



# DIGITAL ACCESS TO SCHOLARSHIP AT HARVARD

## The Influence of Emotion Expression on Perceptions of Trustworthiness in Negotiation

The Harvard community has made this article openly available.  
[Please share](#) how this access benefits you. Your story matters.

<b>Citation</b>	Antos, Dimitrios, Celso De Melo, Jonathan Gratch, and Barbara Grosz. 2011. The influence of emotion expression on perceptions of trustworthiness in negotiation. In Proceedings of the Twenty-Fifth AAI Conference on Artificial Intelligence: August 7-11, 2011, San Francisco, California, ed. American Association for Artificial Intelligence, 772-778. Menlo Park, California: AAI Press.
<b>Published Version</b>	<a href="http://www.aaai.org/ocs/index.php/AAAI/AAAI11/paper/view/3438">http://www.aaai.org/ocs/index.php/AAAI/AAAI11/paper/view/3438</a>
<b>Accessed</b>	February 19, 2015 9:06:31 AM EST
<b>Citable Link</b>	<a href="http://nrs.harvard.edu/urn-3:HUL.InstRepos:5344528">http://nrs.harvard.edu/urn-3:HUL.InstRepos:5344528</a>
<b>Terms of Use</b>	This article was downloaded from Harvard University's DASH repository, and is made available under the terms and conditions applicable to Open Access Policy Articles, as set forth at <a href="http://nrs.harvard.edu/urn-3:HUL.InstRepos:dash.current.terms-of-use#OAP">http://nrs.harvard.edu/urn-3:HUL.InstRepos:dash.current.terms-of-use#OAP</a>

*(Article begins on next page)*

# The influence of emotion expression on perceptions of trustworthiness in negotiation

Dimitrios Antos<sup>1</sup>, Celso De Melo<sup>2</sup>, Jonathan Gratch<sup>2</sup> and Barbara Grosz<sup>1</sup>

<sup>1</sup>Harvard University, Cambridge, MA 02138, USA

<sup>2</sup>Institute for Creative Technologies, University of Southern California, Los Angeles, CA 90094, USA  
{dantos, grosz}@eecs.harvard.edu, demelo@usc.edu, gratch@ict.usc.edu

## Abstract

When interacting with computer agents, people make inferences about various characteristics of these agents, such as their reliability and trustworthiness. These perceptions are significant, as they influence people's behavior towards the agents, and may foster or inhibit repeated interactions between them. In this paper we investigate whether computer agents can use the expression of emotion to influence human perceptions of trustworthiness. In particular, we study human-computer interactions within the context of a negotiation game, in which players make alternating offers to decide on how to divide a set of resources. A series of negotiation games between a human and several agents is then followed by a "trust game." In this game people have to choose one among several agents to interact with, as well as how much of their resources they will trust to it. Our results indicate that, among those agents that displayed emotion, those whose expression was in accord with their actions (strategy) during the negotiation game were generally preferred as partners in the trust game over those whose emotion expressions and actions did not mesh. Moreover, we observed that when emotion does not carry useful new information, it fails to strongly influence human decision-making behavior in a negotiation setting.

## Introduction

It has been well-documented that humans treat computers as social agents, in that they perceive in them human traits, expect socially intelligent responses, and act toward them in a socially appropriate manner (Nass and Moon 2000; Nass 2004). This leaves open the possibility that an agent may influence how humans perceive it through its presence and behavior. The agent might find it advantageous to do so if such perceptions stand to influence the interaction outcome, or perhaps the likelihood that similar interactions will take place in the future. In this paper we focus on human perceptions of *trustworthiness*, and study the extent that these may be influenced by agent expressions of *emotion* in a computer-human negotiation environment.

Negotiation is a commonly-used method for parties with diverging interests to reach a mutually-beneficial agreement.

Copyright © 2011, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

People negotiate over how to schedule activities when participants have different time constraints and priorities, to efficiently allocate valuable resources across individuals or corporations with varying needs and preferences, or to resolve international conflicts without resorting to violence. The significance of negotiation has led to a large literature in the fields of psychology, economics, sociology and computer science (e.g., see (Lewicki, Barry, and Saunders 2010), (Raiffa 1985) and (Jennings et al. 2001)). Computer agents have been used in many negotiation settings, sometimes acting on behalf of humans (Jennings et al. 2001) and sometimes negotiating with them (Lin and Kraus 2010). As the scale and complexity of the domains in which negotiation is employed are expected to increase, we foresee a growth in the use of computer agents as negotiators; examples of such domains might include traffic management (Kamar and Horvitz 2009) and commerce (Maes, Guttman, and Moukas 1999), among others. Furthermore, computer agents have shown potential (compared with human negotiations) for improving negotiation outcomes in some cases (Lin, Oshrat, and Kraus 2009).

Yet negotiation studies with computer agents have largely overlooked the fact that humans use significant verbal and non-verbal cues when they negotiate (Drolet and Morris 2000). The expression of emotion, in particular, has been shown to significantly influence negotiation outcomes (Barry, Fulmer, and Goates 2006; Van Kleef, De Dreu, and Manstead 2010). For instance, displaying anger was shown to be effective in forcing larger concessions out of the other party, whereas positive emotion was found to be helpful in exploring and achieving mutually beneficial (integrative) solutions (Van Kleef, De Dreu, and Manstead 2004; Carnevale and Pruitt 1992). The effects of emotion during business negotiations were also shown to be modulated by culture (Leung et al. 2005). In most studies in which emotion was expressed by computer agents, this emotion was conveyed by means of text sent by the computer agent to its human partner (e.g., in (Van Kleef, De Dreu, and Manstead 2004) the computer would say "this offer makes me really angry"). More recent implementations of virtual agents tested for differences among different modalities of emotion expression (e.g., text or virtual face) and found no significant differences (de Melo, Carnevale, and Gratch 2010).

We make two contributions in this paper. First, we inves-

tigate the effects of emotion on people’s perceptions of an agent’s *trustworthiness*. Such perceptions are formed in the course of human-computer interactions (in our case, negotiations; see (Bos et al. 2002)). We obtain robust results indicating that the expression of emotion does influence perceptions of trustworthiness. In particular, an agent is being perceived as more trustworthy when its expressed emotion is in accord with its actions during the interaction.

Second, we look at the effect of expressed emotion on negotiation outcomes. Our results reveal only small influences on most negotiation metrics, in contrast with many well-documented findings (e.g., anger resulting in higher concessions by the human). We hypothesize that this is a consequence of the fact that the negotiation strategies of our agents were too “scripted” and predictable. As a result, the expressed emotion did not convey new information that could be used by humans in their decision-making. Without the emotion carrying an important informational signal, people seem to ignore it.

These observations show that the expression of emotion holds potential benefits in human-computer decision-making, and should be viewed as a key part of agent design. However, deploying emotion successfully is not a simple matter, because the appropriate expression is very much dependent on context. Moreover, our results suggest that emotion can influence human decision-making behavior only when it has the ability to convey new information beyond the observed decisions of the agent.

## Experiment Design

To investigate the effect of emotion in human-computer negotiation and perceptions of trustworthiness we developed the following game: A human subject ( $h$ ) is paired with a computer agent ( $a$ ), and they must negotiate on how to divide a set of resources amongst themselves. The resources consist of virtual “coins” of four types: gold, silver, bronze and iron. In each game, there are three (3) coins of each type. Before the game starts, people are told the *relative* value of each coin type, i.e., that gold coins are more valuable than silver, which are more valuable than bronze, which are more valuable than iron coins. However, people are not given the exact numeric value (in points) of each coin type.<sup>1</sup> Subjects are also informed that the relative valuation of the items by the agents might be different than their own, e.g., computers might prefer silver coins over gold ones. Notationally, we refer to the four item types with numbers  $j \in \{1, 2, 3, 4\}$ . We also use  $w_j^i$  to denote the number of points player  $i \in \{h, a\}$  receives by possessing a coin of type  $j$ . In all experiments  $\mathbf{w}^h = \langle 10, 7, 5, 2 \rangle$  and  $\mathbf{w}^a = \langle 10, 2, 5, 7 \rangle$ . Notice how item types 1 and 3 (gold and bronze) present a “distributive problem,” i.e., players need to decide how to split items of common value, but items 2 and 4 (silver and iron) present “integrative potential,” i.e., there are exchanges of items that lead to mutual benefit. Moreover, it must be pointed out that computer agents have full knowledge of vectors  $\mathbf{w}^h$  and  $\mathbf{w}^a$ .

<sup>1</sup>A preliminary experiment with  $N = 25$  showed that people’s behavior is not affected by them knowing the point value of every coin type.

The game proceeds by means of alternating offers, and participants play in turns, with the human always making the first offer in a game. An offer by player  $i \in \{h, a\}$  consists of a complete allocation of all coins between the two participants. We use the notation  $c_j^i(t)$  to denote how many items of type  $j \in \{1, 2, 3, 4\}$  player  $i$  was given in the offer made at round  $t \in \{1, 2, \dots\}$ . Hence, allowing only for complete allocations means that in every round  $t$ , offers must satisfy  $c_j^h(t) + c_j^a(t) = 3, \forall j$ . In every round  $t > 1$  the negotiator whose turn it is may *accept* an offer, in which case the game ends and both participants make their corresponding points; for the human player, these would be  $\pi^h = \sum_j w_j^h c_j^h(t - 1)$ . Alternatively, she may *reject* the offer, and counter-offer a different split of the items. Finally, at any point in the game, either participant may *drop out*. A game consists of a maximum of 15 rounds.<sup>2</sup> If no agreement is reached in any of these rounds, or if either player drops out, both players make zero points.

The agents in our experiment differed in two ways: with respect to their *strategy*, and with respect to their *emotion expression*. The strategy of an agent encompasses when offers are accepted or rejected, and what counter-offers are made by it. Likewise, the emotion expression of an agent defines whether emotion is expressed, and what type of emotion is displayed in each circumstance. Below we discuss the strategies and emotion expression policies we used in our agents; we also explain how emotion was displayed to the human.

## Strategies of computer agents

The strategy of an agent prescribes how the agent behaves as a negotiator. Although the literature on effective negotiation strategies is extensive (Sycara and Dai 2010), we limited ourselves to simple, intuitive strategies for this experiment. Our goal was not to exhaustively explore the effect of emotion expression given complex strategies, but to assess whether emotion has any effect on people’s behavior in computer-human negotiation, and whether this effect is dependent upon the agent’s strategy.

The strategies we used varied along two dimensions: “flexibility” and “self-interestedness.” An agent follows a *flexible* strategy if its offers change from round to round throughout the game; its strategy is *inflexible* if it always makes the same offer (or very similar ones) in every round. In a similar fashion, an agent is said to follow a *self-interested* strategy if it attempts to keep for itself almost all the points being negotiated; its strategy is *non-self-interested* if the agent seeks agreement on balanced offers, which provide a more or less equal split of the points. We used four simple strategies, described below. (Table 1 groups them according to flexibility and self-interestedness.)

1. *Selfish*: The selfish agent in every round chooses a single coin at random (but never a gold one) to counter-offer to the human, and keeps everything else for itself. Thus selfish agents are inflexible and self-interested.

<sup>2</sup>Notice how, if the human always makes the first offer, the agent always makes the last offer. If the game reaches the 15th round, then the human can either *accept* the computer’s offer, or *drop out*, since there are no more rounds for counter-offers to be made.

2. *Nash*: This agent computes the Nash bargaining point (N.B.P.) of the game, which is the allocation that maximizes the product of both players' payoffs, and offers that in every round. The N.B.P. presents the theoretically most efficient point in the negotiation, as it is Pareto-optimal and satisfies a series of axiomatic constraints (Nash 1950). N.B.P. allocations split the points in a very balanced fashion, thus this agent is inflexible but non-self-interested.
3. *Conceder*: This agent performs concessions in a constant rate. In particular, no matter how the human behaves, at round  $t$  the agent offers her  $\frac{3t}{2}$  points and keeps everything else for itself. In other words, the first time it plays it will offer the human 3 points (round 2), the second time 6 points, etc. Since this agent starts from very imbalanced offers (only 3 or 6 of a total of 72 points) and concedes slowly, it is categorized as self-interested but flexible.
4. *Tit-For-Tat*: This is an agent implementing reciprocity. In round  $t$  it offers the human  $0.8 \times \sum_j w_j^a(t-1)$  points. Hence, the more concessionary the human has been in her last offer, the more concessionary the agent becomes; likewise, if the human has been selfish, the agent would reciprocate this. The 0.8 coefficient represents a degree of "toughness" by the agent, i.e., it reciprocates slightly less than what it is being offered. This agent is both flexible and non-self-interested, as agreement will only be reached when the two players start conceding, eventually "meeting" somewhere in the middle.

	Inflexible	Flexible
Non-self-interested	Nash	Tit-For-Tat
Self-interested	Selfish	Conceder

Table 1: Strategies used in the negotiation game, grouped according to their flexibility and self-interestedness.

All agent types accept an offer made by the human if and only if the points they would request in their counter-offer (according to their strategy) are no greater than the points the human is currently giving them. Agents never drop out of the game. Also, whenever agents wish to give a certain number of points to the human, they choose the most integrative way of doing so (i.e., of all possible counter-offers that would give the human  $c$  points, they choose the offer that maximizes their own points).

### Emotion expression by agents

The emotion expression policy of an agent denotes whether and how it displays affect. Affect in our game was displayed on a "face" the agent was given. Faces were all male, and were randomly assigned to the various agents from a pool of 15 templates, such that no subject would interact with two agents bearing the same face during the experiment. The face of the agent was rendered to the side of the game board, on which the items were being negotiated. We used five emotion expression policies, described below:

1. *No-Face*: This is the baseline case, in which there is no visible face to the agent.

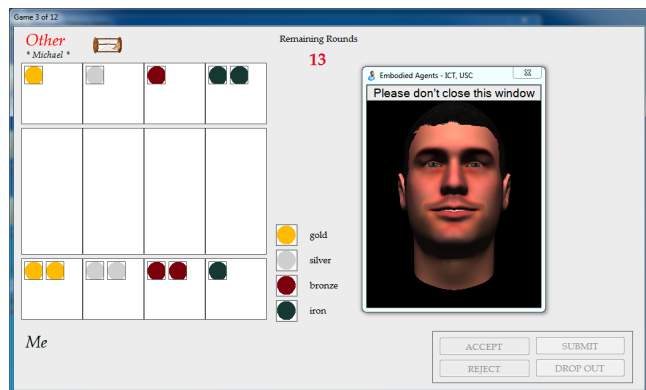


Figure 1: The Negotiation Game

2. *Poker-Face*: This agent shows a face, but never displays any emotion on it, always keeping a neutral expression. We differentiate between this agent and the No-Face to assess how much of any effect comes from displaying emotions, or merely from the presence of a face (even if it displays no emotions).
3. *Always-Smile*: This agent displays a face and smiles to all the offers made by the human, independently of what these offers look like.
4. *Always-Angry*: This agent displays a face and expresses anger toward all the offers made by the human, again, independently of what these offers look like.
5. *Appraisal*: This agent would smile or show anger depending on the human's actions, instead of following a fixed strategy. If at round  $t$  it was offered by the human at least  $\frac{3}{2}t$  points it would smile, otherwise it would show anger.

All agents that display emotion follow the same pattern of expression: First, they "look" to their right, where the coins are, to "see" the offer made by the human; they then "look back" toward the human (straight gaze), perform their expression (of joy or anger), and send their counter-offer (or acceptance notification). Joy is expressed by a smile across the face, which forms in moderate speed (1 sec). Anger is expressed by contraction of the corrugator muscle (frowning) as well as an aggressive twitching of the mouth. Expressions dissipate linearly towards normal (expressionless face) after the counter-offer is sent, while the human is deciding her move. These particular displays were shown to be successful in conveying the desired emotions in (de Melo, Carnevale, and Gratch 2010). Also, no "gradients" of the expressions were employed (e.g., more or less intense smiles)—all expressions were of the same intensity. A screenshot of the negotiation game can be seen in Figure 1.

It must be pointed out that several factors in the presentation of emotion have been overlooked in order to keep the experiment simple, although they could presumably be carrying strong effects which are well-documented in the literature. In particular, we did not test for gender effects, as all our agents were male. We also did not test for the effect of

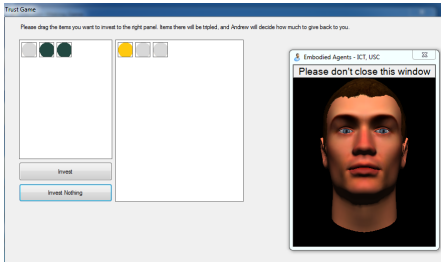


Figure 2: The Trust Game

race, age, or physical attractiveness, which could all mediate how expressed emotions are interpreted. In all these cases, we tried to keep these variables constant (using agents of the same age and gender) whenever possible, or randomize uniformly (for race).

### The trust game

Each subject in the experiment plays several negotiation games with different agents. After every three such games, in which the agents differ in their emotion expression policy, but not in their strategy, subjects are asked to play a ‘trust game.’ As an example, someone might play with three Tit-For-Tat agents, the first of whom always smiles, the second always shows anger, and the third maintains an expressionless face. After these three games, the subject is asked to play a trust game with one of them. The trust game is a variant of the popular public goods game. To play it, people first have to select which of these three agents they want to be paired with. (Agents are given names at random, like “Michael” or “James” to assist recall.) After they indicate their choice, they are presented with a *trust* problem. In particular, they are shown their total gains (in coins) from the previous three negotiation games, and are asked what fraction of these they are willing to trust to the agent they have chosen as their partner. If the agent’s policy includes showing a face, it also is displayed, but no emotions are ever shown on it.

If the subject chooses to trust her partner with a non-empty subset of her gains  $\mathbf{t} = \langle t_1, t_2, t_3, t_4 \rangle$ , where  $t_j$  is the count of coins of type  $j$  she has trusted, this subset is tripled at the hands of the agent. Then, however, the agent chooses what percentage  $p \in [0, 1]$  of the value of these coins it will return to the human, and how much it will keep for itself. Agents all choose  $p$  uniformly at random from  $[\frac{1}{3}, \frac{2}{3}]$ , but subjects are told nothing about this. The subject at the end of the trust game keeps the points she has not trusted to the agent, to which  $p \times \sum_j (3t_j)w_j^h$  points is added. The trust game looks like the screenshot in Fig. 2.

### Experiment process

Each subject in the experiment played twelve negotiation games, with each triad followed by a trust game. The order of games was randomized, and each subject faced triads of agents that differed in terms of strategy of emotion expression policy, but not both. Instructions were delivered to

the subjects over video, and they were all debriefed after the end of the experiment. After each negotiation game, subjects were asked to answer, in Likert 1-7 scales, four short questions regarding their experience with the agent they had just played with. Subjects were paid \$20 for their participation, which lasted about 45 minutes. They were also told that the person who would score the highest number of points would be awarded an extra \$100. We had  $N = 88$  subjects, for a total of 1,056 negotiation games and 352 trust games.

## Hypotheses

Our experiments were designed to test two hypotheses. The first concerns the influence of emotion on perceptions of trustworthiness, while the second relates to the effect of emotion on negotiation outcomes. We formulate our main hypothesis (H1) using the notion of *action-consistency*. We call an emotion expression *action-consistent* with a strategy if the emotion emphasizes the characteristics of the agent that manifest in its actions. Positive emotion typically emphasizes kindness and goodwill, whereas negative emotion is usually illustrative of selfishness and intransigence. Hence, positive emotion is more *action-consistent* with non-self-interested strategies, and negative emotion is more *action-consistent* with self-interested strategies. In the same spirit, positive emotion is more *action-consistent* with flexible than with inflexible strategies. Alternative notions of consistency, which we do not discuss here, have been introduced in the literature before (e.g., (Nass et al. 2005)).

**H1.** *People’s perceptions of an agent’s trustworthiness are influenced by the action-consistency between the agent’s emotional expression and its strategy.*

Agents whose expressions are consistent with their strategies will be preferred as partners for the trust game. Thus, when faced with a choice among self-interested or inflexible agents, people will tend to prefer angry ones; and when faced with a choice among non-self-interested or flexible agents, they will tend to prefer smiling ones.

**H2.** Within the negotiation, we expect that:

- (a) *The expression of anger will result in higher concession rates by humans* (Van Kleef, De Dreu, and Manstead 2004).
- (b) *Agents who smile will help foster more integrative outcomes* (Carnevale and Pruitt 1992). Recall that integrative outcomes are those in which the sum of the two players’ payoffs is high. These can be achieved if the players realize there are mutually beneficial exchanges (like one silver coin for one iron coin) that increase both their payoffs.
- (c) *Positive emotion will cause humans to also concede more points.* The theory of “emotion contagion” (Hatfield, Cacioppo, and Rapson 1992) predicts that humans in those situations will be more cooperative.

## Results

To assess people’s behavior in the trust game, we used as a metric their choice of partner in the trust game. As was mentioned before, a trust game always came after three negotiation games, in which the agents differed in their emotion expression policy but not in their strategy. Therefore people were presented with a choice among three candidate emotion display policies (keeping strategy fixed). To investigate whether smiling agents or angry agents were preferred, we tested for  $\Delta\mu = \mu_{smile} - \mu_{angry} = 0$  (in which  $\mu_x$  denotes the fraction expression  $x$  was chosen), which would be the case if smiling and angry agents were equally preferred. We observe that, among selfish agents,  $\Delta\mu = -0.25, p < 0.05$ , denoting a preference toward angry agents. Among conceiver agents, we similarly have  $\Delta\mu = -0.125$ , although this is not statistically significant. A similar trend was seen with Nash agents ( $\Delta\mu = -0.153$ ). On the other hand, among tit-for-tat agents we observe preference toward smiling agents ( $\Delta\mu = +0.25, p < 0.05$ ). Notice that angry agents are being more preferred over smiling ones the more the strategy of the agent becomes inflexible or self-interested, confirming hypothesis H1 (results are summarized in Table 2). We found no effect of emotion on the amount of resource trusted.

	Inflexible	Flexible
<b>Non-self-interested</b>	-0.153	+0.25*
<b>Self-interested</b>	-0.25*	-0.125

Table 2:  $\Delta\mu = \mu_{smile} - \mu_{angry}$ , where  $\mu_x$  denotes the fraction of times emotional expression  $x$  was preferred, for the various strategies in the trust game. Asterisk denotes that the mean is significantly different from zero.

To assess the influence of emotion expression on negotiation outcomes, we used both behavioral and subjective measures. Behavioral measures include variables relating to the negotiation game, such as the players’ payoff at the end of the game, drop-out rates (i.e., the percentage of games with an agent in which humans dropped out, ending the game without a deal and awarding zero payoff to both participants), and measures indicating whether the games’ integrative potential was exploited. Subjective measures, on the other hand, come from people’s questionnaire responses after each negotiation game. Table 3 lists all the measures used. We now turn to our hypotheses.

Surprisingly, we did not observe anger having any effect of the average human concession rate ( $\kappa$ ), thus disconfirming hypothesis H2(a). Figure 3 plots the average concession rate across emotion expressions for all four strategies examined. As can be seen, the average concession rate of the people was strongly influenced by the strategy of the agent, but very little by its emotion expression. Similarly, we found no support that the drop-out rate ( $d$ ), or the integrative potential, as measured by  $\pi_\Sigma$  or  $\pi_\Pi$ , was influenced by the choice of emotion across all strategies (thus H2(b) is also rejected). Finally, our hypothesis that positive emotion will influence

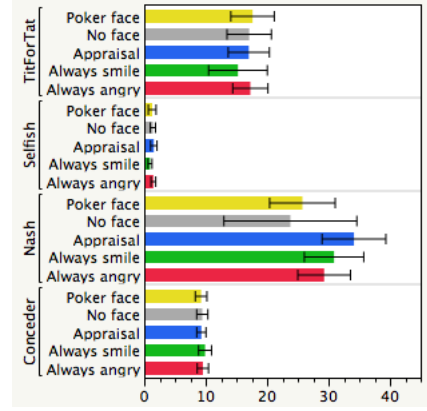


Figure 3: Average concession rate ( $\kappa$ ) per emotion expression and strategy in the negotiation game (error bars denote 95% confidence intervals).

concession rates according to the theory of “emotion contagion,” according to H2(c), was also not supported by our measurements.

It must be noted that action-consistency seems to also play a role in people’s subjective reports. With respect to “liking” ( $q_1$ ) people showed an aversion towards always-angry agents for flexible strategies, but not for inflexible ones. With respect to how much people felt the agent cared about them ( $q_2$ ), always-smile agents were preferred under non-self-interested strategies. Also, with respect to people’s expressed desire to play with an agent again in the future ( $q_4$ ), we saw an aversion toward always-angry agents only among agents playing the tit-for-tat strategy. All the above results indicate that action-consistent expressions are preferred over others. Finally, looking at how much the agent was perceived to be human-like ( $q_3$ ) we noticed no effects of emotion.

## Supplementary results

This section presents further findings of the experiment that were not directly related to the hypotheses examined. We report these for two reasons: (a) because they might be useful

Behavioral Measures	
$\pi^a$	: agent points at the end of the game
$\pi^h$	: human points at the end of the game
$\pi_\Sigma = \pi^a + \pi^h$	: sum of payoffs
$\pi_\Pi = \pi^a \pi^h$	: product of payoffs
$\kappa$	: average human concession rate between two rounds
$d$	: average drop-out rate (%)
Subjective Measures (Likert 1-7)	
$q_1$	: how much did you like this agent?
$q_2$	: how much did you feel this agent cared about you?
$q_3$	: how human-like was the agent?
$q_4$	: would you play with this agent again?

Table 3: Behavioral and subjective measures.

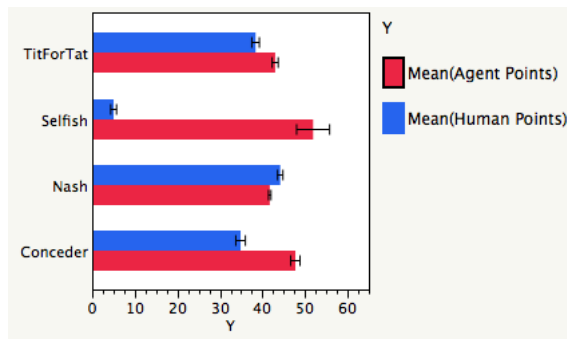


Figure 4: Average human and agent points in the end of the game per strategy (error bars denote 95% confidence intervals).

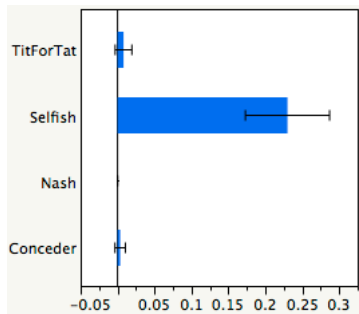


Figure 5: Average human drop-out rate per strategy (error bars denote 95% confidence intervals).

in comparing with other negotiation experiments in the literature, and (b) as evidence that the various strategies our agents employed did make a difference in people’s behavior (in other words, the strategy choice was not superfluous). Hence Figure 4 displays the points in the possession of the agent and the human at the end of the game. Here it can be seen that the conceder and selfish agents (self-interested) fare better in terms of final score than the non-self-interested agents. Also, the Nash agent causes very equitable outcomes to be obtained.

Figure 5 show the effect of the selfish agents on humans who choose to drop-out, hence punishing the agent at a cost (they both receive zero points). As can be seen in the chart, no player drops out against any other agent, but up to a quarter of the participants do when their partner shows such intransigence and self-interestedness.

Finally, Figure 6 shows the average duration of the game in rounds. It is clear that people agree to a deal sooner with non-self-interested agents, but try harder against self-interested ones. (Interestingly enough, we observed a small but statistically significant difference of emotion expression on the game’s duration in the case of the selfish agent. In that subset of the data, it can be seen that smiling causes people to play on average for two more rounds with the selfish agent, trying to forge a deal, before they accept its terms or drop out. We believe that this is because smiling conveys an interest—on the agent’s behalf—to be cooperative and work

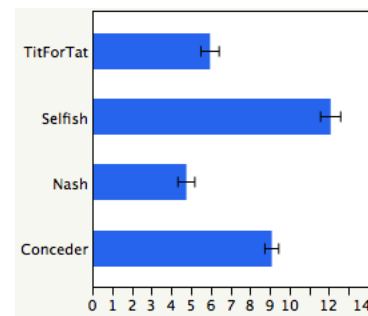


Figure 6: Average duration of game in rounds per strategy (error bars denote 95% confidence intervals).

together with the human, which keeps her trying for longer before conceding to the fact that the agent will not change its mind.)

## Summary & Discussion

Our experiment has illustrated that agents can use the expression of emotion to influence how trustworthy humans perceive them to be. These human perceptions were reflected by people’s choice of partner in a “trust game,” for which they had to select one among three agents. These agents had previously interacted with them in the context of a negotiation and differed in their emotional expressions, but not in their negotiation strategies. Second, we observed that, given a set of very “scripted” and predictable negotiation strategies, emotion seems to be ignored and well-documented effects of it disappear. In interpreting this finding, we hypothesize that when the emotion does not carry useful new information, it fails to influence human decision-making in a strong and consistent manner.

These results imply that emotion expression can confer advantages to an agent and can be seen as an essential part of agent design, perhaps just as important as the strategy of the agent, if the goal is to ensure “returning customers.” As our experiment has indicated, agents with the appropriate emotion expressions can be perceived as more trustworthy and be selected more often for repeated interactions, even if the algorithms used for their decision-making do not differ from other agents’. Moreover, this effect of emotion is not necessarily mediated by the outcome of the interaction between a person and an agent. In other words, it is not necessarily the fact that emotion expressions change the outcome of an interaction, which then in turn influences human perceptions. Even when the outcome of the negotiation in our experiment remained the same across emotion expression conditions, people’s perceptions of the agents’ trustworthiness in practice differed significantly.

Furthermore, our findings indicate that emotion is a signal people attend to in their decision-making by considering an contextual interpretation of it. Hence, it is not merely the “presence” of *some* emotion (e.g., statically positive emotion) that matters, but what this emotion is capable of telling people about the situation. Therefore, in cases where the emotion expressed is not easily or unambiguously at-

tributable to some (real or imagined) trait or intention of the agent, its presence might fail to introduce any effect. This complicates the design of emotion expressions for artificial agents, but also establishes a link between emotions and the tasks and roles the agents are expected to perform.

In the future we plan to extend these investigations. In particular, we wish to explore a wider set of negotiation strategies, which will be more “natural” and differ along more dimensions than just self-interestedness and flexibility. We also plan to entertain more varied emotional expressions, in terms of the emotion types the agent is capable of displaying, as well as the appraisal functions used to compute what emotion should be expressed in every circumstance. Finally, we are interested in exploring how emotional signals impact strategic reasoning, and whether game-theoretic models of decision-making can be developed that receive affective signals to update belief distributions and compute best responses given such non-verbal cues.

### Acknowledgements

This work has been partly supported by NSF grant IIS-0705406.

### References

- Barry, B.; Fulmer, I.; and Goates, N. 2006. Bargaining with feeling: Emotionality in and around negotiation. *Negotiation Theory and Research, New York Psychology Press* 99–127.
- Bos, N.; Olson, J.; Gergle, D.; Olson, G.; and Wright, Z. 2002. Effects of four computer-mediated communications channels on trust development. In *Proceedings of the SIGCHI conference on Human factors in computing systems: Changing our world, changing ourselves, CHI '02*, 135–140. New York, NY, USA: ACM.
- Carnevale, P., and Pruitt, D. 1992. Negotiation and mediation. *Annual Review of Psychology* (43):531–582.
- de Melo, C.; Carnevale, P.; and Gratch, J. 2010. The influence of emotions in embodied agents on human decision-making. In *Intelligent Virtual Agents*.
- Drolet, A., and Morris, M. 2000. Rapport in conflict resolution: Accounting for how face-to-face contact fosters cooperation in mixed-motive conflicts. *Journal of Experimental Social Psychology* (36):26–50.
- Hatfield, R.; Cacioppo, J.; and Rapson, R. 1992. *Emotion and social behavior*. Sage, Newbury Park, CA. chapter Primitive emotional contagion, 151–177.
- Jennings, N.; Faratin, P.; Lomuscio, A.; Parsons, S.; Wooldridge, M.; and Sierra, C. 2001. Automated negotiation: Prospects, methods and challenges. *Group Decision and Negotiation* 10(2):199–215.
- Kamar, E., and Horvitz, E. 2009. Collaboration and shared plans in the open world: Studies of ridesharing. In *IJCAI*.
- Leung, K.; Bhagat, R.; Buchan, N.; Erez, M.; and Gibson, C. 2005. Culture and international business: recent advances and their implications for future research. *Journal of International Business Studies* (36):357–378.
- Lewicki, R.; Barry, B.; and Saunders, D. 2010. *Negotiation*. McGraw-Hill.
- Lin, R., and Kraus, S. 2010. Can automated agents proficiently negotiate with humans? *Communications of the ACM* 53(1):78–88.
- Lin, R.; Oshrat, Y.; and Kraus, S. 2009. Investigating the benefits of automated negotiations in enhancing people’s negotiation skills. In *AAMAS*.
- Maes, P.; Guttman, R.; and Moukas, A. 1999. Agents that buy and sell. *Communications of the ACM* 42(3):81–91.
- Nash, J. 1950. The bargaining problem. *Econometrica* (18):155–162.
- Nass, C., and Moon, Y. 2000. Machines and mindlessness: Social responses to computers. *Journal of Social Issues* 56(1):81–103.
- Nass, C.; Jonsson, I.-M.; Harris, H.; Reaves, B.; Endo, J.; and Brave, S. 2005. Improving automotive safety by pairing driver emotion and car voice emotion. In *Human Factors in Computing Systems Conference (CHI 2005)*.
- Nass, C. 2004. Etiquette equality: Exhibitions and expectations of computer politeness. *Communications of the ACM* 47(4):35–37.
- Raiffa, H. 1985. *The Art and Science of Negotiation*. Harvard University Press.
- Sycara, K., and Dai, T. 2010. *Handbook of Group Decision and Negotiation*. Springer Netherlands. chapter Agent Reasoning in Negotiation, 437–451.
- Van Kleef, G.; De Dreu, C.; and Manstead, A. 2004. The interpersonal effects of anger and happiness in negotiations. *Journal of Personality and Social Psychology* (86):57–76.
- Van Kleef, G.; De Dreu, C.; and Manstead, A. 2010. An interpersonal approach to emotion in social decision making: The emotions as social information model. *Advances in Experimental Social Psychology* (42):45–96.