New Challenges for Character-Based AI for Games

Damián Isla and Bruce Blumberg

The Synthetic Characters Group, MIT Media Lab 20 Ames St. Cambridge, MA 02139 {naimad, bruce}@media.mit.edu

Abstract

Based on the recent surge in interest in both academia and the games industry in character-based AI, and inspired by work from a range of literature from animal learning and psychology to classic AI, we suggest some potential "next steps" towards the creation of virtual autonomous characters that are lifelike, intelligent and convey empathy. We believe that character-based AI systems can draw from *and* contribute to the state of the art in such areas as pattern recognition and planning.

Introduction: AI's Next Steps

First things first: what is meant by character-based AI?

It is a category that is especially meaningful in games, to distinguish systems that seek to simulate the behavior of a single agent from strategic AIs or turn-based game-opponent AIs. These latter categories might be considered attempts to codify and emulate high-level logical human thinking. Character-based AI, on the other hand, is an exercise in creating *complete* brains. Strategic and logical thinking in this type of work usually takes a back seat to issues of low-level perception, reactive behavior and motor control. Since the creatures in these character-based systems often have graphical bodies (a sort of *virtual embodiment*) the work is often rendered with an eye toward recreating life-like behavior, and emotion-modelling and robustness are often also central issues.

There are many examples of character-based AI work in both academia and the computer game industry. Well known academic projects include the Biomechanical fish of Tu and Terzopoulos (Tu & Terzopoulos 1994), Bates' Oz Project on interactive story-telling (Bates 1993) and Hayes-Roth's Virtual Theater Project (Hayes-Roth *et al.* 1995). The work of Blumberg ((Blumberg 1994) and (Blumberg 1996))and his Synthetic Characters Group at the MIT Media Lab has for some years explored the application of animal learning and psychology to computational decision-making systems in the form of virtually-embodied animals.

From the computer games industry there has been a recent surge in interest in character-based AI. An early example of this was the game *Creatures* (Grand, Cliff, & Malhotra 1997), whose "Norns" were driven by simple synthetic perception, a neural-network learning mechanism and a simulated biochemical system. Two recent games, impressive in both their technical sophistication and their commercial success are Will Wright's *The Sims* and Peter Molyneux's *Black and White*, the latter of which also featured impressive learning.

All of these projects, academic and industrial, show a concern for the same general issues: how do the simulated creatures perceive their world, and how is their perception of the word-state limited realistically? How are the creature's reactions handled, and how do they go about satisfying their goals? What *are* their goals?

Having gotten to where we are, the principle question we now face is: where do we go from here?

This is a singularly appropriate time to be asking this question, since both academia and industry are closer than ever to be being in a position to attempt an answer. Graphics and animation, long a stumbling block for academic researchers, have finally become cheap, due not only to the recent abundance of graphics-related tools and platforms ((Laird 2000), for example, was implemented on top of the Quake engine) but also to the recent abundance of graphics and animation talent. At the same time, game-designers are finally getting to the point where they can afford to spend significant numbers of execution cycles on AI-processing. Perhaps most crucially of all, the game industry now provides an economic impetus to conduct this type of research. Computer games, after all, sell.

This paper makes some concrete – and very subjective – suggestions about where character-based AI, whether in academia or in the games industry, can go. At the heart of all these suggestions lies the important idea that tackling the many problems of intelligent behavior in an integrated way – in a single brain that handles everything from perception to memory to action selection – makes those problems easier to solve, and results in more compelling behavior.

The topics that are covered in this paper might prompt some to observe that game AI is fated to recapitulate the development of classical academic AI. We agree with this assessment, and consider it a necessary but fruitful process. For in recapitulating classical AI from a new perspective and with different goals, the ultimate effect will be not just recapitulation but reinvention.

Copyright © 2002, American Association for Artificial Intelligence (www.aaai.org). All rights reserved.

Perception

It is through perception that a character receives input from the world. It is on the basis of this input that the character will make decisions and act. As we will discuss, there are a number of interesting issues that arise from the design of a character's perceptual system.

Sensory Honesty

One of the fundamental issues of perception is "what should the character be able to perceive?" As the complexity of our characters' brains has increased, we have found it increasingly useful to treat the mechanisms of perception in a more principled manner. This entails an attempt at realism in terms of both *what* things a character can perceive (visual events that occur behind a wall should not be perceived) and *how* it perceives them (locations of visual events might be passed to a character in that character's eyecoordinates, rather than in world-coordinates to which the character should rightfully have no access). We call this the principle of *sensory honesty*.

A fundamental aspect of sensory honesty is that it forces a separation between the actual state of the world and a particular character's view of the state of the world. The most obvious reason is that a character's perception can only cover a small part of that world-state (Black and White's creatures, for example, cannot see events that take place behind them). By keeping separate the state of the world and the character's beliefs about that state, we open the door for these two representations to occasionally diverge. This forces a character to make intelligent assumptions about aspects of the worldstate that it cannot directly observe. A character should not lose track of an object that it has been observing, simply because it has been hidden by an occluder. Rather, a reasonable guess as to the object's location should be maintained. As will be shown in later sections (for example, Anticipation and Imagination) this ability to maintain reasonable assumptions can in turn lead to some very interesting and life-like behavior.

Pattern Recognition

Perception also subsumes the larger problem of pattern recognition, where the patterns in question might be visual scenes, sounds, gestures, etc. In (Tu & Terzopoulos 1994), the artificial fish recognize other fish through shape recognition and extraction from a rendering of the fishtank from each individual's point of view. In *Sheep—Dog*, a virtual sheep-herding installation built by the Synthetic Characters Group, the user (playing the part of a shepherd) would communicate with Duncan (the autonomous sheepdog) through a voice interface. In both these cases, the pattern recognition mechanism is integrated as a part of the creature's perception system.

Luckily there are many existing algorithms for all forms of pattern recognition that can be used for these types of problems – the challenge, and a very surmountable one, is simply to provide an architecture into which these algorithms can be seamlessly integrated. The structure that we use is a Percept Tree (Burke et al. 2001).

The utility of these algorithms might be increased by taking into account high-level feedback from the rest of the brain. Two forms that this feedback might take include:

- Reward: Pattern classifiers are data-driven subdivisions of a multi-dimensional perceptual input space. While there exist unsupervised techniques for automatically finding useful divisions of this space, some amount of supervision can usually improve an algorithm's performance by indicating roughly how important certain regions of the input space are. In the case of an embodied character, that supervision might take the form of reward (the reward depending on the character's success at satisfying its drives). Thus the classification mechanism might model important regions of the input space (i.e. regions critical to the discrimination of reward and non-reward situations) at a higher resolution than other regions that are unimportant in predicting reward (though the data itself might support high-resolution modelling of the region). The result is a subdivision of the input space that is both data-driven and reward-driven. See (Burke et al. 2001) for a description of the Sheep-Dog acoustic pattern matcher that operates on this principle using a hierarchical clustering technique.
- **Priming:** Priming corresponds roughly to the incorporation of high-level expectation information into the lowlevel recognition process. These expectations cause perception to be more forgiving of ambiguity for certain strongly-anticipated sensory experiences, resulting in *a priori* biases towards certain data classifications.

If, for example, a certain utterance has been repeated over and over again recently, then the introduction of a new utterance that would normally be too vague to classify might nonetheless be taken as a new instance of the same utterance. Note that the converse phenomenon is also acceptable – that perceptual categories that are long out of use might take longer to function, just as it might take us slightly longer to recognize a friend whom we have not seen for years (and whom we do not *expect* to see) than to recognize a friend we see every day.

The important lesson is that expectations and emotions can strongly influence the perceptual experience of a character. Thus the paranoid character jumps in fright at the slightest motion, and the character in love momentarily mistakes every woman he passes as the object of his affection. Priming is an interesting phenomenon to explore for Character-based AI systems precisely because it requires a behavioral/emotional infrastructure that an abstract machine learning system does not provide.

It is notable that these two feedback channels are related in their overall effect, since the first leads to long-term perceptual adaptation and the second is a form of short-term perceptual adaptation.

Anticipation and Imagination

A large part of appearing intelligent is not only the ability to be predictable (in the sense of Dennett's *intentional stance*, (Dennet 1987)) but also the ability to *form* predictions, and to act in anticipation of those events predicted. If an animal is attacked by a predator whenever it visits a given food source, the animal will, before long, begin avoiding that food source. Similarly, if a character is hit in the head by a flying brick every time it opens the window, then it would simply appear *broken* if the character continued to open the window with no anticipation of being hit. As can be imagined the capacity to form predictions can be essential in maintaining the appearance of common sense.

Note that these rules about what should be anticipated given a sequence of events could be hand-coded (as was done in (Laird 2000)). In other cases, we would like our characters to form these types of expectations on the fly, especially if we foresee the player being the source of some of these event correlations. If the player always acts in a certain predictable way in a certain situation, we might like the simulated character to pick up on that and react appropriately.

How do we go about forming predictions about the future? If the specific quality or state that we are tracking takes the form of a scalar or vector, we might use standard function-extrapolation techniques to predict future values in the near-term. Those values might be bounded or modified by higher-level concepts of what we know to be reasonable (even if an object was last observed flying upwards, we should still expect it to eventually come down).

One very important form of expectation formation is provided by the classical or Pavlovian conditioning paradigm. If a stimulus A reliably precedes stimulus B, where B is some salient event (such as the appearance of food), then the appearance of A should be enough to cause the anticipation of B. In Pavlov's famous example, a group of dogs were conditioned to expect food upon the ringing of a bell (the expectation being shown by salivation in response to the ringing of the bell). (Burke 2001) describes a system in which statistically reliable correlations between salient stimuli are used to infer apparent causality rules. These expectations are represented by cause-effect pairs with an associated temporal delay (if the causal conditions are observed, the effect conditions are expected after the temporal delay). Thus if every time the button is pushed (causal condition) the elevator doors open (the effect condition) then the character can form a very specific expectation that the elevator doors will open when it sees the button pushed, and it can react appropriately before the fact (approach if it wants to use the elevator, flee if it knows that there is a lion inside, etc.).

As already discussed in the section *Sensory Honesty, any* form of expectation – an assumption about something that will happen in the world – can sometimes prove false. When this occurs, the result is an *expectation violation*. These violations can play an extremely important role in focusing a character's attention (such that an explanation of why the expectation turned out false might be sought) and in refining cause-effect hypotheses. (Kline 1999) contains an excellent discussion of expectation theory.

Note also that the ability to anticipate and occasionally be mistaken has the potential to open up a whole new range of potential interactions. A character that can make mistakes (and be surprised when it does so) is a character that can be tricked, teased, out-smarted, or even lied to. The hope is that this could add an additional layer of sophistication to the autonomous characters, friend or foe, that a player might encounter.

Imagination, or Planning

In confronting a problem, we often think through multiple scenarios before actually attempting a solution. This amounts to running simplified simulations in our heads. These simulations are governed by the expected results of actions that we consider taking, expectations based both on common sense and on observed reliable cause-effect correlations (as described above). If the elevator doors reliably opened whenever the button was pushed, then that cause-effect correlation can probably be incorporated into a plan. There already exist, of course, many algorithms for performing planning. Simply incorporating these existing technologies is problematic, however, due to run-time considerations – with problems of any non-trivial complexity, the cause-effect search spaces become enormous and hence very difficult to search efficiently.

Planning in character-based AIs might be made easier by bringing to bear some constraints of character and situation to prune the search space. Here, the psychology of perception and attention can be used to expand only regions of the tree that seem relevant. Emotional state might also influence the search path, perhaps resulting in "optimistic" and "pessimistic" plans, depending on the character's state of mind. Just as evolution has provided us with certain "cheats" for focusing our mental power, embodied characters with perceptual, attentional, emotional and motivational models should be able to find convenient ways to avoid searching likely uninteresting regions of the search space. Clearly, the plans that result will sometimes not be optimal. But then so are many of the plans that *we* come up with.

Theory of Mind

The plans and expectations that a synthetic character might form would be further improved through special treatment of a special class of objects in the world whose behavior is less easy to predict based on observed event-correlations, namely *other* characters. In this case however, we have another source of predictive power available to us, namely the character's *own* behavior given a context and emotional state. In other words, the character can make use of a theory of mind, predicting the actions of other characters by in effect imagining "what would I do in that situation, given the state of mind that I presume that character to be in?". Incorporating theory of mind into the types of plans described in the previous section would amount to a sort of character-based Alpha-Beta algorithm (see (Nilsson 1998) for an overview).

A theory of mind might be useful in other types of problems as well. It might be used in inferring intentionality ("Character B kicked a lamppost. I only kick lampposts when I am angry. Therefore Character B must be angry") or in imitation learning ("Character B has kicked the lamppost and Character B is angry. Perhaps when I am angry, kicking the lamppost will help"). Theory of mind represents a significant jump in sophistication in terms of how a character deals with predictions for objects in the world. Clearly, it should not be applied to *all* objects, since the results would be both inaccurate and computationally expensive. On the other hand, one of the important aspects of the *intentional stance* (Dennett's take on theory of mind (Dennet 1987)) is that it *can* occasionally be applied fruitfully to objects that we know do not have minds of their own – the point being that as long as the attribution of intentionality is *useful* in predicting the object's behavior, then it is as good a model as any. Initially, we expect, it will be convenient to tag explicitly objects in the world for which a theory of mind should be assumed. However, it would be interesting in the long-run to see whether a character can *learn* which types of objects should be so treated.

Emotion Modelling

In much of the work done in characters, emotion modelling has been used primarily as a diagnostic channel. Emotions, after all, are convenient indicators of overall system state. Emotional levels are easily routed to the facial or emotionparameterized animation engine in order to let both developers and users know what, at a high-level, is going on inside a character.

However, emotions clearly play a far larger role in our behavior than simply coloring our physical motion and stance. Emotions influence the way that we make decisions, the way we think about and plan for the future and even the way we perceive the world.

Similarly, for synthetic characters, emotions should be used to their full effect, influencing the way that actionselection functions, the way that plans are formed, the salience of different sensory stimuli, and so on. Levels of frustration, for example, could be used as a signal to the action-selection system that "this strategy is not working, try another". Since most action-selection schemes incorporate some form of anti-dithering mechanism (to keep a character from flipping too rapidly between alternative actions), levels of "urgency" or "desperation" might influence the degree of dithering allowed (a desperate character might be expected to dither more in its choice of actions) just as they might influence the depth to which a multi-step sequence of actions might be planned (and thus the speed with which an answer is found). A very interesting use of emotions is found in (Burke 2001), where a *curiosity* emotion is used to influence exploration/exploitation decisions made by the character (a curious character is one more likely to test out its ideas about action-stimulus correlations, whereas a non-curious one will tend to pick at each moment the action that is expected to render the greatest amount of reward).

Beyond Happiness and Sadness

Curiosity, mentioned in the example above, is a good example of a *secondary* emotion. These emotions go beyond the traditional happiness/sadness models (such as in (Ekman 1992)) to express more subtle aspects of a character's mental state. In this case, curiosity represents the character's overall willingness to experiment.

There are a number of emotions that might be derived from a system that can form expectations. (Isla 2001), which describes a system for forming and maintaining expectations about location, notes that emotions such as surprise (when an unexpected value is observed) and confusion (when an expected value is observed to be false) can be derived from explicit representations of expectations. These emotions could further be modulated according to emotional attitudes towards the objects whose location are being anticipated, such that "surprise" can become "fear surprise" or "pleasant surprise", and "confusion" can become "relief" or "worry". All these emotions, again, should have effects not just on the quality of the character's motion, but also on the character's basic decision-making processes.

Learning and Memory

The heading of this section is somewhat of a misnomer, since almost every section so far has dealt in some way with the problem of learning. As should already be obvious, "learning" does not refer to any one process, but rather to many individual adaptation processes that occur throughout a character's brain. "Learning at multiple levels" is a primary feature of the C4 brain (Isla *et al.* 2001), where perceptual refinement, behavioral adaptation, environment mapping and motor learning all act concurrently. Another example are the creatures of *Black and White*, who learn at the same time "how to be nice" and who to be nice *to* (Evans 2001). Both of these models go significantly beyond the classic reinforcement learning paradigm of "action A in state S leads to reward R and state S".

One relationship that has yet to be explored thoroughly is the one between learning and explicit memory formation.

Learning through Episodic Memory

There are many forms of memory (see (Reisberg 1997) for an overview). Procedural memory allows us to practice and improve specific skills over time. Short-term memory is a recent perceptual history of the world and working memory is a slightly higher-level recent history of the objects relevant to the current task or behavioral context. The colloquial use of the word "memory" refers to long-term or Episodic Memory, in which explicit sequences of events are stored for later conscious recall. These memories often detail archetypal event-sequences whose outcomes can be used to inform subsequent similar experiences. This is, of course, the very definition of learning. However, few behavior simulation systems to date have made explicit use of episodic memory as a learning mechanism.

Learning through episodic memory might be considered a variant on the observation-based rule inference of the section on *Anticipation and Imagination*. However, rather than simply keeping statistical models of causes, effects and cause-effect probabilities, an Episodic memory-based mechanism might keep around specific example *episodes*. These episodes would be exemplary event sequences that would be used to predict future sequences with similar starting conditions. Thus if the last time the button was pushed the elevator doors opened, then next time the door must be opened, pushing the button might be a good way to do this.

The advantage of this type of learning is speed: useable hypotheses about causality could be formed after just one observation. Furthermore, the causality models that are formed could be accompanied by specific remembered instances in which the hypothesis succeeded or failed, thereby providing a pool of data to support generalization or discrimination of the cause and effect models.

A number of issues quickly arise when we begin to consider an episode-based learning mechanism: how are episodes recognized? How are they matched? How are the circumstances of one episode generalized to affect similar episodes? When can we safely *forget* an episode? These questions can likely only be answered by building more characters with functional episodic memory capabilities and seeing what kinds of heuristics are useful.

Conclusions

While the authors are excited about the developments that lie on the horizon, there are three crucial questions that have yet to be addressed:

- How are these capabilities integrated into a game? To what extent must they be "designed for"?
- How do they impact gameplay? What new varieties of gameplay are possible?
- Do they result in "deeper" gameplay experiences?

Indeed, this last question represents best the motivation for following this line of research at all. Unfortunately, we must leave its definitive answer to our colleagues in the game industry, whom we wish the best of luck.

Acknowledgements

Thanks to the members of the Synthetic Characters Group: Matt Berlin, Robert Burke, Marc Downie, Scott Eaton, Yuri Ivanov, Michael Johnson, Jesse Grey, Chris Kline, Ben Resner and Bill Tomlinson.

References

Bates, J. 1993. The nature of characters in interactive worlds and the oz project. In Loeder, C., ed., *Virtual Realities: Anthology of Industry and Culture.*

Blumberg, B. 1994. Action selection in hamsterdam: Lessons from ethology. In Cliff, D.; Husbands, P.; Meyer, J.; and Wilson, S., eds., *Proceedings of the Third International Conference on Simulation of Adaptive Behavior*, number 3 in From Animals to Animats. MIT Press.

Blumberg, B. 1996. Old Tricks, New Dogs: Ethology and Interactive Creatures. PhD dissertation, MIT Media Lab.

Burke, R.; Isla, D.; Downie, M.; Ivanov, Y.; and Blumberg, B. 2001. CreatureSmarts: The art and architecture of a virtual brain. In *Proceedings of the Game Developers Conference*.

Burke, R. 2001. It's about time: Temporal representations for synthetic characters. Master's thesis, MIT Media Lab.

Dennet, D. 1987. *The Intentional Stance*. Cambridge, MA: The MIT Press.

Downie, M. 2001. Behavior, animation and music: The music and movement of synthetic characters. Master's thesis, MIT Media Lab.

Ekman, P. 1992. An argument for basic emotions. In Stein, N., and Oatley, K., eds., *Basic Emotions*. UK: Hove. 169–200.

Evans, R. 2001. The future of ai in games: A personal view. In *Game Developers Magazine*.

Grand, S.; Cliff, D.; and Malhotra, A. 1997. Creatures: Artificial life autonomous software agents. In *Proceedings of the First International Conference on Autonomous Agents*.

Hayes-Roth, B.; Sincoff, E.; Brownston, L.; Huard, R.; and Lent, B. 1995. Directed improvisation with animated puppets. In *Proceedings of CHI '95*. Denver, CO: CHI.

Isla, D.; Burke, R.; Downie, M.; and Blumberg, B. 2001. A layered brain architecture for synthetic creatures. In *The Proceedings of the International Joint Conference on Artificial Intelligence*. Seattle, WA: IJCAI.

Isla, D. 2001. The virtual hippocampus: Spatial common sense for synthetic creatures. Master's thesis, Massachusetts Institute of Technology.

Johnson, M. P.; Wilson, A.; Blumberg, B.; Kline, C.; and Bobick, A. 1999. Sympathetic interfaces: Using a plush toy to direct synthetic characters. In *CHI*, 152–158.

Kline, C. 1999. Observation-based expectation generation and response for behavior-based artificial creatures. Master's thesis, Massachusetts Institute of Technology.

Laird, J. 2000. It knows what you're going to do: Adding anticipation to a quakebot. In *The AAAI 2000 Spring Symposium Series: Artificial Intelligence and Interactive Entertainment*.

Nilsson, N. 1998. Artificial Intelligence: A New Synthesis. Morgan Kaufman.

Reisberg, D. 1997. *Cognition: Exploring the Science of the Mind*. W.W. Norton and Company.

Tomlinson, W., and Blumberg, B. 2001. Social behavior, emotion and learning in a pack of virtual wolves. In *AAAI Fall Symposium*. Falmouth, MA: AAAI.

Tomlinson, W. 1999. Interactivity and emotion through cinematography. Master's thesis, Massachusetts Institute of Technology.

Tu, X., and Terzopoulos, D. 1994. Artificial fishes: Physics, locomotion, perception, behavior. In *SIGGRAPH* '94 Conference Proceedings. Orlando, FL: ACM Siggraph.

Yoon, S.-Y.; Burke, R.; Blumberg, B.; and Schneider, G. 2000. Interactive training for synthetic characters. In *Proceedings of AAAI*.