# CHAPTER 2: INTELLIGENT AGENTS

An Introduction to Multiagent Systems

`http://www.csc.liv.ac.uk/~mjw/pubs/imas/`

| What is an Agent? |
|---|

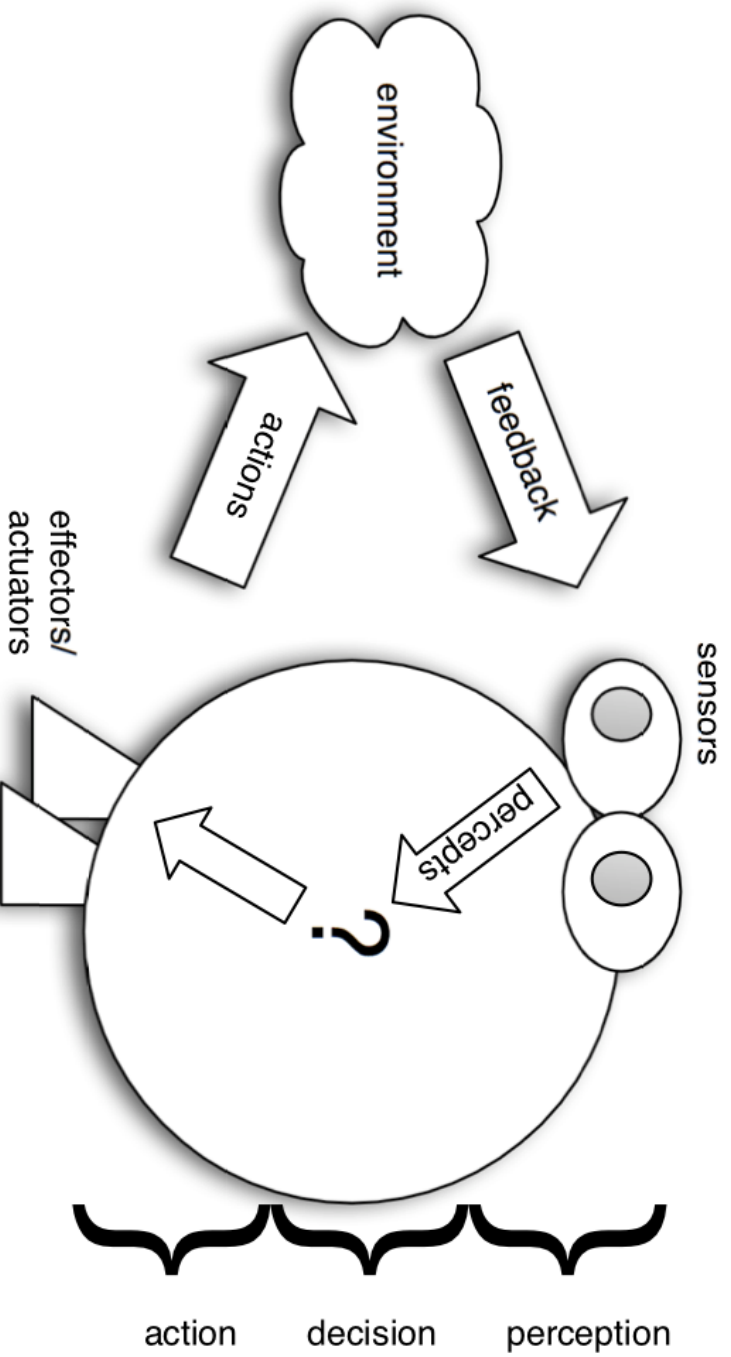- The main point about agents is they are *autonomous*: capable independent action.

- Thus:

  *an agent is a computer system capable of autonomous action in some environment, in order to achieve its delegated goals.*

- We think of an agent as being in a close-coupled, continual interaction with its environment:

  sense – decide – act – sense – decide . . .

http://www.csc.liv.ac.uk/~mjw/pubs/imas/

# Agent and Environment

## Simple (Uninteresting) Agents

- Thermostat
  - delegated goal is maintain room temperature
  - actions are heat on/off

- UNIX biff program
  - delegated goal is monitor for incoming email and flag it
  - actions are GUI actions.

They are trivial because the *decision making* they do is trivial.

# Intelligent Agents

We typically think of as intelligent agent as exhibiting 3 types of behaviour:

- *reactive*;
- *pro-active*;
- *social*.

## Reactivity

- If a program's environment is guaranteed to be fixed, a program can just execute blindly.

- The real world is not like that: most environments are *dynamic*.

- Software is hard to build for dynamic domains: program must take into account possibility of failure — ask itself whether it is worth executing!

- A *reactive* system is one that maintains an ongoing interaction with its environment, and responds to changes that occur in it (in time for the response to be useful).

# Proactiveness

- Reacting to an environment is easy (e.g., stimulus → response rules).

- But we generally want agents to *do things for us*.

- Hence *goal directed behaviour*.

- Pro-activeness = generating and attempting to achieve goals; not driven solely by events; taking the initiative.

- Recognising opportunities.

# Social Ability

- The real world is a *multi*-agent environment: we cannot go around attempting to achieve goals without taking others into account.

- Some goals can only be achieved by interacting with others.

- Similarly for many computer environments: witness the INTERNET.

- *Social ability* in agents is the ability to interact with other agents (and possibly humans) via *cooperation*, *coordination*, and *negotiation*.

At the very least, it means the ability to communicate...

## Social Ability: Cooperation

- Cooperation is *working together as a team to achieve a shared goal.*

- Often prompted either by the fact that no one agent can achieve the goal alone, or that cooperation will obtain a better result (e.g., get result faster).

# Social Ability: Coordination

- Coordination is *managing the interdependencies between activities*.

- For example, if there is a non-sharable resource that you want to use and I want to use, then we need to coordinate.

# Social Ability: Negotiation

- Negotiation is *the ability to reach agreements on matters of common interest.*

- For example: You have one TV in your house; you want to watch a movie, your housemate wants to watch football.

  A possible deal: watch football tonight, and a movie tomorrow.

- Typically involves offer and counter-offer, with compromises made by participants.

# Some Other Properties....

- *Mobility*
- *Veracity*
- *Benevolence*
- *Rationality*
- *Learning/adaption:*

# Agents and Objects

- Are agents just objects by another name?
- Object:
  - encapsulates some state;
  - communicates via message passing;
  - has methods, corresponding to operations that may be performed on this state.

# Differences between Agents & Objects

● *Agents are autonomous:*

agents embody stronger notion of autonomy than objects, and in particular, they decide for themselves whether or not to perform an action on request from another agent;

● *Agents are smart:*

capable of flexible (reactive, pro-active, social) behavior – the OO model has nothing to say about such types of behavior;

● *Agents are active:*

not passive service providers.

Objects do it for free....

- *agents do it because they want to;*

- *agents do it for money.*

## Agents and Expert Systems

- Aren't agents just expert systems by another name?

- Expert systems typically disembodied 'expertise' about some (abstract) domain of discourse.

- Example: MYCIN knows about blood diseases in humans.

  It has a wealth of knowledge about blood diseases, in the form of rules.

  A doctor can obtain expert advice about blood diseases by giving MYCIN facts, answering questions, and posing queries.

# Differences between Agents & Expert Systems

- agents are *situated in an environment*:
  MYCIN is not aware of the world — only information obtained is by asking the user questions.

- agents *act*:
  MYCIN does not operate on patients.

Some *real-time* (typically process control) expert systems *are* agents.

## Intelligent Agents and AI

- Aren't agents just the AI project?
  Isn't building an agent what AI is all about?

- AI aims to build systems that can (ultimately)
  understand natural language, recognise and
  understand scenes, use common sense, think
  creatively, etc — all of which are very hard.

- So, don't we need to solve all of AI to build an
  agent. . . ?

- When building an agent, we simply want a system that can choose the right action to perform, typically in a limited domain.

- We *do not* have to solve *all* the problems of AI to build a useful agent:

  *a little intelligence goes a long way!*

- Oren Etzioni, speaking about the commercial experience of NETBOT, Inc:

  We made our agents dumber and dumber and dumber . . . until finally they made money.

# Properties of Environments

- *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.

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.

# Agents as Intentional Systems

- When explaining human activity, we use statements like the following:

  Janine took her umbrella because she *believed* it was raining and she *wanted* to stay dry.

- These statements make use of a *folk psychology*, by which human behaviour is predicted and explained by attributing *attitudes* such as believing, wanting, hoping, fearing, ... ...

# Dennett on Intentional Systems

Daniel Dennett coined the term *intentional system* to describe entities 'whose behaviour can be predicted by the method of attributing belief, desires and rational acumen'.

'A *first-order* intentional system has beliefs and desires (etc.) but no beliefs and desires *about* beliefs and desires.

... A *second-order* intentional system is more sophisticated; it has beliefs and desires (and no doubt other intentional states) about beliefs and desires (and other intentional states) — both those of others and its own'.

# Can We Apply the Intentional Stance to Machines?

'To ascribe *beliefs*, *free will*, *intentions*, *consciousness*, *abilities*, or *wants* to a machine is *legitimate* when such an ascription expresses the same information about the machine that it expresses about a person. It is *useful* when the ascription helps us understand the structure of the machine, its past or future behaviour, or how to repair or improve it. It is perhaps never *logically required* even for humans, but expressing reasonably briefly what is actually known about the state of the machine in a particular situation may require mental qualities or qualities isomorphic to them. Theories of belief, knowledge and wanting can be constructed for machines in a simpler setting than for humans, and later applied to humans. Ascription of mental qualities is *most straightforward* for machines of known structure such as thermostats and computer operating systems, but is *most useful* when applied to entities whose structure is incompletely known'. (John McCarthy)

## What can be described with the intentional stance?

Consider a light switch:

'It is perfectly coherent to treat a light switch as a
(very cooperative) agent with the capability of
transmitting current at will, who invariably
transmits current when it believes that we want it
transmitted and not otherwise; flicking the switch
is simply our way of communicating our desires'.
(Yoav Shoham)

- Most adults would find such a description absurd!

- While the intentional stance description is consistent,

  ... it does not *buy us anything*, since we essentially understand the mechanism sufficiently to have a simpler, mechanistic description of its behaviour. (Yoav Shoham)

- The more we know about a system, the less we need to rely on animistic, intentional explanations of its behaviour.

- But with very complex systems, a mechanistic, explanation of its behaviour may not be practicable.

- *As computer systems become ever more complex, we need more powerful abstractions and metaphors to explain their operation — low level explanations become impractical.*

  *The intentional stance is such an abstraction.*

- The intentional notions are thus *abstraction tools*, which provide us with a convenient and familiar way of describing, explaining, and predicting the behaviour of complex systems.

- Remember: most important developments in computing are based on new *abstractions*:

  – procedural abstraction;

  – abstract data types;

  – objects.

  Agents, and agents as intentional systems, represent a further, and increasingly powerful abstraction.

- Points in favour of this idea:

  $\boxed{\text{Charcterising Agents}}$

  - It provides us with a familiar, non-technical way of *understanding & explaing* agents.

  $\boxed{\text{Nested Representations}}$

  - It gives us the potential to specify systems that *include representations of other systems.*

  It is widely accepted that such nested representations are essential for agents that must cooperate with other agents.

# Post-Declarative Systems

- in procedural programming, we say exactly *what* a system should do;

- in declarative programming, we state something that we want to achieve, give the system general info about the relationships between objects, and let a built-in control mechanism (e.g., goal-directed theorem proving) figure out what to do;

- with agents, we give a high-level description of the delegated goal, and let the control mechanism figure out what to do, knowing that it will act in accordance with some built-in theory of rational agency.

## An aside. . . .

- We find that researchers from a more mainstream computing discipline have adopted a similar set of ideas in *knowledge based protocols*.

- The idea: when constructing protocols, one often encounters reasoning such as the following:

  IF process $i$ knows process $j$ has received message $m_1$

  THEN process $i$ should send process $j$ the message $m_2$.

## Abstract Architectures for Agents

- Assume the environment may be in any of a finite set $E$ of discrete, instantaneous states:

$$E = \{e, e', \ldots\}.$$

- Agents are assumed to have a repertoire of possible actions available to them, which transform the state of the environment.

$$Ac = \{\alpha, \alpha', \ldots\}$$

- A *run*, $r$, of an agent in an environment is a sequence of interleaved environment states and actions:

$$r : e_0 \xrightarrow{\alpha_0} e_1 \xrightarrow{\alpha_1} e_2 \xrightarrow{\alpha_2} e_3 \xrightarrow{\alpha_3} \cdots \xrightarrow{\alpha_{u-1}} e_u$$

# Runs

## Let...

- $\mathcal{R}$ be the set of all such possible finite sequences (over $E$ and $Ac$);

- $\mathcal{R}^{Ac}$ be the subset of these that end with an action; and

- $\mathcal{R}^E$ be the subset of these that end with an environment state.

# Environments

- A *state transformer* function represents behaviour of the environment:

$$\tau : \mathcal{R}^{Ac} \to \wp(E)$$

- Note that environments are...

  – *history dependent.*

  – *non-deterministic.*

- If $\tau(r) = \emptyset$, there are no possible successor states to $r$, so we say the run has *ended*. ("Game over.")

- An environment $Env$ is then a triple $Env = \langle E, e_0, \tau \rangle$ where $E$ is set of environment states, $e_0 \in E$ is initial state; and $\tau$ is state transformer function.

# Agents

- Agent is a function which maps runs to actions:

$$Ag : \mathcal{R}^E \to Ac$$

- Thus an agent makes a decision about what action to perform based on the history of the system that it has witnessed to date.

- Let $Ag$ be the set of all agents.

# Systems

- A *system* is a pair containing an agent and an environment.

- Any system will have associated with it a set of possible runs; we denote the set of runs of agent $Ag$ in environment $Env$ by $\mathcal{R}(Ag, Env)$.

- Assume $\mathcal{R}(Ag, Env)$ contains only runs that have ended.

- **Formally, a sequence**

$$(e_0, \alpha_0, e_1, \alpha_1, e_2, \ldots)$$

**represents a run of an agent Ag in environment**
$Env = \langle E, e_0, \tau \rangle$ **if:**

1. $e_0$ is the initial state of $Env$

2. $\alpha_0 = Ag(e_0)$; and

3. for $u > 0$,

$$e_u \in \tau((e_0, \alpha_0, \ldots, \alpha_{u-1}))$$
$$\alpha_u = Ag((e_0, \alpha_0, \ldots, e_u))$$

where

## Purely Reactive Agents

- Some 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.

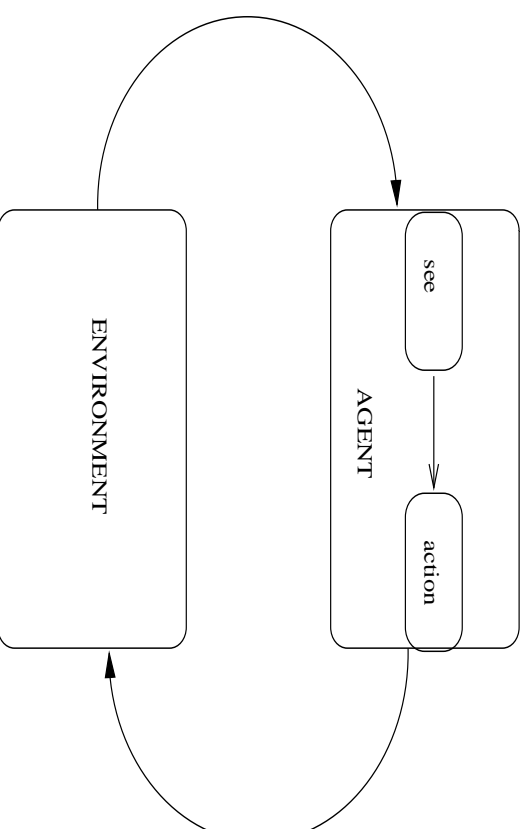- We call such agents *purely reactive*:

$$action : E \rightarrow Ac$$

- A thermostat is a purely reactive agent.

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

## Perception

- Now introduce *perception* system:

- The *see* function is the agent's ability to observe its environment, whereas the *action* function represents the agent's decision making process.

- *Output* of the *see* function is a *percept*:

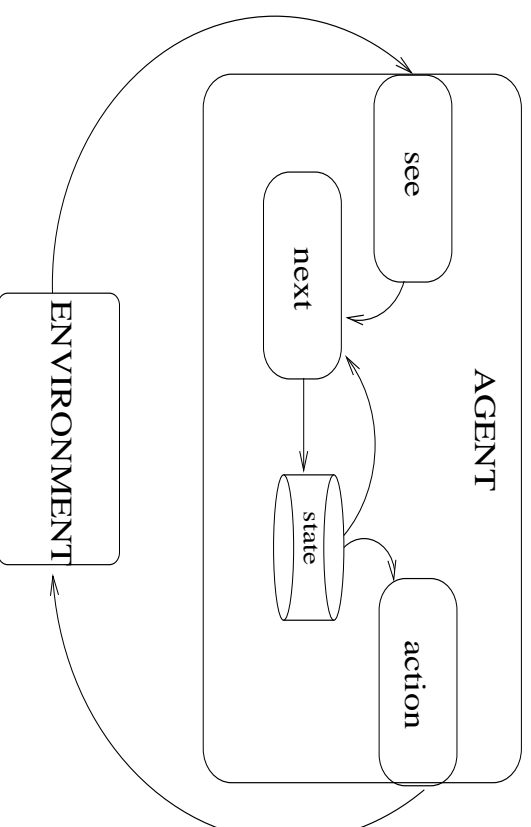$$see : E \rightarrow Per$$

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

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

which maps sequences of percepts to actions.

## Agents with State

- We now consider agents that *maintain state*:

# Perception

- 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 of all internal states of the agent.

- The perception function *see* for a state-based agent is unchanged:

$$see : E \rightarrow Per$$

# Action

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

$$action : I \rightarrow Ac$$

from internal states to actions.

## Next State Function

A function *next* is introduced, which maps an internal state and percept to an internal state:

$$next : I \times Per \to I$$

# Agent control loop

1. Agent starts in some initial internal state $i_0$.

2. repeat forever:

- Observe environment state, and generate a percept through $see(\ldots)$.

- Update internal state via $next$ function,

- Select action via $action(\ldots)$.

- Perform action.

# Tasks for Agents

- We build agents in order to carry out *tasks* for us.

- The task must be *specified* by us. . .

- But we want to tell agents what to do *without* telling them how to do it.

## Utilities Functions over States

- One possibility: associate *utilities* with individual states — the task of the agent is then to bring about states that maximise utility.

- A task specification is a function

$$u : E \to \mathbb{R}$$

which associated a real number with every environment state.

- But what is the value of a *run*...

  – minimum utility of state on run?

  – maximum utility of state on run?

  – sum of utilities of states on run?

  – average?

- Disadvantage: difficult to specify a *long term* view when assigning utilities to individual states.

  (One possibility: a *discount* for states later on.)

## Utilities over Runs

- Another possibility: assigns a utility not to individual states, but to runs themselves:

$$u : \mathcal{R} \to \mathbb{R}$$

- Such an approach takes an inherently *long term* view.

- Other variations: incorporate probabilities of different states emerging.

# Problems with Utility-based Approaches

- "Where do the numbers come from?" (Peter Cheeseman)

- People don't think in terms of utilities — it's hard for people to specify tasks in these terms.

- Nevertheless, works well in certain scenarios. . . .

## Utility in the Tileworld

- Simulated two dimensional grid environment on which there are agents, tiles, obstacles, and holes.

- An agent can move in four directions, up, down, left, or right, and if it is located next to a tile, it can push it.

- Holes have to be filled up with tiles by the agent. An agent scores points by filling holes with tiles, with the aim being to fill as many holes as possible.

- TILEWORLD changes with the random appearance and disappearance of holes.

## Utility in the Tileworld

- Utility function defined as follows:

$$u(r) \;\hat{=}\; \frac{\text{number of holes filled in } r}{\text{number of holes that appeared in } r}$$

- Thus:
  if agent fills *all* holes, utility = 1.
  if agent fills *no* holes, utility = 0.

## Expected Utility

- Write $P(r \mid Ag, Env)$ to denote probability that run $r$
  occurs when agent $Ag$ is placed in environment $Env$.
  Note:

$$\sum_{r \in \mathcal{R}(Ag,Env)} P(r \mid Ag, Env) = 1.$$

- The *expected utility* of agent $Ag$ in environment $Env$
  (given $P$, $u$), is then:

$$EU(Ag, Env) = \sum_{r \in \mathcal{R}(Ag,Env)} u(r)P(r \mid Ag, Env). \quad (1)$$

$$\boxed{\textsf{An Example}}$$

Consider the environment $Env_1 = \langle E, e_0, \tau \rangle$ defined as follows:

$$E = \{e_0, e_1, e_2, e_3, e_4, e_5\}$$

$$\tau(e_0 \xrightarrow{\alpha_0}) = \{e_1, e_2\}$$

$$\tau(e_0 \xrightarrow{\alpha_1}) = \{e_3, e_4, e_5\}$$

There are two agents possible with respect to this environment:

$$Ag_1(e_0) = \alpha_0$$

$$Ag_2(e_0) = \alpha_1$$

# The probabilities of the various runs are as follows:

$P(e_0 \xrightarrow{\alpha_0} e_1 \mid Ag_1, Env_1) = 0.4$

$P(e_0 \xrightarrow{\alpha_0} e_2 \mid Ag_1, Env_1) = 0.6$

$P(e_0 \xrightarrow{\alpha_1} e_3 \mid Ag_2, Env_1) = 0.1$

$P(e_0 \xrightarrow{\alpha_1} e_4 \mid Ag_2, Env_1) = 0.2$

$P(e_0 \xrightarrow{\alpha_1} e_5 \mid Ag_2, Env_1) = 0.7$

# Assume the utility function $u_1$ is defined as follows:

$u_1(e_0 \xrightarrow{\alpha_0} e_1) = 8$

$u_1(e_0 \xrightarrow{\alpha_0} e_2) = 11$

$u_1(e_0 \xrightarrow{\alpha_1} e_3) = 70$

$u_1(e_0 \xrightarrow{\alpha_1} e_4) = 9$

$u_1(e_0 \xrightarrow{\alpha_1} e_5) = 10$

# What are the expected utilities of the agents for this utility function?

# Optimal Agents

- The optimal agent $Ag_{opt}$ in an environment $Env$ is the one that *maximizes expected utility*:

$$Ag_{opt} = \arg \max_{Ag \in AG} EU(Ag, Env) \qquad (2)$$

- Of course, the fact that an agent is optimal does not mean that it *will* be best; only that *on average*, we can expect it to do best.

## Bounded Optimal Agents

- Some agents cannot be implemented on some computers

- Write $AG_m$ to denote the agents that can be implemented on machine (computer) $m$:

$$AG_m = \{Ag \mid Ag \in AG \text{ and } Ag \text{ can be implemented on } m\}.$$

- The *bounded optimal* agent, $Ag_{bopt}$, with respect to $m$ is then...

$$Ag_{bopt} = \arg \max_{Ag \in AG_m} EU(Ag, Env) \qquad (3)$$

# Predicate Task Specifications

- A special case of assigning utilities to histories is to assign 0 (false) or 1 (true) to a run.

- If a run is assigned 1, then the agent *succeeds* on that run, otherwise it *fails*.

- Call these *predicate task specifications*.

- Denote predicate task specification by $\Psi$:

$$\Psi : \mathcal{R} \rightarrow \{0, 1\}$$

## Task Environments

- A *task environment* is a pair $\langle Env, \Psi \rangle$, where $Env$ is an environment, and

$$\Psi : \mathcal{R} \rightarrow \{0, 1\}$$

is a predicate over runs.

Let $\mathcal{TE}$ be the set of all task environments.

- A task environment specifies:

  – the properties of the system the agent will inhabit;

  – the criteria by which an agent will be judged to have either failed or succeeded.

- Write $\mathcal{R}_{\Psi}(Ag, Env)$ to denote set of all runs of the agent $Ag$ in environment $Env$ that satisfy $\Psi$:

$$\mathcal{R}_{\Psi}(Ag, Env) = \{r \mid r \in \mathcal{R}(Ag, Env) \text{ and } \Psi(r) = 1\}.$$

- We then say that an agent $Ag$ succeeds in task environment $\langle Env, \Psi \rangle$ if

$$\mathcal{R}_{\Psi}(Ag, Env) = \mathcal{R}(Ag, Env)$$

# The Probability of Success

- Let $P(r \mid Ag, Env)$ denote probability that run $r$ occurs if agent $Ag$ is placed in environment $Env$.

- Then the probability $P(\Psi \mid Ag, Env)$ that $\Psi$ is satisfied by $Ag$ in $Env$ would then simply be:

$$P(\Psi \mid Ag, Env) = \sum_{r \in \mathcal{R}_{\Psi}(Ag, Env)} P(r \mid Ag, Env)$$

## Achievement & Maintenance Tasks

- Two most common types of tasks are *achievement tasks* and *maintenance tasks*:

1. *Achievement tasks* Are those of the form "achieve state of affairs $\phi$".

2. *Maintenance tasks* Are those of the form "maintain state of affairs $\psi$".

- An achievement task is specified by a set $G$ of "good" or "goal" states: $G \subseteq E$.

  The agent succeeds if it is guaranteed to bring about at least one of these states (we do not care which one — they are all considered equally good).

- A maintenance goal is specified by a set $B$ of "bad" states: $B \subseteq E$.

  The agent succeeds in a particular environment if it manages to *avoid* all states in $B$ — if it never performs actions which result in any state in $B$ occurring.

## Agent Synthesis

- *Agent synthesis* is automatic programming: goal is to have a program that will take a task environment, and from this task environment automatically generate an agent that succeeds in this environment:

$$syn : \mathcal{TE} \rightarrow (AG \cup \{\perp\}).$$

(Think of $\perp$ as being like `null` in JAVA.)

# Soundness and Completeness

- Synthesis algorithm is:

  – *sound* if, whenever it returns an agent, then this agent succeeds in the task environment that is passed as input; and

  – *complete* if it is guaranteed to return an agent whenever there exists an agent that will succeed in the task environment given as input.

- Synthesis algorithm *syn* is sound if it satisfies the following condition:

$$syn(\langle Env, \Psi \rangle) = Ag \text{ implies } \mathcal{R}(Ag, Env) = \mathcal{R}_\Psi(Ag, Env).$$

and complete if:

$$\exists Ag \in AG \text{ s.t. } \mathcal{R}(Ag, Env) = \mathcal{R}_\Psi(Ag, Env) \text{ implies } syn(\langle Env, \Psi \rangle) \neq \bot.$$