

Modelling human emotions for tactical decision-making games

Gillian C. Visschedijk, Ard W. Lazonder, Anja van der Hulst, Nathalie Vink and Henny Leemkuil

Gillian Visschedijk, Anja van der Hulst and Nathalie Vink work at the Departments "Training & Performance Innovations" and "Military Operations," TNO Defence, Security and Safety. They specialise in serious games for all kinds of military training. Ard Lazonder and Henny Leemkuil work at the Department of Instructional Technology, University of Twente. They specialise in technology-enhanced learning and have a particular interest in learning with games and simulations. Address for correspondence: Ms Gillian C. Visschedijk, TNO Defence, Security and Safety, P.O. Box 23, 3769 ZG Soesterberg, The Netherlands. Email: gillian.visschedijk@tno.nl

Abstract

The training of tactical decision making increasingly occurs through serious computer games. A challenging aspect of designing such games is the modelling of human emotions. Two studies were performed to investigate the relation between fidelity and human emotion recognition in virtual human characters. Study 1 compared five versions of a virtual character that expressed emotions through different combinations of posture, facial expression, and tone of voice. Results showed that emotion recognition was best when all three behavioural cues were present; posture + face and posture + tone of voice were joint second best. In study 2, these three versions were supplemented with contextual information. Cross-variant comparisons yielded marginal differences in emotion recognition and no differences in tactical decision making. Together, these findings suggest that the combination of posture with either facial expression or tone of voice is sufficient to ensure recognition of human emotions in tactical decision-making games.

Introduction

Tactical decision making denotes the ability to choose which actions or solutions should best be taken to accomplish a goal or task. The decisions that emanate from this process can literally be of vital importance to professionals such as police officers, fire fighters, security guards and military commanders who operate under dangerous or threatening conditions. The development of tactical decision-making skills increasingly occurs through serious games that, due to advanced computer technology, enable commanders-in-training to make tactical decisions in situations that are impossible in the real world for reasons of safety, cost and time (Kiili, 2007; Knerr, 2006).

Tactical decisions are made by matching features of the current situation to previously acquired patterns, selecting the most appropriate course of action and mentally simulating how these actions would work out in the current situation (Klein, 2008). The first step in the decision-making process thus hinges on situational awareness, which often requires the ability to recognise human emotions. Imagine, for instance, a riot police officer controlling a right extremists' demonstration. His decision to either intervene or refrain from action depends almost exclusively on his assessment of the crowd's level of aggression (Moya, McKenzie & Nguyen, 2008). Gaining proficiency in human emotion recognition thus seems to require maximum training fidelity: the more realistically human emotions are modelled in a game, the better they are recognised. Paradoxically, however, high-fidelity games are particularly expensive to develop while they do

Practitioner Notes

What is already known about this topic

- Tactical decision making often involves the recognition of human emotions.
- Humans exhibit their emotions through facial expressions, body movement and posture and tone of voice.
- Some of these behavioural cues are more predominant than others.
- Environmental cues convey important additional information to help recognise or infer emotional states.

What this paper adds

- Not all three behavioural cues need to be present to recognise human emotions in a known context.
- The combination of posture with either facial expression or tone of voice is sufficient to recognise human emotions.
- Both cue combinations lead to equally high recognition rates and qualitatively comparable tactical decisions.

Implications for practice and/or policy

- Designers of tactical decision-making games can lower the quality of the game characters' facial expressions or omit their vocalisations.
- As high-end graphics involve high development costs, the former option seems the most cost-effective.

not necessarily enhance learning (Feinstein & Cannon, 2002; Mania, Wooldridge, Coxon & Robinson, 2006). The present research therefore aspired to establish whether and how the degree of realism (ie., fidelity) in representing human emotions influences emotion recognition.

The modelling and recognition of emotions is key to a wide range of applications of virtual human characters. Examples include avatars in multi-user virtual environments such as Second Life and animated pedagogical agents in e-learning environments. Particularly noteworthy is the work of Baylor, who examined the design and effects of pedagogical agents on motivational and learning-related outcomes (eg., Baylor, 2011; Baylor & Kim, 2009; Kim, Baylor & Shen, 2007). However, even though the modelling of emotions is pivotal to this research, their recognition is assumed rather than assessed. The present research is thus complementary to Baylor's in that it involved a direct assessment of emotion recognition while using a similar interpretation of fidelity and research methodology.

Human beings have a wide range of emotions, which they exhibit mainly through facial expressions, bodily movements (mostly posture) and tone of voice (Argyle, 1988; Baylor, 2011). Whether these behavioural cues are intentionally communicative or not, they often suggest considerable information about a person's emotional arousal (Gratch & Marsella, 2001). Attempts to efficiently model emotional states in virtual human characters are often based on the cue dominance approach. This approach rests on the notion that some cues are more relevant than others, even if they represent the same information. The cues that are most relevant are called dominant cues; weaker cues are neglected if they appear simultaneously with dominant cues representing the same information (Warren & Riccio, 1985).

Cue dominance can be established from research. Various affective computing specialists have pointed out that facial expressions of virtual human characters tend to be somewhat ambiguous

(eg., Donath, 2001). Argyle (1988) conjectured that this ambiguity is reduced when facial expressions are accompanied with congruent bodily cues. This was substantiated by Vinayagamoorthy, Brogni, Steed and Slater (2006), who showed that posture was a more important indicator of a virtual character's emotional state than facial expression but challenged by Clavel, Plessier, Martin, Ach and Morel (2009) who found the opposite to be true. Tone of voice, the third behavioural cue, is omnipresent in multi-user virtual environments and can facilitate emotion recognition from facial cues (Sebe, Cohen, Gevers & Huang, 2006). However, Bailenson, Yee, Merget and Schroeder (2006) found comparable recognition rates in their "voice" and "voice + face" condition.

These inconsistent findings have given insufficient grounds to establish a cue dominance hierarchy. A possible explanation is that the cited studies did not take the impact of environmental cues into account. Carroll and Russell (1996) asserted that contextual information is crucial to identify or infer emotional states. To illustrate, a different emotion is attached to a crying avatar in Second Life when the accompanying chat message explains that the individual whom the avatar represents has either lost his job or got promoted (Noël, Dumoulin & Lindgaard, 2009). It thus seems plausible that emotion recognition is sensitive to the meaning of the scene in which a virtual human character appears. This, in turn, implies that contextual cues should be taken into consideration in determining which behavioural cues are needed to adequately train human emotion recognition in tactical decision-making games.

The present research investigated this issue in two studies. Based on the cue dominance approach, it was assumed that maximum fidelity of all three behavioural cues may not be needed for an observer to recognise the emotional state of a virtual human character. Study 1 therefore sought to establish a cue dominance hierarchy by assessing the relative effectiveness of different combinations of behavioural cues to represent the six emotional states relevant for tactical decision making, which were derived from a task analysis in the fields of infantry tactics and crowd and riot control tactics (Visschedijk, 2010). In keeping with Baylor and Kim (2009), fidelity was defined by the mere presence or absence of behavioural cues, which was deemed more feasible than trying to define and compare distinct fidelity levels for each individual behavioural cue. Study 2 investigated whether this hierarchy would hold if emotion recognition occurs in context.

For the sake of experimental rigour, both studies were performed outside the context of a tactical decision-making game. While this admittedly lowers external validity, it increases internal validity in that research participants can be asked to recognise the exact same emotions in the exact same order. Such controlled conditions are very difficult to achieve in an actual game where the emotions shown by a game character depend on the player's actions. The implications of this choice of research setting are addressed in the general discussion.

Study 1

Method

Participants

Twenty-eight adult volunteers participated in this study. There were 16 males and 12 females with a mean age of 32 years. Participants were neither trained in tactical decision making nor in human emotion recognition.

Materials

The study used 30 computer animations of a virtual human character against a light grey background that were designed with Moviestorm (2005). The character could show an emotion through facial expression, posture or tone of voice. As the character's gender might influence how emotions are expressed or classified by the observer, the study used a single male character;

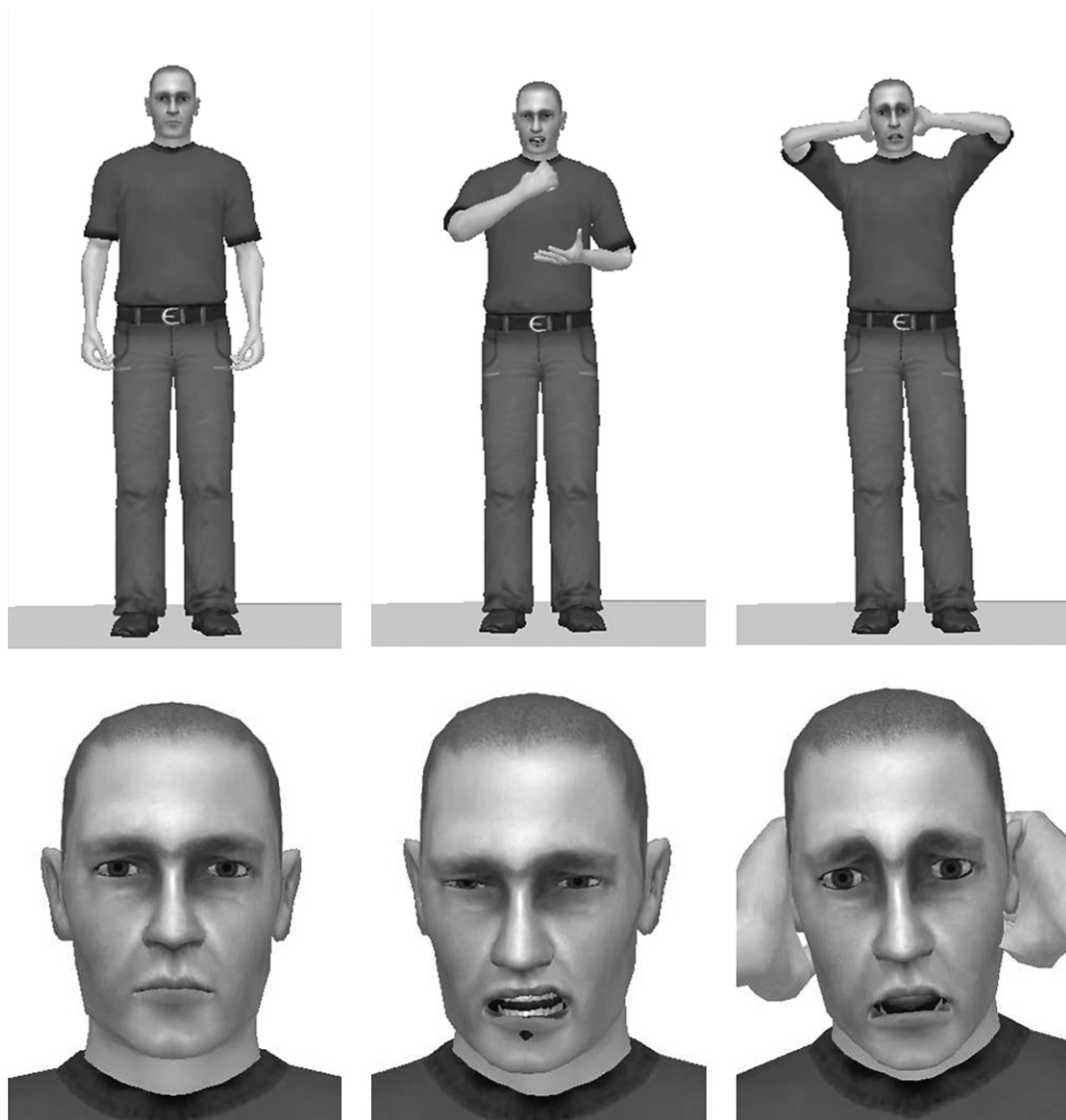


Figure 1: Still images of the virtual human character in a neutral (left), aggression (middle), and panic (right) emotional state. These stills are taken from the 5-second animations that were displayed on a 17-inch monitor (Study1) or via a beamer (study 2)

the choice of this gender was arbitrary. Six emotional states relevant to tactical decision making were modelled: neutral, anger, aggression, fear, panic and elation (see Figure 1 for examples). Facial and postural expressions of each emotion were taken from Moviestorm's library and checked for accuracy against existing standards (Coulson, 2004; Fabri, Moore & Hobbs, 2002; Gunes & Piccardi, 2007; Kleinsmith, De Silva & Bianchi-Berthouze, 2006). Tone of voice was added from Internet-retrieved crowd sound samples. Crowd sounds were preferred because they are prevalent in tactical decision-making situations and do not reveal the information an individual seeks to convey (which would facilitate the recognition of his or her emotional state). Draft versions of the 30 animations were subjected to a pilot test with four individuals who did not participate in the study. Final minor improvements were made based on their comments.

entirely bridged. Assuming that this postulation holds true, it can be tentatively concluded from Figure 2 that most emotional states are best recognised in the P + F + V variant and that the P + F and P + V variants are joint second best.

This provisional conclusion is substantiated by planned contrasts among the variants (this datum was not reported in the results section because of the significant interaction effect but deemed informative for this specific purpose). It was found that the proportion of correct recognition in the P + F + V variant was significantly higher compared with all other variants. The difference between the P + F and P + V variant was not significant, while both variants had significantly higher recognition rates than the P variant. The F + V variant was less effective than the P + F variant (but not the P + V variant) and as effective as the P variant.

Based on these interpretative analyses, it was decided to include the P + F + V, P + F and P + V variant in Study 2. The first goal of this second study was to assess whether contextual cues could reduce cross-variant differences in recognition. A second goal was to find out whether contextual cues would cause all emotional states to be comparably recognisable in all variants.

Study 2

Method

Participants

Twenty-four students from the Dutch military police academy participated in this study. The sample was composed of all males with a mean age of 28 years. These students had some experience in tactical decision making and human emotion recognition and were typical of the target audience for tactical decision-making games.

Materials

Stimulus materials were the 18 animations from the P + F + V, P + F and P + V variants used in study 1. Variants were left intact except that every animation was supplemented with an introductory context description of the setting (the type of people involved, the area where they have gathered and the reason for their gathering), the assignment (the commander's order, such as "control the crowd and prevent escalation") and the penultimate event (e.g., "some crowd members have just been arrested") (cf. Klein, 2008). The former two aspects are part of a standard briefing; the latter was added because emotions are often triggered by some preceding event that, unlike in tactical decision-making games, would otherwise remain covert. Six contextual descriptions were used in total, one for every emotion. The context descriptions were checked by four subject matter experts; minor adjustments were made based on their comments.

An answer form containing two open-ended questions was designed to gauge participants' interpretation of each animation. The first question asked the participants to describe which emotion the animation represented, the second question inquired after their commanding decisions in case the whole crowd would be in this emotional state.

Procedure

The study was conducted in one session that took place in a regular classroom. Instructions were similar to those of study 1 except that the participants were directed to fill out the answer form after each animation. The two questions were clarified and the participants were instructed to answer them as if they were a platoon commander. After the instruction, the experimenter read aloud the first context description and showed the associated animation twice via a beamer. The participants then had 2 minutes maximum to complete the answer form. This time limit was imposed for practical reasons (the experiment had to be performed within one lesson); experiences gained in study 1 proved that 2 minutes would be more than enough time to answer each question. The remaining 17 animations were administered similarly. Variants appeared in ascending order based on the recognition scores from study 1.

the prediction that a subset of these behavioural cues is sufficient to properly recognise an emotion in a known context. This conclusion has direct implications for the design of tactical decision-making games, and provides concrete directions for future research.

Results from study 1 confirm the notion that emotion recognition tends to improve with multiple behavioural cues (cf. Argyle, 1988; Clavel *et al*, 2009; Crane, 2009; Vinayagamoorthy *et al*, 2006). Recognition rates were highest in the P + F + V variant and lowest for the P variant; the differences between the three bimodal variants was less pronounced and depended on the type of emotion. Study 2 yielded superior recognition compared with study 1, but this result can not be attributed solely to the contextual cues because the military police academy students were more experienced in emotion recognition than the adult volunteers in study 1. This indistinctness is of minor importance here because the present research focused on tactical decision-making games that, by definition, contain contextual cues and are intended for trainees with some prior knowledge. It would nevertheless be interesting for future research to investigate the effects of contextual cues and subject knowledge in isolation. Study 2 further showed comparable recognition rates in the P + F + V and P + F variants. The P + V variant was as effective as the P + F variant but less effective than the P + F + V variant. These findings suggest that postural and facial expressions are dominant cues for recognising human emotions in known circumstances, whereas tone of voice is a weak cue with little added value.

However, suboptimal recognition in the P + V variant was mainly due to the emotion “fear”; the other emotions were as well recognised as in the P + F + V variant. This result implies that facial cues are dominant in the recognition of “fear” only. A possible explanation is that “fear,” which is difficult to recognise anyway (Argyle, 1988; see also study 1), is mainly expressed through the face and associated with rather subtle body movements (slight crouching) and ambiguous sounds (soft, shivery breathing). The other emotions, by contrast, rely less heavily on facial expressions and have more explicit bodily movements and sounds that facilitate recognition in absence of facial cues.

More importantly, the quality of the participants’ tactical decisions was comparable among variants. Even though the overall scores were somewhat low (which is quite understandable given that participants were tactical decision-making trainees), the least well-recognised emotions “fear” and “panic” showed substantial internal consistency. This result suggests that the relatively poor recognition in the P + V variant did not impinge on the eventual tactical decision. And because making accurate decisions is the ultimate goal of tactical decision-making training, it seems fair to consider the P + V variant as viable low-fidelity alternative to represent most human emotions.

Taken together, this research provides evidence that posture with either facial expression or tone of voice are dominant cue combinations that suffice to inform tactical decision making if trainees are briefed on the context and background of the target situation. Future research should reveal whether facial expression and tone of voice qualify as dominant combination as well because the F + V variant was not included in Study 2. This conclusion is generally consistent with the cue dominance approach (Warren & Riccio, 1985) and demonstrates that the results found in fundamental emotion recognition experiments extend to more realistic settings. The latter claim is supported by the fact that the present research addressed all emotions relevant to tactical decision making and assessed their recognition through a free response format (which is more valid than selecting emotions from a predefined list). Additionally, as this research went beyond the mere recognition of emotions by considering participants’ tactical decisions, its results generalise to situations in which emotion recognition is a means rather than an end.

However, generalisability might be challenged by the use of a male character to express the emotions. Would the same results be obtained if a female character had been used? It seems

plausible that men and women express the same emotions differently, or with different intensity, and this might influence recognition. Likewise, similar expressions may be variously classified depending on the gender of the virtual character. A screaming woman, for instance, is presumably more often associated with the emotion “panic” than a screaming man. A related issue concerns the virtual character’s clothing. While uniformly dressed in the present research, different outfits and headwear could either facilitate or complicate emotion recognition. In addition, demographic and cultural characteristics such as race, nationality and foreign accents might influence emotion recognition as well. Future research should therefore establish how a virtual character’s appearance influences emotion recognition.

Future studies could also strengthen the validity of the present findings. One suggestion would be to compare the results from study 2 to emotion recognition and commanding decisions in an actual tactical decision-making game. The design of this research, although seemingly straightforward, will need some careful consideration because developing three versions of the same game is a rather expensive endeavour. Another interesting avenue for further research would be to replicate study 2 with a sample of decision-making experts. Their decisions would probably be more appropriate than those from the trainees in study 2. Whether their decisions will also be more consistent across variants remains to be shown.

Practical implications pertain to the design of digital tactical decision-making games. Due to their specific contents, these games are usually custom-made and not freely available. Their development costs are not made public either, but estimates are that their commercial equivalents (war games, strategy games) require over 20 million euros to develop (Prensky, 2008; Takatsuki, 2007). Serious game designers looking for ways to cut development costs while maintaining training effectiveness could either lower the quality of the game characters’ facial expressions or omit their vocalisations. As high-end graphics typically involve high development costs, the former option seems the most cost-effective. Both recommendations might generalise to other serious game genres that hinge on human emotion recognition (eg., role playing games, medical simulation games) and the design of virtual coaches and pedagogical agents.

Note

1. As the sphericity assumption was violated, Greenhouse–Geisser corrected degrees of freedom are reported.

References

- Argyle, M. (1988). *Bodily communication*. New York: Methuen.
- Bailenson, J. N., Yee, N., Merget, D. & Schroeder, R. (2006). The effect of behavioral realism and form realism of real-time avatar faces on verbal disclosure, nonverbal disclosure, emotion recognition, and copresence in dyadic interaction. *Presence*, 15, 359–372.
- Baylor, A. L. (2011). The design of motivational agents and avatars. *Educational Technology Research and Development*, 59, 291–300.
- Baylor, A. L. & Kim, S. (2009). Designing nonverbal communication for pedagogical agents: when less is more. *Computers in Human Behavior*, 25, 450–457.
- Carroll, J. M. & Russell, J. A. (1996). Do facial expressions signal specific emotions? Judging emotion from the face in context. *Journal of Personality and Social Psychology*, 70, 205–218.
- Clavel, C., Plessier, J., Martin, J. C., Ach, L. & Morel, B. (2009). Combining facial and postural expressions of emotions in a virtual character. In Z. Ruttkay, M. Kipp, A. Nijholt & H. H. Vilhjálmsón (Eds), *Intelligent virtual agents* (pp. 287–300). Berlin: Springer.
- Coulson, M. (2004). Attributing emotion to static body postures: recognition accuracy, confusions, and viewpoint dependence. *Journal of Nonverbal Behavior*, 28, 117–139.
- Crane, E. A. (2009). Measures of emotion: how feelings are expressed in the body and face during walking. (Unpublished doctoral dissertation, The University of Michigan).

