

DARE*: An Emotion-based Agent Architecture

Márcia Maçãs, Rodrigo Ventura, Luis Custódio and Carlos Pinto-Ferreira

Institute of Systems and Robotics, Instituto Superior Técnico, Portugal

{marcia,yoda,lmmc,cpf}@isr.ist.utl.pt

http://www.isr.ist.utl.pt/~yoda/~aims

Abstract

The role of emotions in intelligent behaviour has often been discussed: are emotions an essential part of the human intelligence machinery? Recent research on the neurophysiology of human emotions suggests that human decision-making efficiency depends deeply on the emotions mechanism. In particular, António Damásio has proposed that alternative courses of action in a decision-making problem are emotionally (somatic) marked as good or bad and only the positive ones are considered for further reasoning and decision purposes. In this paper an emotion-based agent architecture supported on the Damásio's somatic marker hypothesis is presented.

Introduction

The way human beings use rules with exceptions, handle exceptions, generalize from instances, realize what is relevant, abstract what is accessory, focus attention into what is important for himself, and interact with others, is usually associated with the concept of intelligence and has been a research subject in many scientific communities, particularly in the Artificial Intelligence (AI) one. The hypothesis that intelligent behaviour can be artificially achieved through logical reasoning has been pursued by many researchers (McCarthy 1979) and simultaneously disputed by others (Brooks 1991). Although there are many important contributions in this AI research direction, most of the main problems involved when dealing with human-like decision-making have not been satisfactorily solved. On the other hand, the symbolic AI critics argue that human daily life requires a level of flexibility and efficiency that is not usually compatible with pure (logical) reasoning.

Therefore, some researchers have proposed completely different approaches, for instance the connectionism research field, and others are trying to study some human behaviour characteristics that may be used to help a reasoning mech-

anism to yield more efficient results. The research on the emotions topic is an example of this perspective.

Recent research findings on the neurophysiology of human emotions suggest that human decision-making efficiency depends deeply on the emotions machinery. In particular, the neuroscientist António Damásio (Damásio 1994) claims that alternative courses of action in a decision-making problem are previously (somatic) marked as good or bad, based on an emotional evaluation, and afterwards only the positive ones (a smaller set) are used for further reasoning and decision purposes. This constitutes the essence of the Damásio's somatic marker hypothesis. It has also been proposed that emotions may establish a link between the human reactive systems, which are responsible for instinctive reactions, and the higher level cognitive systems responsible for reasoning (Ventura & Pinto-Ferreira 1999).

The purpose of our research is to study new methodologies for developing agents based on the Damásio's hypothesis. An emotion-based agent is an entity whose behaviour is guided by taking into account a rough evaluation of the goodness or badness of the current stimulus and also an identification of the stimulus based on past experiences.

Our work is supported on the processing of external stimuli under a double perspective: a perceptual, immediate one, which allows the agent to quickly react to urgent situations, and a cognitive, elaborate one, which allows the agent to identify what is happening based on what it already knows about the world. At the perceptual level, the information extracted from the stimulus¹ is simple, basic and easily handled based on a set of built-in characteristics, which allows a fast assessment of the stimulus (is it positive/negative, desirable/avoidable, relevant/irrelevant, urgent/not urgent, and so on.). At the cognitive level, a more complex, rich, divisible, structured and hardly handled information² is extracted based on visual, audio and other sensors. This representation should be sufficient for a comprehensive identification of the stimulus by the higher level cognitive systems.

Using these two sets of information, perceptual and cognitive, the architecture incorporates a marking mechanism that

*DARE stands for (inverse order) "Emotion-based Robotic Agent Development". This work has been supported by the project PRAXIS/C/EEI/12184/1998 entitled "DARE: Study of Methodologies for the Development of Emotion-based Robotic Agents" funded by the FCT - Portuguese Foundation for Science and Technology under the PRAXIS program.
Copyright © 2001, American Association for Artificial Intelligence (www.aaai.org). All rights reserved.

¹Hereafter called perceptual image. In this context, an image is a generalized representation of information got from all available sensors (vision, audio, smell, etc.).

²Hereafter called cognitive image.

Cognitive Evaluation The cognitive evaluation differs from the perceptual in the sense that it uses past experience, stored in memory. The idea is to retrieve from memory the DV associated with cognitive images similar to the present stimulus. Recall that a cognitive image is a stimulus representation including all significant elements of the stimulus; a possible example of a cognitive representation of a visual stimulus is the corresponding bitmap. So, two such stimulus can be compared using the respective bitmaps and an adequate pattern matching method.

This process allows the agent to use past experience in decision making. After obtaining the perceptual and cognitive images for the current stimulus, when the evaluation of the I_p does not reveal urgency, *i.e.*, the resulting DV is not so imperative that would demand an immediate response, a cognitive evaluation is performed. It consists in using the I_p as a memory index⁴ to search for past obtained cognitive images similar to the present I_c . This similarity is measured using all the features that form the cognitive images, working as a pattern matching process. Each I_c in memory, besides having an associated I_p , also has the resulting DV from past evaluation. If the agent already has been exposed to similar stimulus in the past, then it will recall its associated DV, being this the result of the cognitive evaluation. This means that the agent associates with the current stimulus the same desirability that is associated with the stimulus in memory. If the agent has never been exposed to a similar stimulus, no similar I_c will be found in memory, and therefore no DV will be retrieved. In this case the cognitive evaluation does not produce any match and the DV coming from the perceptual evaluation is the one which will be used for the rest of the processing.

Body State Evaluation In this architecture, the notion of body corresponds to an internal state of the agent, *i.e.*, a set of pre-defined variables. The body state consists in the values of the state variables at a particular moment (*e.g.* the level of nutrients available in the body). The internal state may change due to the agent's actions or by direct influence of the environment.

The stimulus assessment proceeds to a body state evaluation, given the DV from the previous evaluation step and a built-in tendency of the body's agent. The innate tendency defines the set of body states considered ideal for the agent. It includes the definition of the Homeostatic Vector (HV) containing the equilibrium values of the state variables. The built-in tendency can be oriented towards the maintenance of those values or the maximization/minimization of some of them. In other words, this comprises a representation of the agent's needs. If the value of the nutrients variable is below the corresponding HV value, this suggests that the agent needs food.

Based on the DV determined for the current stimulus, the body state evaluation consists in an estimation of the possible effects of the alternative courses of action. This corre-

sponds to an anticipation of the possible action outcomes for the agent in the future: "will this action help to re-balance a particular unbalanced body variable, or will get it even more unbalanced?". Given this anticipation, the DV previously obtained may be altered, reflecting the desirability of the stimulus according to the agent's current needs. For instance, if the anticipation approximates the body state to the HV, then the desirability associated with the stimulus must increase.

Considering explicitly the body state and its influence on the stimulus evaluation, it allows to introduce some flexibility in the agent's behaviour. As the agent's decisions depend on a particular body state — the one existing when the agent is deciding — it will not respond always in the same manner to a given stimulus. On the other hand, the existence of a body state forces the agent to behave with pro-activeness, because its internal state "drives" its actions, and autonomy, because it does not rely on an external entity to satisfy its needs.

Decision Making

After the evaluation process being finished, the desirability of the stimulus is determined and the agent will select an adequate action to be executed. In the former step of evaluation — body state evaluation — the effects of all possible actions were anticipated based on the expected changes in the body state. The action with the best contribution for the agent's overall welfare will be selected as the one to be executed next. As interaction with the environment is essential, when being created the agent must start with a collection of basic actions, such as move forward, turn, run, jump, grab, drop, and so on. So, it is assumed here that there is a set of built-in elementary actions that the agent can execute. After the selected action being executed, the changes in the environment will generate a new stimulus to be processed.

Using the feedback from the environment the agent may learn how to compose elementary actions in order to establish new and more effective actions. Also, supposing the agent is able to perceive that a stimulus is a direct effect of a previous executed action, it is possible to implement a simple learning mechanism based on the anticipation of the action by comparing it *a posteriori* with its real consequences⁵. For instance, suppose that a starving agent finds in the world an object that has all the features characterizing food, which in fact is rotten food. Without *a priori* knowledge, the agent will anticipate the positive "effects" of eating that food and it will decide to eat it. However, the real effects will be very different: the agent's welfare may decrease dramatically. In this case, the DV previously associated with the cognitive image of that particular food will be changed, either gradually, if the food was not so bad, or radically, if it was dangerous for agent's survival. Any case, next time the agent sees that special food it will behave differently.

⁴The purpose of using perceptual information to index cognitive images in memory is to reduce search. It is here hypothesized that it is more probable to have the current I_c similar to those in memory having the same I_p (dominant features).

⁵One of the most difficult problems related with this learning issue is to establishing cause-effect relationships, *i.e.*, to perceive which action or actions, if any, have been the cause of the current stimulus (credit-assignment problem). Several ideas are being studied to further develop this action learning.

This kind of learning has been one of the experiments performed with an implementation of this architecture on a pre-defined world. For the sake of paper length, the implementation, the chosen world and the results obtained will be briefly presented in the next section.

Implementation

In order to evaluate it, an implementation of the proposed architecture was performed. Using a simulated maze environment, represented by a grid, experiments were carried out. The main goal of the implementation for this particular environment was the agent's survival while seeking the maze exit. Figure 2 shows the user interface.

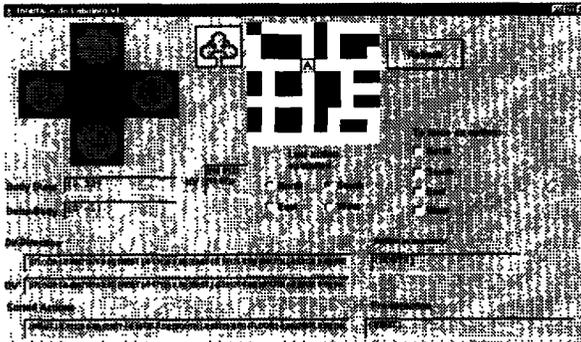


Figure 2: The interface of the maze implementation.

When the agent reaches a new position in the maze, it perceives new images (bitmaps) related with five directions. The agent can move in the maze in four directions (N, S, E, W) and when food is present (floor direction) the agent may eat it. The agent has two basic needs defined by its body state, light and energy. The need for energy arises because every action executed by the agent leads to a reduction of energy in its body, therefore the agent must seek food in the maze. Inside the maze there is light only near the exit. The need for light will make the agent search for the exit. The relevant features used in the perceptual processing are the colors of pixels in the stimulus image. Some colors are innately defined as undesirable and others are considered to be good.

The agent starts at the top left position in the maze, with a balanced level of energy and with a very unbalanced level of light. The experiments performed in this environment using the proposed architecture revealed that the agent is able to find the maze exit (therefore fulfilling its need for light), to search for food and eat it when its level of energy is below a threshold, and to identify and learn to avoid eating rotten food (that reduces energy when eaten).

With these experiments it was exemplified the learning capability introduced by changing the DVs in memory, its utilization by the cognitive evaluation, the pursue of need satisfaction using DVs information, and the decision-making based on the anticipation of action outcomes.

Related work

The discussion concerning the relevance of emotions for artificial intelligence is not new. In fact, AI researchers as Aaron Sloman (Sloman & Croucher 1981) and Marvin Minsky (Minsky 1988) have pointed out that a deeper study of the possible contribution of emotion to intelligence was needed. Recent publications of psychology (Goleman 1996) and neuroscience research results (Damásio 1994; LeDoux 1996) suggest a relationship between emotion and rational behaviour, which has motivated an AI research increase in this area. The introduction of emotions as an attempt to improve intelligent systems has been made through different approaches. Some researchers use emotions (or its underlying mechanisms) as a part of architectures with the ultimate goal of developing autonomous agents that can cope with complex dynamic environments. In this set is included the Velásquez work (Velásquez 1998a; 1998b), who developed a pet-robot based on Damásio's ideas; and an architecture for emotion-based agents proposed by Ventura, Pinto-Ferreira and Custódio (Ventura & Pinto-Ferreira 1998a; 1998b; Ventura, Custódio, & Pinto-Ferreira 1998a; 1998b; Ventura & Pinto-Ferreira 1999). Another architecture (Tabasco) was proposed by Staller and Petta (Staller & Petta 1998), which is based on psychological theories of emotions. Other researchers focused their work on the adaptation aspects of emotions, using it in reinforcement learning (Gadanhó & Hallam 1998). There are researchers who defend that emotion is a side effect of an intelligent system (Sloman 1998), others defend the opposite, *i.e.*, emotion is the basis of emergent intelligent behaviour (Cañamero 1997). The social role of emotion has also been explored by several researchers using it to improve societies of intelligent agents (Cañamero & de Velde 1999; Staller & Petta 1998).

Conclusions

The proposed architecture allowed the implementation of an autonomous agent, (i) where the goal definition results from the agent's behaviour and needs, *i.e.*, it is not imposed or pre-defined; (ii) where the agent is capable of quickly reacting to environment changes due to the perceptual level processing; (iii) where the agent reveals adaptation capabilities due to the cognitive level processing; and finally (iv) where the agent is capable of anticipating the outcomes of its actions, allowing a more informed process of decision making.

The mechanism of emotions seems to play two main roles: a *decisional*, influencing the way agents assess situations and make decisions, and a *communicational*, affecting the way agents express their internal state when confronted with their environment. In what concerns the decisional aspect — the one that is partially covered by the present evolving architecture — the relevant issue is to determine whether the framework suggested by Damásio (and interpreted by the authors) helps in the development of more competent agents.

Future Work

This architecture is now being tested in more complex and dynamic environments, namely using real robots interacting with semi-structured environments. Moreover, the maze world is being extended in order to accommodate different needs, objects and environment reactions. Furthermore, the study of the social role of emotions in a multi-agent system (e.g., the RoboCup environment), the study of non-verbal reasoning mechanisms (e.g., pictorial reasoning) and its relation with the emotion processes, have been addressed.

References

- Brooks, R. A. 1991. Intelligence without representation. *Artificial Intelligence* 139–159.
- Cañamero, D., and de Velde, W. V. 1999. Socially emotional: using emotions to ground social interaction. In Dautenhahn, K., ed., *Human Cognition and Social Agent Technology*. John Benjamins Publishing Company.
- Cañamero, D. 1997. Modeling motivations and emotions as a basis for intelligent behaviour. In Johnson, W. L., ed., *Proc. of the first International Conference on Autonomous Agents*, 148–155. ACM Press.
- Damásio, A. 1994. *Descartes' Error: Emotion, Reason and the Human Brain*. Picador.
- Gadano, S., and Hallam, J. 1998. Emotion-triggered learning for autonomous robots. In Cañamero, D., ed., *Emotional and Intelligent: the Tangled Knot of Cognition. 1998 AAAI Fall Symposium*. AAAI. 31–36.
- Goleman, D. 1996. *Emotional Intelligence*. Bloomsbury.
- LeDoux, J. 1996. *The Emotional Brain*. Simon and Schuster.
- McCarthy, J. 1979. Ascribing mental qualities to machines. In Ringle, M., ed., *Philosophical Perspectives in Artificial Intelligence*. Oxford: Harvester Press.
- Minsky, M. 1988. *The society of mind*. Touchstone.
- Sloman, A., and Croucher, M. 1981. Why robots will have emotions. In *Proc. 7th Int. Joint Conference on AI*, 197–202.
- Sloman, A. 1998. Damasio, descartes, alarms and meta-management. In *Proc. on Int. Conference on systems, Man, and Cybernetics (SMC98)*, 2652–7. IEEE.
- Staller, A., and Petta, P. 1998. Towards a tractable appraisal-based architecture. In Cañamero, D.; Numaoka, C.; and Petta, P., eds., *Workshop: Grounding Emotions in Adaptive Systems Int. Conf. of Simulation of Adaptive Behaviour, from Animals to Animats, SAB'98*. 56–61.
- Velásquez, J. 1998a. Modeling emotion-based decision-making. In Cañamero, D., ed., *Emotional and Intelligent: the Tangled Knot of Cognition. 1998 AAAI Fall Symposium*. AAAI. 164–169.
- Velásquez, J. 1998b. When robots wheep: Emotional memories and decision-making. In *Proc. of AAAI-98*, 70–75. AAAI.
- Ventura, R., and Pinto-Ferreira, C. 1998a. Emotion-based agents. In *Proceedings of AAAI-98*, 1204. AAAI.
- Ventura, R., and Pinto-Ferreira, C. 1998b. Meaning engines - revisiting chinese room. In Cañamero, D.; Numaoka, C.; and Petta, P., eds., *Workshop: Grounding Emotions in Adaptive Systems Int. Conf. of Simulation of Adaptive Behaviour, from Animals to Animats, SAB'98*. 69–70.
- Ventura, R., and Pinto-Ferreira, C. 1999. Emotion based agents: Three approaches to implementation. In Velásquez, J., ed., *Workshop on Emotion-based Agents Architectures, EBAA'99*. 121–129.
- Ventura, R.; Custódio, L.; and Pinto-Ferreira, C. 1998a. Artificial emotions - goodbye mr. spock! In *Proc. of the 2nd Int. Conf. on Cognitive Science*, 938–941.
- Ventura, R.; Custódio, L.; and Pinto-Ferreira, C. 1998b. Emotions - the missing link? In Cañamero, D., ed., *Emotional and Intelligent: the Tangled Knot of Cognition. 1998 AAAI Fall Symposium*. AAAI. 170–175.