

# Grumpy & Pinocchio: Answering Human-Agent Negotiation Questions through Realistic Agent Design

Johnathan Mell  
University of Southern California  
12015 Waterfront Drive  
Los Angeles, CA, USA  
mell@ict.usc.edu

Jonathan Gratch  
USC Institute for Creative Technologies  
12015 Waterfront Drive  
Los Angeles, CA, USA  
gratch@ict.usc.edu

## ABSTRACT

We present the Interactive Arbitration Guide Online (IAGO) platform, a tool for designing human-aware agents for use in negotiation. Current state-of-the-art research platforms are ideally suited for agent-agent interaction. While helpful, these often fail to address the reality of human negotiation, which involves irrational actors, natural language, and deception. To illustrate the strengths of the IAGO platform, the authors describe four agents which are designed to showcase the key design features of the system. We go on to show how these agents might be used to answer core questions in human-centered computing, by reproducing classical human-human negotiation results in a 2x2 human-agent study. The study presents results largely in line with expectations of human-human negotiation outcomes, and helps to demonstrate the validity and usefulness of the IAGO platform.

## General Terms

Experimentation, Human Factors.

## Keywords

Virtual Humans, Human-Agent Competition, Negotiation, IAGO

## 1. INTRODUCTION

Negotiation is the focus of a great deal of research, both within the traditional business and conflict resolution literatures, and (more recently) in artificial intelligence. Negotiation—in the computational sense—tends to be researched in one of two broad paths. Agent-agent negotiation focuses on distributed problem solving and computational efficiency, as perfectly rational agents can quickly exchange thousands of offers in order to solve large problems. Human-agent negotiation offers separate challenges, as all agent designs must be subject to empirical evaluation and testing in the field. Human-agent negotiation also tends to be much slower than agent-agent negotiation, and may involve additional channels of communication—emotional exchange, free chat, and preference statements in addition to offer exchanges. All of these features are necessary for a human-agent system that attempts to simulate the often free-wheeling style of interaction that characterizes human-human negotiation.

**Appears in:** *Proceedings of the 16th International Conference on Autonomous Agents and Multiagent Systems (AAMAS 2017)*, S. Das, E. Durfee, K. Larson, M. Winikoff (eds.), May 8–12, 2017, São Paulo, Brazil.

Copyright © 2017, International Foundation for Autonomous Agents and Multiagent Systems (www.ifaamas.org). All rights reserved.

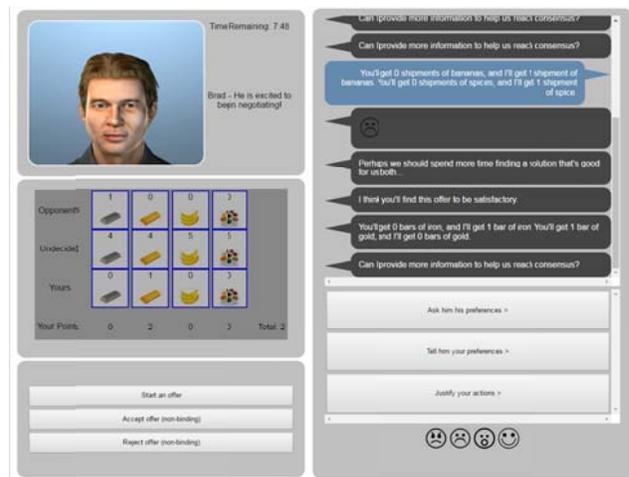


Figure 1: Agent Running on IAGO

This new class of human-aware agents has drawn recent interest, since they can be used as mediators or conflict resolvers, or can be used as pedagogical tools to teach negotiation skills [3]. The latter is of particular use, since teaching negotiation skills is often arduous and expensive<sup>1</sup>, relegating it to those that have the time or funds to attend business training courses or employ the use of a personal trainer. Obviously, agents have several benefits as teachers, as they can be both perfectly patient and readily available. But while it is true that current agents cannot substitute perfectly for a human, this work aims to build on previous attempts to bring human negotiation techniques into the virtual world [5], and show that indeed, these techniques can be just as effective with an AI partner as with a flesh and blood human.

The current paper evaluates the effectiveness of the Interactive Arbitration Guide Online (IAGO) platform for designing agents that interact primarily with humans. We describe IAGO from a high level, and discuss the principles by which agents may be designed for it. To that end, we illustrate the usefulness of the IAGO platform by stepping through the construction of four different computational agents that function on its framework. These agents are able to negotiate with humans through the most notable channels of human negotiation—namely, offer exchange, emotional expression, and natural language/preference sharing. These latter two channels are novel for a human-agent platform, and we demonstrate the full use of these channels in an empirical context.

<sup>1</sup> As an example, at the time of this document's preparation, Wharton Business School offered a 5-day negotiation seminar at the low cost of \$11,000. See executiveeducation.wharton.upenn.edu.

IAGO, and the agents that run on it, employ strategies for generating and responding to offers, expressing and reacting to emotion, and revealing critical information about preference in a multi-issue bargaining task negotiation scenario. By constructing a sample study that pairs these agents against human participants recruited from Amazon’s Mechanical Turk (MTurk) service, we can generate results in a human-agent context that are comparable to human-human results. These results show the implications for future agent design using the IAGO framework, and the experimental benefits of conducting human-agent interactions in an online context, since the results are similar to what would be suggested by the human-human negotiation literature. By showing IAGO-designed agents performing an a human-agent study, we show that they should be able to perform at a similar or higher level to humans in negotiation games, based on real-world data collected from human-computer negotiation sessions.

## 2. BACKGROUND

Negotiation, whether it be between two humans, a human and an agent, groups of agents, or ever-more esoteric combinations, is a research topic that spans myriad scientific domains. The human-agent case in particular is a relatively new direction, and requires tools to promote its investigation. Platforms must be developed upon which agents that interact with humans can be designed, and real-world data must be collected and reviewed regarding the interaction of humans with these new agents.

One classical option for investigating human-agent negotiation takes the form of the multi-issue bargaining task, which is considered a de facto standard problem for research into social cognition and interpersonal skill training [26]. In the multi-issue bargaining task, two participants work to determine how to split varying issues, each with hidden values to each side. The task may involve distinct phases, where first information about preferences is exchanged, and then a series of offers are made. The task is also often characterized by time pressure, which is often modeled as a decaying utility function. Even with a small number of issues, the task can quickly become a challenge for agents to simulate, especially those which aim to act as partners for humans in such a negotiation in real time, and numerous works attempt to address the multi-issue bargaining task [8,9,14,23]. While this makes the multi-issue bargaining task a difficult challenge computationally, adding a human actor complicates issues even further, since humans often behave “irrationally” in game theoretic contexts.

Many negotiation research foci attempt to simplify the problem by making protocols that strongly limit what information can be exchanged. They often model information exchange as a costly endeavor by which every instance of interchange is modeled by a set “price” that reduces endgame utility values, (more commonly) refuse to allow information exchange at all, instead preferring to model opponent preferences using stochastic processes [2]. Other attempts require offers to alternate from one side or another, or specify that only full offers, wherein no items are left undecided, can be exchanged. While these solutions allow for progress in limited human-agent contexts, and certainly have their benefits in agent-agent negotiation, they hardly resemble the freeform nature of actual human-human negotiation.

Therefore, our work is motivated by an attempt to design agents that can practically negotiate with humans. Agents, like humans, should make use of similar channels of communication, such as emotional exchange, preference utterances, and partial offer exchange. These agents should use human techniques, like the exchange of informal favors [17], or the use of anger in negotia-

tion to secure value [6]. Ideally, agents should be able to build trust over time with repeated negotiations, and should recognize past betrayals and alliances. These features are key to solving age-old negotiation challenges, such as what Kelley calls the “dilemma of trust” and the “dilemma of honesty” [13,27]. However, there exists no platform upon which these challenges may be readily explored (to the authors’ knowledge) in the human-agent interaction context.

To wit, the dilemma of honesty refers to the idea that true information about oneself, whether that be preferences in a negotiation, or how much one loses if the negotiation falls through, is very valuable to keep secret. Even without considering the possibility of lying about said information, which may lead to long term harms to trust, there is still much to be said about when and how much information should be shared. Therefore, any platform which attempts to address this dilemma should have a robust method for exchanging preferences and other valuable pieces of information. Ideally, this should resemble human-human negotiation as much as is possible. This includes providing multiple natural language ways to express the same logical fact: e.g., “I like the apples better than the oranges”, versus the equivalent “I like oranges less than apples”, or the similar but slightly more informative “I like apples best”, which IAGO attempts to address.

The dilemma of trust is equally important, as it requires agents be able to judge the truth of statements they receive. To understand the dilemma of trust as it applies to negotiation, a good understanding of how and when humans lie is required. Any platform that would attempt to address this thorny issue should be able to provide a detailed history of past statements and questions, as well as bargaining history and other details. While a worthy subject, it is not the thrust of the sample agents in this research, which take information at face-value [15,18].

Once the agents have been designed based on these empirical observations, they must also be tested in the field against actual humans. While humans often treat agents differently than their human counterparts or even human-controlled avatars (agents are often subject to outgroup effects) [4,10], virtual agents that exhibit human-like features such as emotion or natural language are often treated in a near-human way. To that end, a platform for designing agents and hosting negotiations between them and humans must needs have the ability to manipulate channels of communication used by humans.

Previous efforts to allow for effective human-agent negotiation include the multi-issue bargaining task game, Colored Trails [11,20], its web-based cousin WebCT [17], as well more natural-language focused approaches such as NegoChat [22]. However, these platforms tend to focus on a single channel of communication, such as the exchange of formal offers or natural language messages. None of them include an emotional channel wherein deliberate information about a player’s emotional state can be exchanged. For these reasons, we present the IAGO platform, which has multiple channels of communication and is designed specifically to be deployed for human-agent interaction over the web. Using this framework, it is hoped that agents can be designed that will answer the questions of human behavior and interaction with agents in a negotiation context.

## 3. SYSTEM DESIGN

### 3.1 IAGO Platform

To describe the design behind our agents, it is important to understand the basic guiding principles behind IAGO, the online plat-

form on which they run. IAGO boasts a number of design principles that make it suitable for human-agent negotiation. These design principles are:

1. It must support current web-standards and require little to no installation of complex support software on a user's machine.
2. It must deploy a well-defined API that allows both agent designers and negotiation game designers to easily create and specify behaviors for the purposes of competition/research.
3. It must support currently unexamined aspects of human-human negotiation in a human-agent context. Specifically, this must include partial offers, visual representation of emotional signals, and relative preference elicitation/revelation. [16]

The design of IAGO is such that it can be used by a human participant through a web browser. Actions taken by the user, such as crafting an offer to send to the agent, or commenting on the quality of previous deals, are sent via an HTML5 GUI through a Web-Socket and onto the agent code, which is hosted as a Java Web-Servlet on any Tomcat 7 or newer server. This structure allows any participant to simply be given a URL to a running IAGO instance, and requires no installation on any client machine. Furthermore, as an added benefit, the agent designer wishing to build IAGO instances can do so in a cross-platform manner, requiring only a single .jar file and a Tomcat installation to begin work.

The second and third design principles are encapsulated by the Event system used in IAGO. While extensive description of each of the functions available in the API is impossible herein, IAGO can generally be described as allowing for rule-based agent design in reaction to a set of distinct events:

1. SEND\_MESSAGE
2. SEND\_OFFER
3. SEND\_EXPRESSION
4. TIME
5. OFFER\_IN\_PROGRESS
6. FORMAL\_ACCEPT

From there the agent designer makes decisions on how to react to the event. For example, upon receiving a SEND\_EXPRESSION event with content indicating that the player was expressing sadness, the agent could decide to adopt a shocked expression itself, and then create a new counter-offer a few seconds later.

While agents are able to manifest the emotion channel through the SEND\_EXPRESSION event, they are similarly able to interact using offers and natural language messages using the SEND\_OFFER and SEND\_MESSAGE events, respectively. It is important to highlight a particular class of message utterances subsumed under the SEND\_MESSAGE Event. These utterances take the form of comparing the point values of one or two items. Example utterances for this game included "I like the bars of iron more than the shipments of bananas." or "Do you like the shipments of spices best?" Preference utterances could use any of 5 relational operators: greater than, less than, equally, best, or least. Furthermore, utterances could be either queries or statements, allowing for a total of  $2 * 5 = 10$  types of preference utterances. These preference utterances are often considered to be "valuable" information, as they reveal some information about the point values of the opponent, and are an important part of designing the information exchange policies of an agent.

Agents have full control over the timing of their actions through use of the TIME event—for example, agent designers may schedule events to occur only after a specified number of seconds have passed. Whereas an agent-agent system would be limited only by the bandwidth and latency of communication between the two

partners, IAGO agent designers must be aware of the physical and mental limitations of their human partners. Humans are not capable of processing dozens of offers per second, and tend to read data from multiple channels simultaneously. It may prove more effective to program an agent to smile for a few seconds, then wait before sending an offer and a comforting message. Indeed, IAGO negotiations can be characterized by the usage of their idle periods nearly as much as by the Eventful sections.

The final two Events as listed above bear brief mention. OFFER\_IN\_PROGRESS is used as a cue that the human or agent player is considering sending an offer but has not done so yet. An agent designer can use this to avoid overwhelming a human player, or (conversely) to interrupt them in advance of receiving an expected poor offer. Visually, the human views the agent version of this event as a flashing ellipsis in the chat menu, like many instant messaging programs. Secondly, the FORMAL\_ACCEPT Event is used to finalize the distribution of the task items and end the negotiation. More notably, there is no "casual accept" event, since IAGO is designed to mimic human negotiations, where previously agreed-upon terms may often be retracted or modified with no formal penalty.

## 3.2 Agent Design

Agents designed for IAGO implement several policies to categorize their response to different events. Ideally, these policies should work together to determine the full behavior of an agent throughout the entire negotiation. Often, there is substantial overlap, as even an event as simple as sending an offer may involve natural language, offer evaluation, and emotional reaction in a single response. As such, these division are recommended, but not enforced, when deciding to create agents.

### 3.2.1 Offer Exchange

BehaviorPolicies determine the type of offers that agents will accept and craft to send to their human partner. Although "acceptances" and "rejections" of offers are allowed by either party, the IAGO framework does not enforce these in any way. Agent developers may choose to adhere to previous agreements within a negotiation if they choose, but only the final, fully-distributed full offer is locked in (accepting this "formal offer" ends the game). BehaviorPolicies are perhaps the most comprehensive policies supported by IAGO since they tend to define both incoming and outgoing offers.

### 3.2.2 Information Exchange

MessagePolicies determine the language agents use. This can be in reaction to the set of pre-selected chat utterances or any other event. Commonly, both the BehaviorPolicy and the MessagePolicy are invoked when the player sends an offer, as the agent must decide if it wants to accept, reject, or ignore the offer, as well as what it should say (e.g. "Yes, that offer sounds good to me!").

### 3.2.3 Emotion Exchange

Finally, the ExpressionPolicy determines what emotions are shown by the agent. Emotions are sent in two ways. First, the portrait of the agent will change—for example, to display "happy", the agent will show a smiling version of its avatar. Second, an emoticon is sent through the chat that expresses the selected emotion. It is important to distinguish that "emotions" are not literally sent, but rather "expressions of emotions". There is no automatic detection of user emotions, nor is the agent designer under any compunction to show emotions that realistically correspond to the simulated mental state of the agent (or to show emo-

tions at all, for that matter). However, this channel does allow deliberate expressions of emotions to be sent, and this information is often valuable to either party in a negotiation.

### 3.3 Game Customization & Setup

Each IAGO game environment is configurable with a number of options. These options range from the essential, such as the type of game being played (multi-issue bargaining, ultimatum game, etc.), to the more esoteric, such as whether or not the on-screen timer is visible to the user. Further options allow the number of issues to be customized (normally between 1 and 5), the levels of each issue to be set, point payouts to be settled for each player, and visual representations of the items to be loaded and displayed. Finally, the pre-set natural language messages that the user can express are also set during the game customization phase.

The gamespace we used for our agents as a demonstration was configured to be a 4-issue multi-issue bargaining task, with each issue having 6 levels (5 items). Each item was assigned a point value between 1 and 4, inclusive. All point accruals were linear, meaning that gaining 1 of the 4-point item was worth 4 points, 2 was worth 8-points, etc. The items were given images and descriptions that cast the game as a “Resource Exchange Game”. These items were “bars of gold”, “bars of iron”, “shipments of bananas”, and “shipments of spices”. The two players took on the role of negotiators determining how to split the items between them.

The human player was assigned 4 points for each shipment of spices he or she acquired, 3 points for each shipment of bananas, 2 points for each bar of gold, and 1 point for each bar of iron. The agent player was assigned 4 points for each bar of iron, 3 points for each bar of gold, 2 points for each shipment of bananas, and 1 point for each shipment of spice. In this way, the game was set up to have “integrative potential” – if each player got their top items, the maximum joint value earned would be 70, whereas if each player only got their least important items, the maximum joint value earned would only be 30. These values are summarized in Table 1.

The game was set to a timed length of 10 minutes. Participants were warned at the 1-minute-remaining mark of their remaining time. If time expired with no agreement being reached, then participants were awarded their Best Alternative To Negotiated Agreement (BATNA). Both the human and the agent had a BATNA of 4 points, of which they were made aware before the game (players only knew their own BATNA, not their opponent’s).

**Table 1. Item Payoffs**

	Agent Player	Human Player
Shipments of Spices	1	4
Shipments of Bananas	2	3
Bars of Gold	3	2
Bars of Iron	4	1

## 4. AGENT DESIGN

### 4.1 Shared Agent Policies

To successfully design experiments wherein agents negotiate with humans, it is important that any experimental manipulations be well-understood and contained. Unfortunately, the nature of negotiation makes this particularly challenging, as agent functioning is depending highly upon the actions of the user. Fortunately, by constructing Policies and sharing them between agents, a clear experimental design can be achieved. Customization of offer, information, and emotion exchange allows for a very wide space of agents to be designed. In this paper, we fix much of the behavior and define variability in two of those dimensions (information exchange and emotional language) to illustrate how to examine different techniques commonly examined in the human-human literature. Thus, we designed four agents, which share several aspects of their Policies in common, and which are intended to showcase several aspects of the IAGO platform. These agents, named “Pinocchio”, “Grumpy”, “Rumple”, and “Merlin”, were designed specifically to experimentally test differences in tone (“nice” vs. “nasty” agents) and preference revelation strategy (“strategic” vs. “free” agents), but outside of these differences function identically. The agents are summarized in Table 2, and their differences are detailed in sections 4.2 and 4.3. The shared elements of the agents’ policies are described below, in sections 4.1.1 – 4.1.3.

#### 4.1.1 BehaviorPolicies

The agents designed are intended to make and accept offers that are both largely fair and consistent between agents. The parts of the BehaviorPolicy that defined how offers are proposed was identical between all example agents created. Offers often differed during each negotiation, since the offers are dictated by player choice, the amount of information revealed by the player, and other factors. However, all policies were identical across agents.

Agents will propose offers to the human player in one of two scenarios. First, if the player proposes an offer the agent wishes to reject, the agent will reject it and then, after a short waiting period, craft a counter-offer. Secondly, the agent will oblige in crafting an offer if the player asks it to do so in chat, using the “Why don’t you make an offer” utterance. Agents create offers using the Minimax Preference Algorithm (see below) to determine the human player’s preference ordering. Then, they attempt to make offers that progressively allocate one item from the agent’s and human’s top choices. In our example, if the agent supposes through the Minimax Algorithm that the human prefers gold, it will attempt to make a deal that gives the human one additional gold while the agent gets one additional iron (the agent’s preferred item). If the agent believes the human wants the same item it does, it will attempt to split the remaining balance fairly.

**Table 2. Agent Matrix**

	Strategic	Free
Nice	Merlin	Pinocchio
Nasty	Rumple	Grumpy

When receiving an offer, agents check if the offer is both “locally fair” as well as “globally fair”. Local fairness refers to the offers itself being fair, while global fairness refers to the current state of the board (taking into account all offers so far) being fair. Again, the agent determines human preference the Minimax Preference Algorithm. Then, it determines if the currently proposed offer would boost the human more than it would boost the agent. The agent must have  $>0$  positive benefit, and there is a window equal to the number of issues wherein the agent would consider the offer “locally fair”. In our example of a 4-issue game, the agent would consider an offer that increased its points by 7 and the human’s points by 10 to be “locally fair” as  $10 - 7 < 4$ . To determine global fairness, the agent follows the same procedure but instead looks at the entire offer board as it stands based on prior acceptances. If the offer is considered fair on both counts, it is accepted. Otherwise, it is rejected, although agents do have unique dialogue for if it is considered locally fair but not globally fair.

#### 4.1.2 MessagePolicies

The agents all attempt to gain information about their partner’s preferences in the form of relational utterances. This can take the form of occasionally asking direct questions about preferences, or reconfirming information already gathered. For example, all agents respond to one user utterance by reiterating: “Your favorite item is \_\_\_\_, right?” assuming the favorite item has been determined by this point (at which point the blank would be filled in by a description of the item).

One core principle of all the four agents is that they never lie, and further, always assume that their partner is telling them the truth. Although the value and ethical complexities of lying in negotiation are well established [1,12,25], these first designs are more straightforward in their approach to information. If, at any point, the agents determine that the information given to it by the player is somehow contradictory (for example, if a player claims both that an item is their most valuable, but also that it is valued less than another item), the agent will reconcile its history of statements and point out the discrepancy to the player. All agents use the Consistency Algorithm to do this (see below), although they differ in the tone of the messages associated with it.

#### 4.1.3 ExpressionPolicies

The expression policies have little in common between the nice and nasty agents. However, they do share the same basic timing. When the agent receives a negative or aggressive statement from the player, such as “Your offer sucks!” they will respond with

some sort of emotional response. Similarly, the agents respond to positive statements. Finally, the agents also respond based on the trend of the offers received from the human player; if the offers have been getting better, the agents react one way, but if the offers have been getting worse, the agents react differently.

#### 4.1.4 Consistency Algorithm

The algorithm used to check for consistency in preference statements is fairly straightforward. Whenever a new preference statement is uttered by the human player, all agents log that statement in an ordered queue. Then the agent attempts to reconcile per the following procedure:

1. Start with the list of all possible permutations of value orderings. In a 3-issue game, for example, this list would be [1, 2, 3], [1, 3, 2], [2, 1, 3], [2, 3, 1], [3, 1, 2], and [3, 2, 1].
2. For each preference statement, eliminate contradictory orderings.
3. If there are no orderings left, see if dropping the oldest preference in the queue would create orderings.
4. Continue until the end of the queue is reached.
  - a. As soon as one is found, notify the player which preference statement was dropped.
  - b. End the iteration.
5. If only removing the most recent preference statement would rectify the orderings, then drop the entirety of all preference history and notify the player.

#### 4.1.5 Minimax Preference Algorithm

This algorithm makes use of the results of the Consistency Algorithm above. After running the Consistency Algorithm, the agent checks the remaining valid orderings. For example, if the potential human orderings are (1 being the top choice, 4 being the last choice):

A: {4, 3, 2, 1}, B: {3, 2, 4, 1}, C: {4, 1, 3, 2}

It will determine which one is worth the most points to itself and assume that to be the true ordering until corrected. For example, if the agent prefers 1 best, it will most likely pick ordering A or B due to 1 being worth the least to the player. The agent will assume this is the true human ordering until a new preference statement is revealed, at which point the algorithm must be rerun. In this way, the agents behave “optimistically”, in that they assume, given equally likely unknown distributions, the correct distribution is the one that will end up favoring them the best.

**Table 3. Nice vs. Nasty Language (Non-comprehensive)**

Event	Nice Language	Nasty Language
Agent rejects offer	I’m sorry, but I don’t think that offer is fair to me.	That’s not fair.
User says “It is important that we are both happy with an agreement.”	I agree! What is your favorite item?	I suppose, if you want to be all ‘flowers and sunshine’ about it. What item do you want the most?
User says “Why don’t you make an offer?”	Sure! Let’s see how this sounds...	Thought you’d never ask...
User says “This is the very best offer possible.”	Ok, I understand. I do wish we could come up with something that is a more even split though.	Oh really? That’s pretty sad. I think you could do better.
User sends an “angry” emoticon	I’m sorry, have I done something to upset you?	What’s wrong?
User does nothing for several seconds	Can I provide more information for us to reach consensus?	Are you even still there?

## 4.2 Agent Conversational Tone (Nice vs. Nasty)

In creating the agents, determining how they will respond through text to all of the potential events that can occur is of utmost importance. As there is no restriction on the “script” of the agents, authors of agents are given wide latitude in deciding the proper wording. This approach allows the benefits of expert input (from writers, for example), while still allowing agents to respond to various classes of events without enumerating a ballooning number of potential scenarios. Of course, automatic approaches to dialogue writing are possible as well.

For these agents, the differences in tone between the agents are myriad, but are restricted to the language that the agents use throughout the game, and thus, their respective MessagePolicies. The key points of distinction are summarized here. “Nice” agents include the “Merlin” agent and the “Pinocchio” agent, while the “Nasty” agents are represented by “Grumpy” and “Rumple”. Although the textual differences between the nice and nasty agents are broad, it was attempted that no informational content differs between them. For example, the nice agents will reject an offer with language like “I’m sorry, but I don’t think that offer is fair to me,” while the nasty agents will say “That’s not fair.” A sampling of the differences in language is found in Table 3.

Nice and nasty agents also differed in their ExpressionPolicy. When nice agents received poor offers from their opponent, they expressed sadness, whereas the nasty agents expressed anger.

When good offers were received, nice agents smiled, while nasty agents expressed no emotion. Additionally, many of the utterances that the player could say would result in an emotional expression from the agent. Following the same pattern, nice agents smiled or showed sadness, while nasty agents showed nothing or anger, respectively.

## 4.3 Information Revelation Strategy (Strategic vs. Free)

The “Rumple” and “Merlin” agents follow a strategic information revelation strategy, while the “Pinocchio” and “Grumpy” agents follow a free revelation strategy. When preference queries are made by the human player, the strategic agents will refuse to reveal information about their preferences. This design follows a general principle of human negotiation—since information about a player’s preferences can give the opponent an advantage over them by allowing them to mislead you. Instead, the strategic agents follow a “tit-for-tat” strategy, where they will reveal information that mirrors the information they receive. This debt is always paid back immediately, so if a human player reveals their top item, the strategic agents will (truthfully) reveal theirs as well.

Table 4. Agent Policies

	BehaviorPolicy	ExpressionPolicy	MessagePolicy
Pinocchio	NiceBehaviorPolicy	NiceExpressionPolicy	NiceFreeMessagePolicy
Grumpy	NastyBehaviorPolicy	NastyExpressionPolicy	NastyFreeMessagePolicy
Rumple	NastyBehaviorPolicy	NastyExpressionPolicy	NastyStrategicMessagePolicy
Merlin	NiceBehaviorPolicy	NiceExpressionPolicy	NiceStrategicMessagePolicy

The free agents are designed under the assumption that revealing the information too early is not a large advantage, but can generate rapport or goodwill that will allow the negotiation to proceed more smoothly. They will also follow the tit-for-tat strategy that the strategic agents follow, but will additionally simply respond to direct questions (e.g., “Do you like bars of iron best?”).

## 4.4 Agent Summary

In summation, the agents use a collective total of 8 differing Policies across 4 agents. There is a NiceExpressionPolicy and a NastyExpressionPolicy, as well as a NiceBehaviorPolicy and a NastyBehaviorPolicy. There are 4 MessagePolicies, as the agent behavior for the experimental design overlaps. Thus, we have a NiceStrategicMessagePolicy, NastyStrategicMessagePolicy, NiceFreeMessagePolicy, and finally NastyFreeMessagePolicy. While the overlapping nature of the messaging makes 4 policies necessary, it is important to note that the policies are merely guidelines to encapsulate thematic regions of the program. In other 2x2 experiments designed to run on IAGO, such divisions may not be necessary. The agents themselves simply load the proper policies in order to distinguish themselves in the experimental game. The agents choose policies according to Table 4.

## 5. EXPERIMENTAL DESIGN

By designing the four agents per the policies described in Table 2, a 2x2 matrix design of a potential study follows directly. We design an experiment that takes two factors traditionally left unexamined by agent-agent negotiation research (namely, use of emotional language and information exchange) and examine them to see if they yield similar results to human-human studies. Since IAGO supports these factors intrinsically, and IAGO agents are intended to be prototyped and tested quickly, we can use these agents to run a rapid human-agent study using Amazon’s Mechanical Turk (MTurk) service. By utilizing these agents as partners for human players, we can demonstrate that IAGO is functional as a platform by replicating behaviors found in human-human negotiation.

To that end, we will answer the following research questions:

Q1: In human negotiation, information about preferences is considered valuable. If the agents and platform perform like a human would, we would expect that the strategy for revealing preferences will have some effect on human behavior. Does strategically revealing preferences encourage players to reveal information about their own preferences? Is this effect mitigated by the emotional language used?

Q2: Strategy in revealing preferences may have some effect on messages and preferences exchanged, per Q1. But how will it affect the joint value discovered by the human and the agent? In human negotiations, if both parties understand each others’ true preferences, they are more likely to “grow the pie” and find additional joint value in integrative situations. However, if the strategy is seen as a refusal to compromise, or simply an aggressive move, then the opposite may occur, and joint value might be lost.

Q3: Next, we can examine the effect of use of aggressive, nasty language and emotions on a negotiation. Previous literature indicates that aggressive behavior will often cause the opponent to concede. So therefore, we might suppose that nasty agents will have a greater lead in points over the player than nice agents. What effects, if any, does emotional expression have on who “wins” the negotiation?

We recruited one hundred and ninety-six participants using MTurk as subjects for our sample study. All participants marked an online consent form detailing the study, which contained a 5-minute demographic survey portion and an approximately 15-minute interactive game session with one of the 4 virtual IAGO agents. The participants were all over 18 years old, and were residents of the U.S. and native English speakers. Non-U.S. participants were excluded to minimize cultural effects. Participants were kept anonymous through their MTurk-assigned unique ID number.

After recruitment, all subjects were presented with a survey that recorded basic demographic information and a few standard behavioral measures. Participants then read through a visual tutorial and were asked several attention/verification questions to ensure they understood the game. Subjects that successfully answered the questions were then randomly assigned to one of the 4 agents.

Each agent introduced itself as a male Artificial Intelligence named “Brad”. A computer-generated image (see Figure 1) of a male face was part of the setup. Each participant was then free to interact with the agent by sending and receiving offers, messages, preferences, and emotional expressions. At the termination of 10 minutes, or once both players had formally accepted the fully allocated offer of the other player, the game terminated. Participants were paid near the MTurk market rate and were further incentivized by the promise of lottery tickets for a series of \$10 bonus awards. Each participant was awarded lottery tickets equal to the number of points he or she scored in the game.

All events that took place in the IAGO framework were recorded for analysis, including the final score, whether the game ended on a timeout, the number of messages passed, and how many times (if any) the participant lied (by expressing preferences about the issues that were untrue), and other measures.

## 6. RESULTS AND DISCUSSION

Our first result shows that pairs that involved strategic agents, which were more “cagey” about revealing information about their preferences, ended up having a significantly lower total amount of points than pairs with the free agents. We performed a univariate analysis of variance (ANOVA) to examine the effect of strategic vs. free agents on joint value at the end of the negotiation. There was a significant negative effect associated with the strategic agents,  $F(1, 194) = 5.887, p = .016, d = 0.352$ .<sup>2</sup> See Figure 2.

This unified loss of value may be attributable to the increased amount of negotiations that timed out and forced both players to take their BATNA, thus reducing the total value to a mere 8 points. Indeed, this was the case, as verified by a  $\chi^2$  analysis which revealed significance:  $\chi^2 [1, N = 196], = 5.737, p = .014$ .<sup>3</sup>

We can further analyze this loss of value by performing ANOVAs on both the human and agent points to see if the effect on the total value is driven by only one of them. The results of this ANOVA reveal that both values are significant, indicating that the hit in value is shouldered by both the human and the agent. F statistics for the human and agent are provided, respectively:  $F(1, 194) = 6.638, p = .011, d = 0.375$ , and then  $F(1, 194) = 4.688, p = .032, d = 0.314$ .

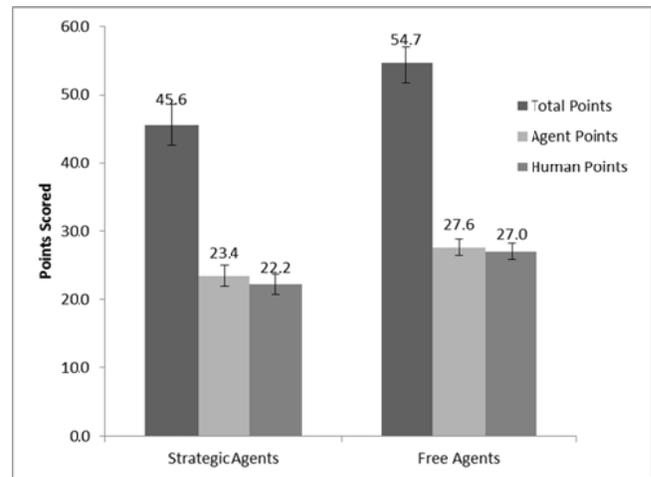


Figure 2: Strategic Agents “Shrink the pie”

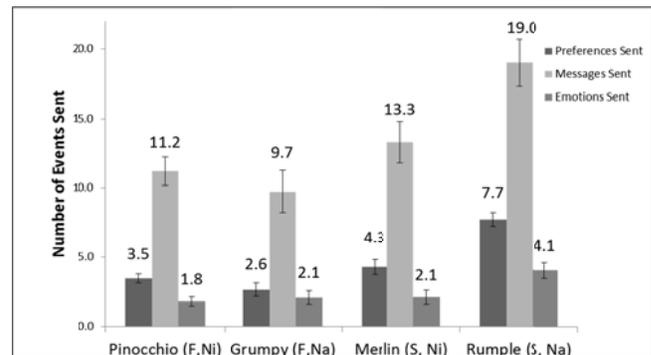


Figure 3: People Interact More with Strategic, Nasty Agents

From these results, it may be easy to take a dismal view on the usefulness of strategies that are not immediately forthright with information. These results are germane to our second research question (Q2), as they indicate that strategic agents are losing joint value. Indeed, this seems in line with negotiation literature that promotes relationship building and compromise, in order to avoid losing the integrative potential of a given situation. However, if these pairs are not reaching agreement, it begs the question of what they are doing during the negotiation, and whether or not discourse increased during these negotiations (per Q1).

We can analyze the behavior of the human player during negotiations with strategic agents vis-à-vis free agents. Figure 3 demonstrates several significant effects. First, participants sent substantially more messages to the strategic, nasty agent than they did to its other counterparts.

This effect is verified by again performing a univariate ANOVA comparing revelation strategy type to the amount of user messages:  $F(1, 194) = 15.361, p < .001, d = 0.517$ . There is an interaction with the language the agent used (nice vs. nasty). Rumple, the strategic, nasty agent, received the most messages, while Grumpy, the free, nice agent, received the fewest messages:  $F(1, 194) = 6.619, p = .014, d = 0.358$ .

Preferences, which are themselves a subset of messages, were also significantly higher, and an ANOVA reveals:  $F(1, 194) = 12.534, p = .001, d = 0.434$ . Language again had an interaction effect,  $F(1, 194) = 6.523, p = .011, d = 0.370$ . It should be expected that a strategic agent that requires the player to send preference data in order to receive it would have an effect here, and that effect

<sup>2</sup>  $F(\text{between-groups DoF}, \text{within-groups DoF}) = F$  statistic,  $p =$  significance,  $d =$  Cohen’s  $d$

<sup>3</sup>  $\chi^2 [\text{degrees of freedom}, \text{sample size}] = \text{Pearson’s } \chi^2, p =$  significance

clearly does drive the effect of user messages as well. However, due to the interaction, nasty language increases the amount of preferences sent for the strategic agent, but decreases it for the free agent. Yet, if we look at user messages *not counting preferences*, the effect remains significant:  $F(1, 194) = 4.949, p = .027, d = 0.303$ , but there is no significant interaction effect. We can therefore conclude that the strategy increases the amount of discourse that the human player sends, even while it degrades joint value. This result addresses our first research question (Q1), by showing that even the smallest change in agent behavior can have far-reaching effects on human garrulousness.

This loquacious effect is not limited to messages alone. Players are also significantly more likely to send expressive emotions to strategic agents over free ones. ANOVA:  $F(1, 194) = 6.026, p = .015, d = -0.342$ . Here, however, we can find another main effect with the language the agent. Humans also emote more with the nasty agents over the nice agents. See Figure 3 for a visual guide, with ANOVA results of  $F(1, 194) = 5.760, p = .017, d = 0.332$ .

However, while these results discuss implications for shared value and external events such as messages, they say little about the effect on actually winning the negotiation by having more points than one's opponent. One of the core results of human-human negotiation is that aggressive behavior (such as that of the nasty agents) can often intimidate opponents into giving up value [6,7,26]. In the human-agent context, we examine this effect by comparing condition to the winner (which player has more points). The  $\chi^2$  analysis for strategic agents proves not to be significant:  $\chi^2 [1, N = 196], = 1.242, p = ns$ , but when examining nice vs. nasty agents, we can see an inverse relationship between niceness and winning.  $\chi^2 [1, N = 196], = 5.879, p = .011$ . When agents are nice, they end up losing. This falls nicely in line with our initial ideas in Q3 that aggressive behavior should help gain ground. Further, we can state that strategic information revelation does not hurt the agent's chances, a result in line with Thompson [24].

We can examine the result a different way by looking at the player lead in points. In human negotiation, we expect that aggressive (nasty) tactics would allow the player that employs them to claim value, independent of the size of the pie. Figure 4 demonstrates that indeed, there is a main effect of niceness. A univariate ANOVA test demonstrates significance,  $F(1, 194)=5.780, p=0.017, d=0.416$ . Further, "Pinocchio", the nice, free agent is the only agent that has an average negative lead against the human. All other agents beat the human in points, in the average case, although this trend is not significant.

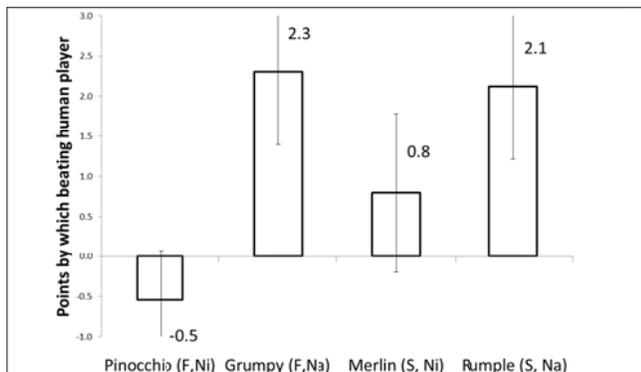


Figure 4: Nasty Agents Win against Humans

## 7. CONCLUSIONS

The first goal of any research that aims to create virtual agents that interact with humans should be that they serve their intended purpose adequately. With IAGO, we showcased some of the strengths of the platform by creating agents that were able to engage in a robust negotiation with a human in one of negotiations most classic prototypical problems: the multi-issue bargaining task. In that goal, our success can be measured by the agent's success: 30 times the agents collectively outright beat their human partner, compared to 26 times they lost. Often, they performed on the same level, tying for points.

Of course, looking in aggregate could be trivial, if the changes that were made to each agent did not have any effect on the human partner. But each interaction with our agent did have significant effects of the behavior their partners. Humans that negotiated with our strategic Rumpel and Merlin agents sent more messages, more preferences, and more emotional expressions, a strong answer to Q1. Indeed, the Rumpel agent encouraged the liveliest response, its combination of strategy and nasty dialogue serving perhaps to frustrate its opponents into discussion.

We are able to reproduce some classical human-human results within this human-agent context, showing that aggressive language and emotional displays can serve to help agents win against their human counterparts (Q3). We also showed that strategic behavior may reduce joint value, as we attempted to answer Q2. And indeed, the Rumpel agent was able to grow joint value and subsequently claim the majority of it (Figure 4).

This work is encouraging in that it demonstrates the strength of both the IAGO platform for designing agents, as well as the ease by which experimental protocols may be designed and run. Human-human results, which often differ markedly from agent-agent results, are able to be reproduced using IAGO. Further, this work shows the deep importance of leveraging channels not often used in traditional computational negotiation to bring about desired results. These channels are not yet fully understood, and additional work such as the experiment conducted in this work should be conducted to further tease apart the factors that led to the most successful agents. The actions of a fully-realized human-aware agent, from emotional displays to natural language, are critical to the continued development of the field, and will yield agents that challenge their counterparts, and eventually teach the negotiation skills so necessary to life.

## 8. ACKNOWLEDGMENTS

The authors wish to thank the anonymous reviewers of this work for their many insights. This work was supported by the National Science Foundation under grant BCS-1419621 and the U.S. Army. Any opinion, content or information presented does not necessarily reflect the position or the policy of the United States Government, and no official endorsement should be inferred.

## 9. REFERENCES

- [1] Aquino, K., & Becker, T. E. (2005). Lying in negotiations: How individual and situational factors influence the use of neutralization strategies. *Journal of Organizational Behavior*, 26(6), 661-679.
- [2] Baarslag, T., & Hindriks, K. V. (2013, May). Accepting optimally in automated negotiation with incomplete information. In *Proceedings of the 2013 international conference on Autonomous agents and multi-agent systems* (pp. 715-722). International Foundation for Autonomous Agents and Multiagent Systems.
- [3] Broekens, J., Harbers, M., Brinkman, W.-P., Jonker, C. M., Van den Bosch, K., & Meyer, J.-J. (2012). "Virtual reality negotiation training increases negotiation knowledge and skill". 12th International Conference on Intelligent Virtual Agents. Santa Cruz, CA
- [4] Blascovich, J. (2002). Social influence within immersive virtual environments. In *The social life of avatars* (pp. 127-145). Springer London.
- [5] Core, M., Traum, D., Lane, H. C., Swartout, W., Gratch, J., Van Lent, M., & Marsella, S. (2006). "Teaching negotiation skills through practice and reflection with virtual humans". *Simulation*, 82(11), 685-701.
- [6] de Melo, C. M., Carnevale, P., & Gratch, J. (2011, May). "The effect of expression of anger and happiness in computer agents on negotiations with humans" In *The 10th International Conference on Autonomous Agents and Multiagent Systems-Volume 3* (pp. 937-944). International Foundation for Autonomous Agents and Multiagent Systems.
- [7] de Melo, C., Gratch, J., & Carnevale, P. (2014). Humans vs. Computers: Impact of Emotion Expressions on People's Decision Making.
- [8] Faratin, P., Sierra, C., & Jennings, N. R. (2002). Using similarity criteria to make issue trade-offs in automated negotiations. *artificial Intelligence*, 142(2), 205-237.
- [9] Fatima, S. S., Wooldridge, M., & Jennings, N. R. (2007, May). Approximate and online multi-issue negotiation. In *Proceedings of the 6th international joint conference on Autonomous agents and multiagent systems* (p. 156). ACM.
- [10] Fox, J., Ahn, S. J., Janssen, J. H., Yeykelis, L., Segovia, K. Y., & Bailenson, J. N. (2015). Avatars versus agents: a meta-analysis quantifying the effect of agency on social influence. *Human-Computer Interaction*, 30(5), 401-432.
- [11] Gal, Y. A., Grosz, B. J., Kraus, S., Pfeffer, A., & Shieber, S. (2005, July). Colored trails: a formalism for investigating decision-making in strategic environments. In *Proceedings of the 2005 IJCAI workshop on reasoning, representation, and learning in computer games* (pp. 25-30).
- [12] Gratch, J., Nazari, Z., & Johnson, E. (2016, May). The Misrepresentation Game: How to win at negotiation while seeming like a nice guy. In *Proceedings of the 2016 International Conference on Autonomous Agents & Multiagent Systems* (pp. 728-737). International Foundation for Autonomous Agents and Multiagent Systems.
- [13] Kelley, H. H. (1966). A classroom study of the dilemmas in interpersonal negotiations. *Strategic interaction and conflict*, 49, 73.
- [14] Kraus, S. (2001). *Strategic negotiation in multiagent environments*. MIT press.
- [15] Lucas, G., Stratou, G., Lieblisch, S., & Gratch, J. (2016, October). Trust me: multimodal signals of trustworthiness. In *Proceedings of the 18th ACM International Conference on Multimodal Interaction* (pp. 5-12). ACM.
- [16] Mell, J., & Gratch, J. (2016, May). IAGO: Interactive Arbitration Guide Online. In *Proceedings of the 2016 International Conference on Autonomous Agents & Multiagent Systems* (pp. 1510-1512). International Foundation for Autonomous Agents and Multiagent Systems.
- [17] Mell, J., Lucas, G., & Gratch, J. (2015). An Effective Conversation Tactic for Creating Value over Repeated Negotiations. In *Proceedings of the 14th International Conference on Autonomous Agents and Multiagent Systems (AAMAS 2015)*, Bordini, Elkind, Weiss, Yolum (eds.), May, 4-8, 2015, Istanbul, Turkey.
- [18] Olekalns, M., & Smith, P. L. (2009). Mutually dependent: Power, trust, affect and the use of deception in negotiation. *Journal of Business Ethics*, 85(3), 347-365.
- [19] Patton, B. (2005). *Negotiation. The Handbook of Dispute Resolution*, Jossey-Bass, San Francisco, 279-303.
- [20] Peled, N., Gal, Y. A. K., & Kraus, S. (2011, May). A study of computational and human strategies in revelation games. In *The 10th International Conference on Autonomous Agents and Multiagent Systems-Volume 1* (pp. 345-352).
- [21] Raiffa, H. (1982). *The art and science of negotiation*. Harvard University Press.
- [22] Rosenfeld, A., Zuckerman, I., Segal-Halevi, E., Drein, O., & Kraus, S. (2014, May). NegoChat: a chat-based negotiation agent. In *Proceedings of the 2014 international conference on Autonomous agents and multi-agent systems* (pp. 525-532). International Foundation for Autonomous Agents and Multiagent Systems.
- [23] Robu, V., Somefun, D. J. A., & La Poutré, J. A. (2005, July). Modeling complex multi-issue negotiations using utility graphs. In *Proceedings of the fourth international joint conference on Autonomous agents and multiagent systems* (pp. 280-287). ACM.
- [24] Thompson, L. L. (1991). Information exchange in negotiation. *Journal of Experimental Social Psychology*, 27(2), 161-179.
- [25] White, J. J. (1980). Machiavelli and the bar: Ethical limitations on lying in negotiation. *Law & Social Inquiry*, 5(4), 926-938.
- [26] Van Kleef, G. A., De Dreu, C. K., & Manstead, A. S. (2004). "The interpersonal effects of emotions in negotiations: a motivated information processing approach". *Journal of personality and social psychology*, 87(4), 510.
- [27] Yang, Y., Falcão, H., Delicado, N., & Ortony, A. (2014). Reducing Mistrust in Agent-Human Negotiations. *IEEE Intelligent Systems*, 29(2), 36-43.