determines the options available to the agent (its desires), on the basis of its current beliefs about its environment and its current *intentions*;
- a set of *current options*, representing possible courses of actions available to the agent;
- a *filter function* (*filter*), which represents the agent’s *deliberation* process, and which determines the agent’s intentions on the basis of its current beliefs, desires, and intentions;
- a set of current *intentions*, representing the agent’s current focus – those states of affairs that it has committed to trying to bring about;
- an *action selection function* (*execute*), which determines an action to perform on the basis of current intentions.

It is straightforward to formally define these components. First, let *Bel* be the set of all possible beliefs, *Des* be the set of all possible desires, and *Int* be the set of

all possible intentions. For the purposes of this chapter, the content of these sets is not important. (Often, beliefs, desires, and intentions are represented as logical formulae, perhaps of first-order logic.) Whatever the content of these sets, it is worth noting that they should have some notion of *consistency* defined upon them, so that one can answer the question of, for example, whether having an intention to achieve x is consistent with the belief that y . Representing beliefs, desires, and intentions as logical formulae permits us to recast such questions as problems of determining whether logical formulae are consistent – a well known and well-understood problem. The state of a BDI agent at any given moment is, unsurprisingly, a triple (B, D, I) , where $B \subseteq Bel$, $D \subseteq Des$, and $I \subseteq Int$.

An agent's belief revision function is a mapping

$$brf : \wp(Bel) \times P \rightarrow \wp(Bel)$$

which on the basis of the current percept and current beliefs determines a new set of beliefs. Belief revision is out of the scope of this chapter (and indeed this book), and so we shall say no more about it here.

The option generation function, *options*, maps a set of beliefs and a set of intentions to a set of desires.

$$options : \wp(Bel) \times \wp(Int) \rightarrow \wp(Des)$$

This function plays several roles. First, it must be responsible for the agent's means-ends reasoning – the process of deciding how to achieve intentions. Thus, once an agent has formed an intention to x , it must subsequently consider options to *achieve* x . These options will be more concrete – less abstract – than x . As some of these options then become intentions themselves, they will also feedback into option generation, resulting in yet more concrete options being generated. We can thus think of a BDI agent's option generation process as one of recursively elaborating a hierarchical plan structure, considering and committing to progressively more specific intentions, until finally it reaches the intentions that correspond to immediately executable actions.

While the main purpose of the *options* function is thus means-ends reasoning, it must in addition satisfy several other constraints. First, it must be *consistent*: any options generated must be consistent with both the agent's current beliefs and current intentions. Secondly, it must be *opportunistic*, in that it should recognize when environmental circumstances change advantageously, to offer the agent new ways of achieving intentions, or the possibility of achieving intentions that were otherwise unachievable.

A BDI agent's deliberation process (deciding *what* to do) is represented in the *filter* function,

$$filter : \wp(Bel) \times \wp(Des) \times \wp(Int) \rightarrow \wp(Int)$$

which updates the agent's intentions on the basis of its previously-held intentions and current beliefs and desires. This function must fulfill two roles. First, it must *drop* any intentions that are no longer achievable, or for which the expected cost

of achieving them exceeds the expected gain associated with successfully achieving them. Second, it should *retain* intentions that are not achieved, and that are still expected to have a positive overall benefit. Finally, it should *adopt* new intentions, either to achieve existing intentions, or to exploit new opportunities.

Notice that we do not expect this function to introduce intentions from nowhere. Thus *filter* should satisfy the following constraint:

$$\forall B \in \wp(\text{Bel}), \forall D \in \wp(\text{Des}), \forall I \in \wp(\text{Int}), \text{filter}(B, D, I) \subseteq I \cup D.$$

In other words, current intentions are either previously held intentions or newly adopted options.

The *execute* function is assumed to simply return any executable intentions – one that corresponds to a directly executable action:

$$\text{execute} : \wp(\text{Int}) \rightarrow A$$

The agent decision function, *action* of a BDI agent is then a function

$$\text{action} : P \rightarrow A$$

and is defined by the following pseudo-code.

1. function *action*($p : P$) : A
2. begin
3. $B := \text{brf}(B, p)$
4. $D := \text{options}(D, I)$
5. $I := \text{filter}(B, D, I)$
6. return *execute*(I)
7. end function *action*

Note that representing an agent's intentions as a *set* (i.e., as an unstructured collection) is generally too simplistic in practice. A simple alternative is to associate a *priority* with each intention, indicating its relative importance. Another natural idea is to represent intentions as a *stack*. An intention is pushed on to the stack when it is adopted, and popped when it is either achieved or else not achievable. More abstract intentions will tend to be at the bottom of the stack, with more concrete intentions towards the top.

To summarize, BDI architectures are practical reasoning architectures, in which the process of deciding what to do resembles the kind of practical reasoning that we appear to use in our everyday lives. The basic components of a BDI architecture are data structures representing the beliefs, desires, and intentions of the agent, and functions that represent its deliberation (deciding *what* intentions to have – i.e., deciding what to do) and means-ends reasoning (deciding how to do it). Intentions play a central role in the BDI model: they provide stability for decision making, and act to focus the agent's practical reasoning. A major issue in BDI architectures is the problem of striking a *balance* between

being committed to and overcommitted to one's intentions: the deliberation process must be finely tuned to its environment, ensuring that in more dynamic, highly unpredictable domains, it reconsiders its intentions relatively frequently – in more static environments, less frequent reconsideration is necessary.

The BDI model is attractive for several reasons. First, it is intuitive – we all recognize the processes of deciding what to do and then how to do it, and we all have an informal understanding of the notions of belief, desire, and intention. Second, it gives us a clear functional decomposition, which indicates what sorts of subsystems might be required to build an agent. But the main difficulty, as ever, is knowing how to efficiently implement these functions.

Sources and Further Reading. Belief-desire-intention architectures originated in the work of the Rational Agency project at Stanford Research Institute in the mid 1980s. The origins of the model lie in the theory of human practical reasoning developed by the philosopher Michael Bratman [5], which focusses particularly on the role of intentions in practical reasoning. The conceptual framework of the BDI model is described in [6], which also describes a specific BDI agent architecture called IRMA. The description of the BDI model given here (and in particular Fig. 5) is adapted from [6]. One of the interesting aspects of the BDI model is that it has been used in one of the most successful agent architectures to date. The Procedural Reasoning System (PRS), originally developed by Michael Georgeff and Amy Lansky [22], has been used to build some of the most exacting agent applications to date, including fault diagnosis for the reaction control system of the space shuttle, and an air traffic management system at Sydney airport in Australia – overviews of these systems are described in [23]. In the PRS, an agent is equipped with a library of *plans* which are used to perform means-ends reasoning. Deliberation is achieved by the use of *meta-level plans*, which are able to modify an agent's intention structure at run-time, in order to change the focus of the agent's practical reasoning. Beliefs in the PRS are represented as PROLOG-like facts – essentially, as atoms of first-order logic.

The BDI model is also interesting because a great deal of effort has been devoted to formalising it. In particular, Anand Rao and Michael Georgeff have developed a range of BDI logics, which they use to axiomatize properties of BDI-based practical reasoning agents [42,46,43,44,45,41]. These models have been extended by others to deal with, for example, communication between agents [24].

4.4 Layered Architectures

Given the requirement that an agent be capable of reactive and pro-active behavior, an obvious decomposition involves creating separate subsystems to deal with these different types of behaviors. This idea leads naturally to a class of architectures in which the various subsystems are arranged into a hierarchy of interacting *layers*. In this section, we will consider some general aspects of layered architectures, and then go on to consider two examples of such architectures: INTERRAP and TOURINGMACHINES.

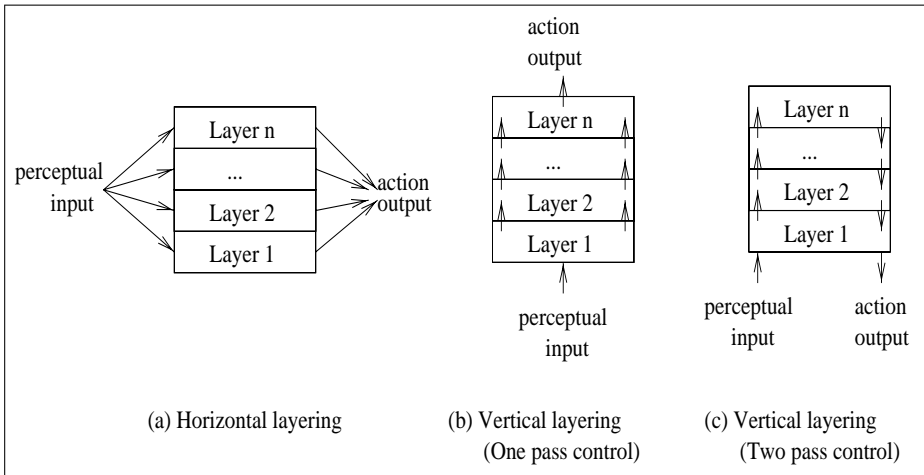


Fig. 6. Information and control flows in three types of layered agent architecture (Source: [38, p263]).

Typically, there will be at least two layers, to deal with reactive and proactive behaviors respectively. In principle, there is no reason why there should not be many more layers. However many layers there are, a useful typology for such architectures is by the information and control flows within them. Broadly speaking, we can identify two types of control flow within layered architectures (see Fig. 6):

– *Horizontal layering.*

In horizontally layered architectures (Fig. 6a), the software layers are each directly connected to the sensory input and action output. In effect, each layer itself acts like an agent, producing suggestions as to what action to perform.

– *Vertical layering.*

In vertically layered architectures (Fig. 6b and 6c), sensory input and action output are each dealt with by at most one layer each.

The great advantage of horizontally layered architectures is their conceptual simplicity: if we need an agent to exhibit n different types of behavior, then we implement n different layers. However, because the layers are each in effect competing with one-another to generate action suggestions, there is a danger that the *overall* behavior of the agent will not be coherent. In order to ensure that horizontally layered architectures *are* consistent, they generally include a *mediator* function, which makes decisions about which layer has “control” of the agent at any given time. The need for such central control is problematic: it means that the designer must potentially consider all possible interactions between layers. If there are n layers in the architecture, and each layer is capable

of suggesting m possible actions, then this means there are m^n such interactions to be considered. This is clearly difficult from a design point of view in any but the most simple system. The introduction of a central control system also introduces a *bottleneck* into the agent's decision making.

These problems are partly alleviated in a vertically layered architecture. We can subdivide vertically layered architectures into *one pass* architectures (Fig. 6b) and *two pass* architectures (Fig. 6c). In one-pass architectures, control flows sequentially through each layer, until the final layer generates action output. In two-pass architectures, information flows up the architecture (the first pass) and control then flows back down. There are some interesting similarities between the idea of two-pass vertically layered architectures and the way that organisations work, with information flowing up to the highest levels of the organisation, and commands then flowing down. In both one pass and two pass vertically layered architectures, the complexity of interactions between layers is reduced: since there are $n - 1$ interfaces between n layers, then if each layer is capable of suggesting m actions, there are at most $m^2(n - 1)$ interactions to be considered between layers. This is clearly much simpler than the horizontally layered case. However, this simplicity comes at the cost of some flexibility: in order for a vertically layered architecture to make a decision, control must pass between *each* different layer. This is not fault tolerant: failures in any one layer are likely to have serious consequences for agent performance.

In the remainder of this section, we will consider two examples of layered architectures: Innes Ferguson's TOURINGMACHINES, and Jörg Müller's INTERRAP. The former is an example of a horizontally layered architecture; the latter is a (two pass) vertically layered architecture.

TouringMachines. The TOURINGMACHINES architecture is illustrated in Fig. 7. As this figure shows, TOURINGMACHINES consists of three *activity producing layers*. That is, each layer continually produces "suggestions" for what actions the agent should perform. The *reactive layer* provides a more-or-less immediate response to changes that occur in the environment. It is implemented as a set of situation-action rules, like the behaviors in Brooks' subsumption architecture (Sect. 4.2). These rules map sensor input directly to effector output. The original demonstration scenario for TOURINGMACHINES was that of autonomous vehicles driving between locations through streets populated by other similar agents. In this scenario, reactive rules typically deal with functions like obstacle avoidance. For example, here is an example of a reactive rule for avoiding the kerb (from [14, p59]):

```
rule-1: kerb-avoidance
  if
    is-in-front(Kerb, Observer) and
    speed(Observer) > 0 and
    separation(Kerb, Observer) < KerbThreshold
  then
    change-orientation(KerbAvoidanceAngle)
```

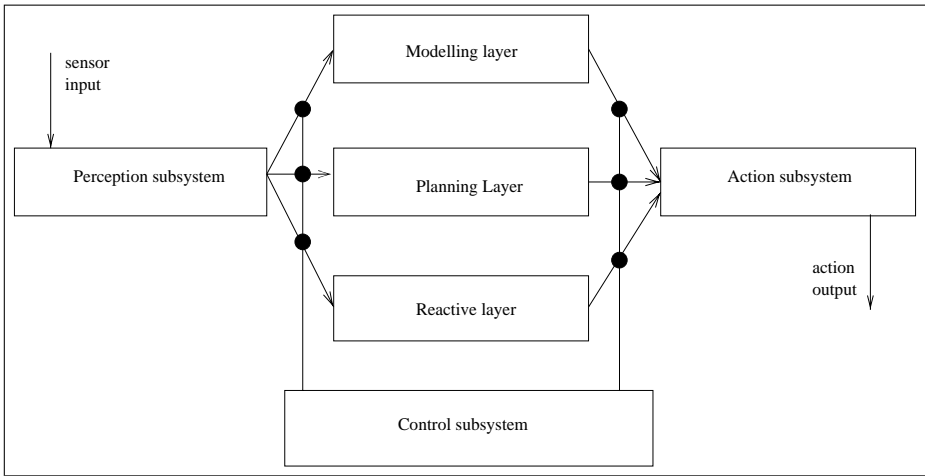


Fig. 7. TOURINGMACHINES: a horizontally layered agent architecture

Here `change-orientation(...)` is the action suggested if the rule fires. The rules can only make references to the agent’s current state – they cannot do any explicit reasoning about the world, and on the right hand side of rules are *actions*, not predicates. Thus if this rule fired, it would not result in any central environment model being updated, but would just result in an action being suggested by the reactive layer.

The TOURINGMACHINES *planning layer* achieves the agent’s pro-active behavior. Specifically, the planning layer is responsible for the “day-to-day” running of the agent – under normal circumstances, the planning layer will be responsible for deciding what the agent does. However, the planning layer does not do “first-principles” planning. That is, it does not attempt to generate plans from scratch. Rather, the planning layer employs a *library* of plan “skeletons” called *schemas*. These skeletons are in essence hierarchically structured plans, which the TOURINGMACHINES planning layer elaborates at run time in order to decide what to do. So, in order to achieve a goal, the planning layer attempts to find a schema in its library which matches that goal. This schema will contain sub-goals, which the planning layer elaborates by attempting to find other schemas in its plan library that match these sub-goals.

The *modeling* layer represents the various entities in the world (including the agent itself, as well as other agents). The modeling layer thus predicts conflicts between agents, and generates new goals to be achieved in order to resolve these conflicts. These new goals are then posted down to the planning layer, which makes use of its plan library in order to determine how to satisfy them.

The three control layers are embedded within a *control subsystem*, which is effectively responsible for deciding which of the layers should have control over the agent. This control subsystem is implemented as a set of *control rules*.

Control rules can either *suppress* sensor information between the control rules and the control layers, or else *ensor* action outputs from the control layers. Here is an example sensor rule [16, p207]:

```

sensor-rule-1:
  if
    entity(obstacle-6) in perception-buffer
  then
    remove-sensory-record(layer-R, entity(obstacle-6))

```

This rule prevents the reactive layer from ever knowing about whether `obstacle-6` has been perceived. The intuition is that although the reactive layer will in general be the most appropriate layer for dealing with obstacle avoidance, there are certain obstacles for which other layers are more appropriate. This rule ensures that the reactive layer never comes to know about these obstacles.

InteRRaP. INTERRAP is an example of a vertically layered two-pass agent architecture – see Fig. 8.

As Fig. 8 shows, INTERRAP contains three control layers, as in TOURINGMACHINES. Moreover, the purpose of each INTERRAP layer appears to be rather similar to the purpose of each corresponding TOURINGMACHINES layer. Thus the lowest (*behavior based*) layer deals with reactive behavior; the middle (*local planning*) layer deals with everyday planning to achieve the agent’s goals, and the uppermost (*cooperative planning*) layer deals with social interactions. Each layer has associated with it a *knowledge base*, i.e., a representation of the world appropriate for that layer. These different knowledge bases represent the agent and its environment at different levels of abstraction. Thus the highest

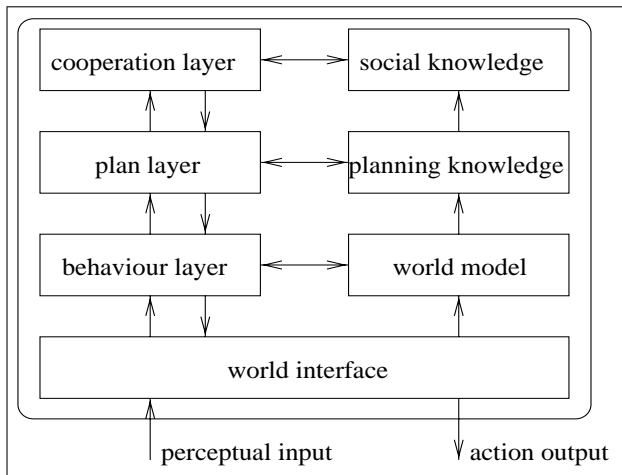


Fig. 8. INTERRAP – a vertically layered two-pass agent architecture.

level knowledge base represents the plans and actions of other agents in the environment; the middle-level knowledge base represents the plans and actions of the agent itself; and the lowest level knowledge base represents “raw” information about the environment. The explicit introduction of these knowledge bases distinguishes TOURINGMACHINES from INTERRAP.

The way the different layers in INTERRAP conspire to produce behavior is also quite different from TOURINGMACHINES. The main difference is in the way the layers interact with the environment. In TOURINGMACHINES, each layer was directly coupled to perceptual input and action output. This necessitated the introduction of a supervisory control framework, to deal with conflicts or problems between layers. In INTERRAP, layers interact with *each other* to achieve the same end. The two main types of interaction between layers are *bottom-up activation* and *top-down execution*. Bottom-up activation occurs when a lower layer passes control to a higher layer because it is not *competent* to deal with the current situation. Top-down execution occurs when a higher layer makes use of the facilities provided by a lower layer to achieve one of its goals. The basic flow of control in INTERRAP begins when perceptual input arrives at the lowest layer in the architecture. If the reactive layer can deal with this input, then it will do so; otherwise, bottom-up activation will occur, and control will be passed to the local planning layer. If the local planning layer can handle the situation, then it will do so, typically by making use of top-down execution. Otherwise, it will use bottom-up activation to pass control to the highest layer. In this way, control in INTERRAP will flow from the lowest layer to higher layers of the architecture, and then back down again.

The internals of each layer are not important for the purposes of this chapter. However, it is worth noting that each layer implements two general functions. The first of these is a *situation recognition and goal activation* function. This function acts rather like the *options* function in a BDI architecture (see Sect. 4.3). It maps a knowledge base (one of the three layers) and current goals to a new set of goals. The second function is responsible for *planning and scheduling* – it is responsible for selecting which plans to execute, based on the current plans, goals, and knowledge base of that layer.

Layered architectures are currently the most popular general class of agent architecture available. Layering represents a natural decomposition of functionality: it is easy to see how reactive, pro-active, social behavior can be generated by the reactive, pro-active, and social layers in an architecture. The main problem with layered architectures is that while they are arguably a *pragmatic* solution, they lack the conceptual and semantic clarity of unlayered approaches. In particular, while logic-based approaches have a clear logical semantics, it is difficult to see how such a semantics could be devised for a layered architecture. Another issue is that of interactions between layers. If each layer is an independent activity producing process (as in TOURINGMACHINES), then it is necessary to consider all possible ways that the layers can interact with one another. This problem is partly alleviated in two-pass vertically layered architecture such as INTERRAP.

Sources and Further Reading. The introductory discussion of layered architectures given here draws heavily upon [38, pp262–264]. The best reference to TOURINGMACHINES is [14]; more accessible references include [15,16]. The definitive reference to INTERRAP is [37], although [18] is also a useful reference. Other examples of layered architectures include the subsumption architecture [7] (see also Sect. 4.2), and the 3T architecture [4].

5 Conclusions

I hope that after reading this chapter, you understand what agents are and why they are considered to be an important area of research and development. The requirement for systems that can operate autonomously is very common. The requirement for systems capable of *flexible* autonomous action, in the sense that I have described in this chapter, is similarly common. This leads me to conclude that intelligent agents have the potential to play a significant role in the future of software engineering. Intelligent agent research is about the theory, design, construction, and application of such systems. This chapter has focussed on the design of intelligent agents. It has presented a high-level, abstract view of intelligent agents, and described the sort of properties that one would expect such an agent to enjoy. It went on to show how this view of an agent could be refined into various different types of agent architecture – purely logical agents, purely reactive/behavioral agents, BDI agents, and layered agent architectures.

References

1. P. Agre and D. Chapman. PENGI: An implementation of a theory of activity. In *Proceedings of the Sixth National Conference on Artificial Intelligence (AAAI-87)*, pages 268–272, Seattle, WA, 1987.
2. P. E. Agre and S. J. Rosenschein, editors. *Computational Theories of Interaction and Agency*. The MIT Press: Cambridge, MA, 1996.
3. H. Barringer, M. Fisher, D. Gabbay, G. Gough, and R. Owens. METATEM: A framework for programming in temporal logic. In *REX Workshop on Stepwise Refinement of Distributed Systems: Models, Formalisms, Correctness (LNCS Volume 430)*, pages 94–129. Springer-Verlag: Berlin, Germany, June 1989.
4. R. P. Bonasso, D. Kortenkamp, D. P. Miller, and M. Slack. Experiences with an architecture for intelligent, reactive agents. In M. Wooldridge, J. P. Müller, and M. Tambe, editors, *Intelligent Agents II (LNAI Volume 1037)*, pages 187–202. Springer-Verlag: Berlin, Germany, 1996.
5. M. E. Bratman. *Intentions, Plans, and Practical Reason*. Harvard University Press: Cambridge, MA, 1987.
6. M. E. Bratman, D. J. Israel, and M. E. Pollack. Plans and resource-bounded practical reasoning. *Computational Intelligence*, 4:349–355, 1988.
7. R. A. Brooks. A robust layered control system for a mobile robot. *IEEE Journal of Robotics and Automation*, 2(1):14–23, 1986.
8. R. A. Brooks. Elephants don't play chess. In P. Maes, editor, *Designing Autonomous Agents*, pages 3–15. The MIT Press: Cambridge, MA, 1990.

9. R. A. Brooks. Intelligence without reason. In *Proceedings of the Twelfth International Joint Conference on Artificial Intelligence (IJCAI-91)*, pages 569–595, Sydney, Australia, 1991.
10. R. A. Brooks. Intelligence without representation. *Artificial Intelligence*, 47:139–159, 1991.
11. Oren Etzioni. Intelligence without robots. *AI Magazine*, 14(4), December 1993.
12. R. Fagin, J. Y. Halpern, Y. Moses, and M. Y. Vardi. *Reasoning About Knowledge*. The MIT Press: Cambridge, MA, 1995.
13. J. Ferber. Reactive distributed artificial intelligence. In G. M. P. O'Hare and N. R. Jennings, editors, *Foundations of Distributed Artificial Intelligence*, pages 287–317. John Wiley, 1996.
14. I. A. Ferguson. *Touring Machines: An Architecture for Dynamic, Rational, Mobile Agents*. PhD thesis, Clare Hall, University of Cambridge, UK, November 1992. (Also available as Technical Report No. 273, University of Cambridge Computer Laboratory).
15. I. A. Ferguson. Towards an architecture for adaptive, rational, mobile agents. In E. Werner and Y. Demazeau, editors, *Decentralized AI 3 – Proceedings of the Third European Workshop on Modelling Autonomous Agents in a Multi-Agent World (MAAMAW-91)*, pages 249–262. Elsevier Science Publishers B.V.: Amsterdam, The Netherlands, 1992.
16. I. A. Ferguson. Integrated control and coordinated behaviour: A case for agent models. In M. Wooldridge and N. R. Jennings, editors, *Intelligent Agents: Theories, Architectures, and Languages (LNAI Volume 890)*, pages 203–218. Springer-Verlag: Berlin, Germany, January 1995.
17. J. A. Firby. An investigation into reactive planning in complex domains. In *Proceedings of the Tenth International Joint Conference on Artificial Intelligence (IJCAI-87)*, pages 202–206, Milan, Italy, 1987.
18. K. Fischer, J. P. Müller, and M. Pischel. A pragmatic BDI architecture. In M. Wooldridge, J. P. Müller, and M. Tambe, editors, *Intelligent Agents II (LNAI Volume 1037)*, pages 203–218. Springer-Verlag: Berlin, Germany, 1996.
19. M. Fisher. A survey of Concurrent METATEM – the language and its applications. In D. M. Gabbay and H. J. Ohlbach, editors, *Temporal Logic – Proceedings of the First International Conference (LNAI Volume 827)*, pages 480–505. Springer-Verlag: Berlin, Germany, July 1994.
20. L. Gasser and J. P. Briot. Object-based concurrent programming and DAI. In *Distributed Artificial Intelligence: Theory and Praxis*, pages 81–108. Kluwer Academic Publishers: Boston, MA, 1992.
21. M. R. Genesereth and N. Nilsson. *Logical Foundations of Artificial Intelligence*. Morgan Kaufmann Publishers: San Mateo, CA, 1987.
22. M. P. Georgeff and A. L. Lansky. Reactive reasoning and planning. In *Proceedings of the Sixth National Conference on Artificial Intelligence (AAAI-87)*, pages 677–682, Seattle, WA, 1987.
23. M. P. Georgeff and A. S. Rao. A profile of the Australian AI Institute. *IEEE Expert*, 11(6):89–92, December 1996.
24. A. Haddadi. *Communication and Cooperation in Agent Systems (LNAI Volume 1056)*. Springer-Verlag: Berlin, Germany, 1996.
25. J. Y. Halpern. Using reasoning about knowledge to analyze distributed systems. *Annual Review of Computer Science*, 2:37–68, 1987.

26. L. P. Kaelbling. An architecture for intelligent reactive systems. In M. P. Georgeff and A. L. Lansky, editors, *Reasoning About Actions & Plans – Proceedings of the 1986 Workshop*, pages 395–410. Morgan Kaufmann Publishers: San Mateo, CA, 1986.
27. L. P. Kaelbling. A situated automata approach to the design of embedded agents. *SIGART Bulletin*, 2(4):85–88, 1991.
28. L. P. Kaelbling and S. J. Rosenschein. Action and planning in embedded agents. In P. Maes, editor, *Designing Autonomous Agents*, pages 35–48. The MIT Press: Cambridge, MA, 1990.
29. D. Kinny and M. Georgeff. Commitment and effectiveness of situated agents. In *Proceedings of the Twelfth International Joint Conference on Artificial Intelligence (IJCAI-91)*, pages 82–88, Sydney, Australia, 1991.
30. K. Konolige. *A Deduction Model of Belief*. Pitman Publishing: London and Morgan Kaufmann: San Mateo, CA, 1986.
31. Y. Lésperance, H. J. Levesque, F. Lin, D. Marcu, R. Reiter, and R. B. Scherl. Foundations of a logical approach to agent programming. In M. Wooldridge, J. P. Müller, and M. Tambe, editors, *Intelligent Agents II (LNAI Volume 1037)*, pages 331–346. Springer-Verlag: Berlin, Germany, 1996.
32. P. Maes. The dynamics of action selection. In *Proceedings of the Eleventh International Joint Conference on Artificial Intelligence (IJCAI-89)*, pages 991–997, Detroit, MI, 1989.
33. P. Maes, editor. *Designing Autonomous Agents*. The MIT Press: Cambridge, MA, 1990.
34. P. Maes. Situated agents can have goals. In P. Maes, editor, *Designing Autonomous Agents*, pages 49–70. The MIT Press: Cambridge, MA, 1990.
35. P. Maes. The agent network architecture (ANA). *SIGART Bulletin*, 2(4):115–120, 1991.
36. J. McCarthy and P. J. Hayes. Some philosophical problems from the standpoint of artificial intelligence. In B. Meltzer and D. Michie, editors, *Machine Intelligence 4*. Edinburgh University Press, 1969.
37. J. P. Müller. *The Design of Intelligent Agents*. Springer-Verlag: Berlin, Germany, 1997.
38. J. P. Müller, M. Pischel, and M. Thiel. Modelling reactive behaviour in vertically layered agent architectures. In M. Wooldridge and N. R. Jennings, editors, *Intelligent Agents: Theories, Architectures, and Languages (LNAI Volume 890)*, pages 261–276. Springer-Verlag: Berlin, Germany, January 1995.
39. J. P. Müller, M. Wooldridge, and N. R. Jennings, editors. *Intelligent Agents III (LNAI Volume 1193)*. Springer-Verlag: Berlin, Germany, 1995.
40. N. J. Nilsson. Towards agent programs with circuit semantics. Technical Report STAN-CS-92-1412, Computer Science Department, Stanford University, Stanford, CA 94305, January 1992.
41. A. S. Rao. Decision procedures for propositional linear-time Belief-Desire-Intention logics. In M. Wooldridge, J. P. Müller, and M. Tambe, editors, *Intelligent Agents II (LNAI Volume 1037)*, pages 33–48. Springer-Verlag: Berlin, Germany, 1996.
42. A. S. Rao and M. P. Georgeff. Asymmetry thesis and side-effect problems in linear time and branching time intention logics. In *Proceedings of the Twelfth International Joint Conference on Artificial Intelligence (IJCAI-91)*, pages 498–504, Sydney, Australia, 1991.

43. A. S. Rao and M. P. Georgeff. Modeling rational agents within a BDI-architecture. In R. Fikes and E. Sandewall, editors, *Proceedings of Knowledge Representation and Reasoning (KR&R-91)*, pages 473–484. Morgan Kaufmann Publishers: San Mateo, CA, April 1991.
44. A. S. Rao and M. P. Georgeff. An abstract architecture for rational agents. In C. Rich, W. Swartout, and B. Nebel, editors, *Proceedings of Knowledge Representation and Reasoning (KR&R-92)*, pages 439–449, 1992.
45. A. S. Rao and M. P. Georgeff. A model-theoretic approach to the verification of situated reasoning systems. In *Proceedings of the Thirteenth International Joint Conference on Artificial Intelligence (IJCAI-93)*, pages 318–324, Chambéry, France, 1993.
46. A. S. Rao, M. P. Georgeff, and E. A. Sonenberg. Social plans: A preliminary report. In E. Werner and Y. Demazeau, editors, *Decentralized AI 3 – Proceedings of the Third European Workshop on Modelling Autonomous Agents in a Multi-Agent World (MAAMAW-91)*, pages 57–76. Elsevier Science Publishers B.V.: Amsterdam, The Netherlands, 1992.
47. S. Rosenschein and L. P. Kaelbling. The synthesis of digital machines with provable epistemic properties. In J. Y. Halpern, editor, *Proceedings of the 1986 Conference on Theoretical Aspects of Reasoning About Knowledge*, pages 83–98. Morgan Kaufmann Publishers: San Mateo, CA, 1986.
48. S. J. Rosenschein and L. P. Kaelbling. A situated view of representation and control. In P. E. Agre and S. J. Rosenschein, editors, *Computational Theories of Interaction and Agency*, pages 515–540. The MIT Press: Cambridge, MA, 1996.
49. S. Russell and P. Norvig. *Artificial Intelligence: A Modern Approach*. Prentice-Hall, 1995.
50. S. Russell and D. Subramanian. Provably bounded-optimal agents. *Journal of AI Research*, 2:575–609, 1995.
51. S. J. Russell and E. Wefald. *Do the Right Thing – Studies in Limited Rationality*. The MIT Press: Cambridge, MA, 1991.
52. M. J. Schoppers. Universal plans for reactive robots in unpredictable environments. In *Proceedings of the Tenth International Joint Conference on Artificial Intelligence (IJCAI-87)*, pages 1039–1046, Milan, Italy, 1987.
53. L. Steels. Cooperation between distributed agents through self organization. In Y. Demazeau and J.-P. Müller, editors, *Decentralized AI – Proceedings of the First European Workshop on Modelling Autonomous Agents in a Multi-Agent World (MAAMAW-89)*, pages 175–196. Elsevier Science Publishers B.V.: Amsterdam, The Netherlands, 1990.
54. M. Wooldridge. Agent-based software engineering. *IEE Transactions on Software Engineering*, 144(1):26–37, February 1997.
55. M. Wooldridge and N. R. Jennings. Intelligent agents: Theory and practice. *The Knowledge Engineering Review*, 10(2):115–152, 1995.