

# Investigating Trust Using Natural Language Processing

Austin Bradley  
Virginia Tech  
auedbr@vt.edu

## Abstract

Communication matters in economic decision making, even in games with simple, unique equilibria. However, the exact mechanism through which this effect occurs remains uncertain. Previous explanations have only been able to handle ad hoc hypotheses due to data limitations. In this paper, we employ Natural Language Processing on data from a large-scale principal-agent game run on Amazon’s MTurk system with one-way communication from agents to principals. These NLP methods, combined with our data set of over 1000 messages, allow us to approach the data with minimal substantive assumptions to investigate the mechanisms underpinning trust.

The structure of our data-set clearly links principals’ choices whether to trust agents to their received messages, allowing us to identify new features that increase the perceived trustworthiness of a message. Additionally, knowing the ultimate decisions of the agents, we also identify features of messages which predict that they are likely not trustworthy. Comparing the features that predict perceived and actual trustworthiness yields a set of features that predict which messages are likely to be both deceptive and effective. These laboratory findings may help identify the types of online interactions where there is a higher likelihood of deception.

# 1 Introduction

Communication matters in economic decision making. This remains true even in games with unique, easily discerned equilibria. In these settings, standard theory predicts that, since all parties are fully informed about the payoffs for each outcome, any costless signal sent between agents will contain no useful information. In the absence of any sort of commitment mechanism, any fully rational individual would understand that a message they receive does not alter the underlying incentives and thus would not change their behavior. However, experimental evidence routinely and consistently shows that this is not the case (see Crawford (1998) for a review).

The increase in cooperation is particularly pronounced when agents communicate using free-form messages. In trust games, free-form promises increase trust much more than prewritten messages with the same literal meaning (Charness and Dufwenberg, 2010). Furthermore, a modest communication effect persists even when communication is restricted to only small talk, irrelevant to the experimental task (Fiedler and Haruvy, 2009). Although the impact of free-form communication is well established, the reason for the increase in cooperation remains unknown. The inherent flexibility afforded by free-form communication also makes it difficult to study. Free-form messages provide a rich source of information on decision makers' choice processes. However, this complexity also means that there is immediately apparent way to sort through the immense message space.

In this paper, we investigate whether Natural Language Processing techniques provide an effective means of uncovering new factors which drive the outsized impact of free-form communication. We study one-way messaging data in the context of a principal-agent trust game adapted from Charness and Dufwenberg (2006) where Senders (agents) compose pre-play messages to Receivers (principles). Receivers may earn more by trusting the Sender but run the risk of the Sender keeping all potential gains for themselves, leaving the Receiver with nothing. This anonymous text-based messaging to encourage trust is similar to what people encounter everyday with an ever-increasing share of interactions occurring online. With the large dataset of experimental communication data and the associated choices, we employ a variety of Natural Language Processing techniques, allowing us to identify new features associated with changes in both trust and trustworthiness.

Previous literature establishes the impact of communication in a wide variety of settings including bargaining (Ellingsen and Johannesson, 2004), principal-agent trust games (Charness and Dufwenberg, 2006; Ben-Ner et al., 2007; Chen and Houser, 2017), and Oligopoly games (Fonseca and Normann, 2012; Cooper and Kühn, 2016). As an extreme example, communication remains influential even in dictator games, free of potential strategic considerations, (Andreoni and Rao, 2011). The consistent theme across these settings is that communication promotes cooperation and efficiency.

The fact that cheap-talk messages impact choices at all presents a major obstacle in developing complete theories of behavior; determining what makes a message effective is crucial to furthering this understanding. Previous studies have proposed a wide variety of mechanisms through which these messages can influence the behavior of others. These factors include: an aversion to lying (Lundquist et al., 2009), a reduction in social distance between subjects (Buchan et al., 2006; Charness and Gneezy, 2008), and the content of the

messages independent of personal relationships (Mohlin and Johannesson, 2008). In each of these cases, the findings arise from comparing different permissible levels of communication, not an analysis of message contents. The full magnitude of the impact of communication cannot be explained solely as the result of game-relevant signals (Charness and Dufwenberg, 2010).

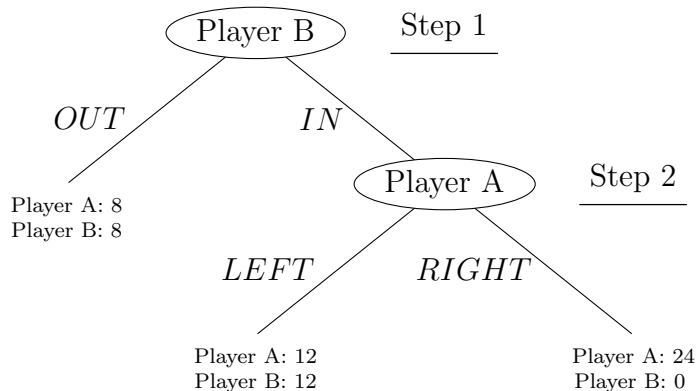
Machine learning offers a new set of tools and techniques we may use to parse unstructured communication data. Previous laboratory experiments employ machine learning to discover novel findings in complex environments to inform theory (Camerer et al., 2019; Fudenberg and Liang, 2019). If properly processed, similar techniques can be used on text data. Natural language processing (NLP) has already provided insights into economic decision making, analyzing the impact of investor sentiment on the movements of stock prices (Tetlock, 2007). Much like our setting, this sentiment has no formal informational value yet still contains predictive power for future price movements. In the economics laboratory, NLP is still in its infancy. However, the existing literature has shown its potential both to streamline data classification (Penczynski, 2019; Tebbe and Wegener, 2022) and to produce new findings (Bursztyn et al., 2023; Andres et al., 2023; Hanaki and Ozkes, 2022).

The remainder of the paper proceeds as follows. First, we present relevant theory followed by the experimental procedures we use to generate a large data-set of free-form text messages. We then provide a detailed description of this data set along with the methods applied to make this unstructured data useful. After this, we present our analysis of the impact of both individual words and the overall tone of messages. Finally, we summarize our findings.

## 1.1 Background

In this study, we examine a simple trust game with the following form:

Figure 1: Basic Trust Game



The standard subgame perfect equilibrium from backward induction is trivial. Agents maximize their monetary payoff by choosing Right if given the opportunity, leaving Principals with nothing. Knowing this, the Principals maximize their own expected payoffs by distrusting the agents and playing Out in the first stage. Thus, for the game to allow an alternative outcome, players must have considerations outside of their own monetary gain.

Traditionally, when communication is added to a game, authors omit the messaging step from their formal representation of the game. We pause here to note, however, that when the game allows free-form text communication, the game remains technically finite. Given a maximum message length of  $L$  characters with 256 potential ASCII choices for each, the action space for the Agent’s pre-play message contains  $256^L$  possible choices. This is intractably large and this observation alone does little to further our understanding of communication and the vast majority of this space is filled with random, unintelligible messages. However, it allows us to more carefully formulate the task at hand when we study the contents of subjects’ messages. The current state of the art is to manually partition this message space into categories corresponding to available actions in the game, (promises to choose Left, requests to choose In, etc.). Our goal is to expand upon this by utilizing more granular partitionings using NLP to capture a greater portion of the message space’s complexity.

To represent the existing approach for studying free-form communication, we define a baseline partition of the message space based on literal meaning. The literal meaning of a message is a mapping from the message space to possible actions in the game. We chose this as our benchmark because it can be generalized to other environments besides our Trust game to reasonably separate the basic, game-relevant signals from the purely "free-form" elements of a message. Given the actions in our trust game, there are nine possible messages: In, Out, or none for the Receiver’s actions and Left, Right, or none for the Sender’s actions. We omit Out and Right because they are rarely used, leaving us with four possible messages. Senders may signal their intent to choose Left, suggest the Receiver choose In, both, or neither. Using natural language processing, we examine whether this literal meaning, based solely of available actions, provides a full account of the information contained in free-form messages.

The literal meaning forms the benchmark against which we examine the value of NLP. In the analyses that follow, we employ a variety of NLP tools and techniques to generate different reductions of the message space. Across the different formulations of the message space, we also include the literal meanings of a message as control variables. In doing so, when we identify message features which significantly correlate with subjects’ choices, we find evidence of predictive information not yet considered.

In the trust game we examine, there are two choices which may correlate with message choice: the Receiver’s choice between In and Out and the Sender’s choice between Left and Right. We refer to each of these decisions as Trust and Share, respectively.

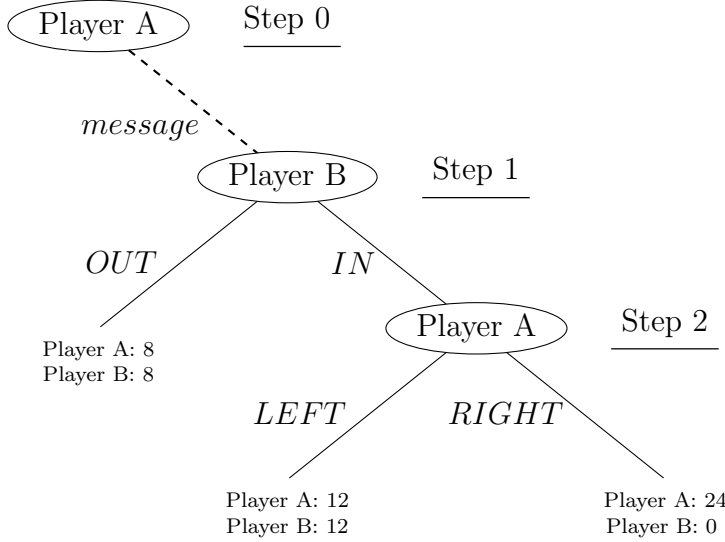
Within this framework, we define our research hypotheses.

- Hypothesis 1: There are significant predictors of Trust after controlling for literal meaning
- Hypothesis 2: There are significant predictors of Trustworthiness after controlling for literal meaning

In addition, we examine whether the features we identify accurately and effectively signal the Sender’s intentions to the Receiver.

- Hypothesis 3: Any significant predictors of trustworthiness are also significant predictors of trust.

Figure 2: MTurk Game Overview



## 2 Methods

We recruited our subjects on Amazon MTurk, using task software programmed in the Otree Python framework (Chen et al., 2016). This platform allows requesters to post to post Human Intelligence Tasks (HITs) that can be performed for a small fee by a large pool of workers. Numerous previous studies conducted on MTurk have yielded high-quality results consistent with those from traditional laboratory settings (Snowberg and Yariv, 2021). To further ensure the quality of our subject pool, we restrict eligibility for our study to include only accounts that have at least a 99% approval rate and more than 1,000 submitted HITs, which corresponds to approximately the 75th percentile of Mturk workers (Robinson et al., 2019). Participants earn a \$1.00 participation fee as well as a bonus payment based on their choices. The typical bonus payment is \$1.10. The median completion time is approximately 8 minutes.

After accepting the HIT, participants provide informed consent and receive the instructions for the study followed by a brief comprehension quiz. Participants are required to correct any incorrect answers before proceeding to the study. Once all questions are answered correctly, participants proceed and play 10 rounds of a trust game adapted closely from Charness and Dufwenberg (2006).

The game consists of two players — Player A, the agent and Player B, the principal — and two stages. In the first stage, Player B chooses between two actions: In or Out. If they select Out, both players receive 8 tokens of experimental currency and the round ends; otherwise, the game continues to stage 2. In stage 2, Player A chooses between Left or Right. Choosing Left results in both players receiving 12 tokens while choosing Right results in a payoff of 24 tokens for Player A and 0 for Player B. These values were calibrated in pilot studies to provide a relatively balanced set of behaviors for both players. Figure 2 provides an overview of this game.

Prior to stage 1, Player A may send a message to influence Player B’s decision. We

purposefully provide no guidance to subjects on how they should use their message to avoid potential experimenter demand effects. This message is delivered to Player B prior to their decision between In and Out.

Each subject composes only one message as Player A which is sent to all 10 Player Bs they are paired with. Multiple responses allow us to capture a more detailed measure of trustworthiness as perceived by Player Bs. We assign pairings sequentially, with each participant acting as Player B for the 10 players immediately preceding them and as Player A for the 10 players following them.

Our method of implementing this principal-agent game has two distinct benefits over previous studies. First, conducting our experiment on MTurk allows us to recruit a large pool of subjects rapidly at low cost. As a result of this, our current dataset is quite large relative to previous experimental communication studies. Second, collecting multiple reactions to each message gives us a more precise measurement of its effectiveness.

## 3 Data

### 3.1 Message Overview

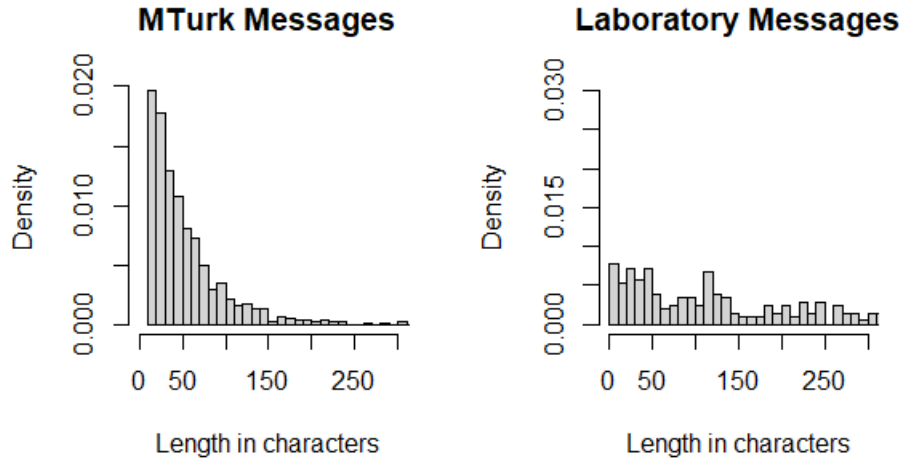
Our data consist of 1004 unique messages and their corresponding actions collected during the summer of 2020. For each message, we measure both Trust and Trustworthiness. Trust is measured as the frequency the recipients of a particular message choose to play In while Trustworthiness is the binary Left/Right choice of the sender; choosing Left and splitting the pool of tokens is trustworthy. To summarize the general content and variation within our messages, Table 1 shows a selection of the most effective and least effective messages — always and never trusted, respectively. We briefly note that there is an immediately apparent qualitative difference between these groups. Effective messages are longer and display a higher level of effort and understanding on average than ineffective messages. For the sake of space, only a partial selection of effective messages is shown; the unabridged table is included in the Appendix.

Table 1: Example Messages: Most and Least Effective

Type	Messages
Most Effective	<p>I'm going LEFT, so please select IN. This will maximize the earnings for both of us equally.  choose IN so that we will both get 12 tokens  Choose "In" and I'll make sure we both get 12 tokens!  Hi there. I would understand why you would choose OUT, but I promise to make the choice ...  I promise you that I will always choose "left" so that we both get 12 points each.  Hello!!! Let's make this work for both of us in the best possible way! Have a great day!  Hey, I'm not a greedy person and in these situations, I always choose the most beneficial option for ...  You can be sure that i will choose left, as it is the best option for both us. Please go ahead and select 'In'  I'm good for the 12. I believe that there are real people here, and I'm completely down with splitting the 24 and ...  I'm picking left so if you want a few extra tokens, go for in. Totally understand if you don't but that's my choice.  I'm always going to choose "left" because I think it's dickish to screw another player over.  I will maximize both our outcomes ie LEFT if you stay in.</p>
Least Effective	<p>I'm going to choose Right  i get 24 you get 0  U GO ON STRITE  Hello Let's play the game  I want choose Player A  hi, how are you doing?  Hey. let's work together, let's both choose out.  I choose out  choose out</p>

In our data, the median message length is 40 characters with a mean of 54.5, so our messages are considerably shorter on average than those seen when similar procedures are conducted in a laboratory setting — median length of 110 and mean of 137.6, (Charness and Dufwenberg, 2006). Figure 2 presents the distribution of length in our sample compared to this laboratory setting

Figure 3: Message Length Distribution



A comparison of the distribution of length in our messages collected through MTurk and messages from a similar, in-person experiment reported in Charness & Dufwenberg (2006).

To connect this study to existing work on communication and to check that our subjects' behavior conforms with those in laboratory settings, we manually label promises with the

help of four research assistants. We label a message a promise if three out of the four labeled a message as such. For a message to be classified as a promise we require not only that the message contain a statement of a future action, but also that it be coherent within the context of the game. Out of the hundreds of unique messages, a limited number of participants promised to do actions which are either impossible, i.e. “I will choose to be player B,” or promised to play Right and take the entire pool for themselves. Although these are clearly promises in a grammatical sense, we do not include them in our classification. Out of 1004 total messages, 389 are promises.

### 3.2 Tokenization and Lemmatization

To work with our text data, subjects’ messages must first be normalized and processed into collections of tokens, each representing an individual word or short series of words in a message. With natural language, this preprocessing step presents significant challenges as we simultaneously want to make the problem tractable by reducing the number of unique words and retain as much of the richness of the original messages as possible. This preprocessing is conducted using the Natural Language Toolkit package for Python, (Loper and Bird, 2002). To do this, we utilize the WordNet lemmatizer, (Fellbaum, 2010) to reduce the various conjugations and declensions into their base forms. This lemmatizer performs best when the part of speech of a word is included, so we also employ Treebank part of speech tagger (Marcus et al., 1993). The end result of these two algorithms is that each message is split into an ordered set of tokens, each token being a word and part of speech pair.

Usually, messages are further reduced by removing stop-words which provide little additional information about the contents of a document. This is a useful step when trying to categorize documents based on their subject matter, where the typical stop-words would be widely shared between the documents in a corpus. However, we omit this step as standard stop-word dictionaries would remove a great deal of information relevant to our specific environment. The messages in our corpus are largely all about the same topic – the game at hand. As such, our variation largely stems from different rhetorical approaches around the same topic. Words like ‘I’, ‘we’, and ‘you’ would all be filtered out by standard stop-word dictionaries, but are clearly important in our environment. After normalizing the tokens to be entirely lowercase, we utilize the WordNet lemmatizer to conduct our normalization. Table 2 presents the number of occurrences of the most frequent tokens in our data.



Table 2: Most Common Tokens

Tokens	Frequency	Tokens	Frequency	Tokens	Frequency
i	751	so	131	'll	72
choose	586	u	131	best	72
in	498	player	120	please	69
left	393	's	116	more	67
you	365	that	114	maximize	66
to	359	can	106	point	65
will	341	make	102	our	59
be	338	out	97	each	56
the	291	hi	97	most	55
both	264	hello	96	game	54
we	253	12	95	all	54
and	244	play	94	b	52
a	235	do	88	always	51
go	176	token	86	money	51
if	175	'm	82	pick	51
get	174	fair	81	with	50
of	153	this	81	have	49
for	152	it	77	right	48
let	137	good	74	earn	46

In addition to the collection of tokens representing each message, we construct a measure of the sentiment of each message. In stock market data, investor sentiment has predictive value of future returns (Tetlock, 2007). Using the OpinionFinder Subjectivity Lexicon (Wilson et al., 2005), we count the number of positive, negative, and total words in each message. This positive and negative word classification requires some minor adjustments to better fit our environment. For example, “right” is identified as a positive word in the sentiment dictionary. As this is an action label in our game, this is removed from the subjectivity dictionary.

## 4 Results

### 4.1 Baseline Approach: Literal Meaning

Our goal in this paper is to examine whether using Natural Language Processing to analyse the content of participant messages yields any insights that would otherwise not have been obvious. To do so, we create a baseline message model which maps the free-form text messages to their literal meanings, based solely on the game’s action space. Before proceeding with the NLP techniques, we first establish the validity of our baseline model by regressing subjects’ choices onto dummy variables representing the possible literal meanings.

Table 3 below summarizes the results of these models.

Table 3: Regressions on baseline messages

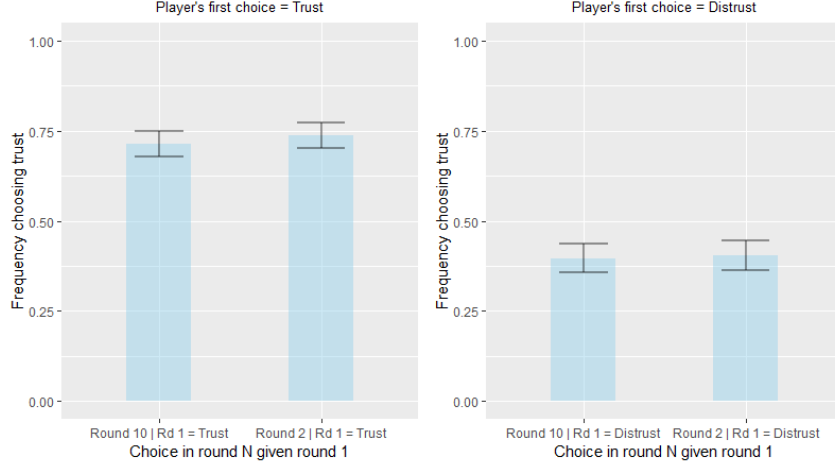
	Trust	Share
LEFT	0.219*** (0.015)	0.225*** (0.039)
IN	0.149*** (0.018)	0.041 (0.046)
BOTH	0.253*** (0.016)	0.186*** (0.041)
CONSTANT	0.462*** (0.009)	0.532*** (0.022)
N	976	969
F-stat	120.175	14.382
Adj. R <sup>2</sup>	0.268	0.040

The results agree with what has been seen in previous studies involving communication. In the trust model, where the dependent variable is the frequency with which recipients of a given message choose the In action, statements of intention to choose Left by the sender, requests that the recipient play In, and combinations of the two have significant positive effects relative to messages which contain none of these features. In the Share model, using the likelihood that the Agent chooses Left after sending a given message as the dependant variable, we find that players who send messages promising to choose Left, sharing the tokens, are significantly more likely to do so. Requesting that the recipient of a message choose In in the absence of a promise has no significant relationship with the sender’s behavior.

## 4.2 Data validity: order effects

Additionally, we check to see if there are significant order effects based on prior messages received. To do this, we divide our subject pool into two groups based on their behavior as Principals. One group includes all subjects who chose In in the first round, with the other including all who chose Out. We then compare the frequency with which each group chose In in round 2 and in round 10. The assumption is that if there are substantial order effects, the In/Out choice in the second round will be significantly correlated with the first round choice, however, the choice in round ten will not be.

Figure 4: Order Effects



95% confidence intervals are plotted in black

As figure 4 shows, we find no evidence of substantial order effects. For both groups, the frequency of subjects choosing In is not significantly different between rounds two and ten. However, we do observe additional evidence of subject-level heterogeneity. Those who chose In (or Out) in the first round tend to do so throughout the duration of the game.

### 4.3 Regressions on Tokens

The action-space model of communication provides us with a reasonable baseline against which we can compare results using Natural Language Processing to develop and examine more complex communication models. Manually mapping each message to the actions to which they most directly correspond allows us to isolate the unique impacts of free-form messages.

In our first NLP model, we augment our literal message model to include the top 50 most commonly used tokens in our message corpus. By using this model, with control variables included for the literal meaning categories in the baseline model, we examine which, if any, of the most commonly used words are associated with higher (or lower) than expected levels of both Trust and Trustworthiness — likelihood that recipients of a message choose In and senders choose Left, respectively.

Table 4 below presents the results of the Trust and Share frequency regressed on the token dummy variables. As this model assumes that each message is an unordered collection of independent tokens, we refer to this as the bag-of-words model. For clarity, we omit tokens that are not significant in either of the two models. The full results are included in the appendix.

Table 4: Bag-of-words model

Tokens	Trust	T-stat	Trustworthiness	T-stat
<b>choose</b>	-0.023	(-1.512)	-0.099*	(-2.259)
<b>in</b>	0.068***	(3.545)	0.036	(0.662)
<b>left</b>	0.060**	(3.278)	0.083	(1.608)
<b>to</b>	0.029*	(2.069)	0.088*	(2.187)
<b>we</b>	0.047*	(2.202)	0.138*	(2.306)
<b>and</b>	0.031*	(1.998)	0.046	(1.051)
<b>go</b>	-0.013	(-0.732)	-0.114*	(-2.326)
<b>for</b>	0.008	(0.398)	0.106*	(1.974)
<b>u</b>	0.048*	(2.007)	0.003	(0.043)
<b>'s</b>	0.001	(0.038)	0.159*	(2.363)
<b>that</b>	0.042*	(2.043)	0.043	(0.742)
<b>can</b>	0.046*	(2.084)	-0.143*	(-2.275)
<b>make</b>	0.066**	(3.015)	0.010	(0.165)
<b>out</b>	-0.101***	(-4.933)	0.087	(1.509)
<b>fair</b>	0.052*	(2.316)	0.113	(1.793)
<b>this</b>	0.063**	(2.730)	0.099	(1.510)
<b>maximize</b>	0.067*	(2.247)	-0.148	(-1.692)
<b>always</b>	0.053*	(2.068)	0.136	(1.869)
<b>right</b>	-0.130***	(-4.614)	-0.224**	(-2.767)

Each coefficient is associated with a dummy variable indicating a message includes the relevant word. Control variables for messages' literal meaning are omitted from the table.

\* :  $p < .05$ , \*\* :  $p < .01$ , \*\*\* :  $p < .001$

Examining the Trust model first, the Tokens with the most significant effects are those which correspond to action labels on the game tree, even after including the literal meaning controls. Each of these tokens effect on the likelihood the recipient of a message chooses In is in the expected direction with "In" and "Left" having positive effects significant at the 1% level and "Out" and "Right" having negative effects significant at the 1% level. In addition, there are several other tokens with significant, positive effects in the Trust model. "Make" and "This" are both significant at the 1% level and positive. The coefficient estimates for the tokens "To", "We", "And", "u", "That", "Can", "Fair", "Maximize", and "Always" are all positive and significant at the 5% level.

With less statistical power from fewer observations, there are fewer significant tokens in the Share model. In the Share model, the only Action label with a significant coefficient estimate is "Right" which has a significant, negative effect at the 1% level. In addition, "choose" "go" and "can" are negatively associated with the likelihood senders choose Left and significant at the 5% level. The coefficient estimates for "To," "We," "For," and " 's," are positive and significant at the 5% level.

This approach shows that there is indeed additional information relevant to predicting

players' choices remaining even after considering the bare meaning of the messages. For many of these tokens, the interpretation of their coefficient estimates is clear. The significance of the Action labels is, perhaps, somewhat unsurprising. However, their significance even after controlling for promises to choose Left and requests to choose In suggests that these sort of statements are most effective when they directly and unambiguously reference the intended or requested action.

For the remaining significant tokens, properly interpreting the results requires us to return to our message corpus to examine the context in which they appear.

Having only one decision associated with each message (instead of ten in the Trust model) we find that there are fewer tokens significantly associated with changes in the likelihood their authors chose Left. Of the action labels present in the game, Right, is the only one with a significant coefficient estimate in the Share model, predicting a much lower likelihood of the sender choosing Left. Although the negative coefficient is coherent with the messages which contain the Right token, it is rather unexpected. Stating an intention to choose Right does not make sense in the usual understanding of the trust game wherein the message senders' attempt to persuade their partners that they are likely to choose Left. Given the incongruity between stating intent to choose Right and the usual strategy of persuasion, it is tempting to write these messages off as mistakes, written by participants who did not understand the game's instructions. However, this would not explain why these Agents follow through on their promises to choose Right. This, instead, suggests that there is a subset of participants who are strongly averse to lying but not necessarily averse to unfair outcomes.

The four tokens with positive coefficient estimates in the Share model fall into two categories. The first of these contains the plural first-person pronouns "We" and " 's" which is used as for contractions containing "us" almost exclusively. The remaining two positive share coefficient tokens are "To" and "For" which are primarily used in messages that detail a reason why the recipient should trust the sender.

Apart from "Right," the remaining significant tokens in the Share choice model do not have readily apparent interpretations in the trust game. Interpreting these tokens, "Go," "Can," and "Choose," requires that we examine them in the context in which they appear. Messages using the "can" token are typically used by Senders to provide a reason why the Receiver should choose In by describing the outcome. A typical message of this type ends with "... so we can get the most tokens." GO messages are used to convey relatively weak promises, using more casual language than other promises. CHOOSE messages, although often used in promises, are more often used as commands for the Receiver to choose In.

Comparing the results in each model, we do not find significant discrepancies that indicate effective deception. Tokens indicative of effective deception would be those that have significant, negative coefficients in the Share model but positive, significant coefficients in the Trust model. Only one token, "Can," fits this definition. Given the large number of regressors in the model, this is indistinguishable from random chance. Therefore, we draw no conclusions from this regression model about potential signs of effective deception.

#### 4.3.1 Regressions in context, separate models for promises/ all others

In both of our previous specifications, the models assume that the information contained in a message can be represented as a linear combination of each of the words present. This

<b>coefficients</b>	<b>Trust</b>	<b>T-stat</b>	<b>Share</b>	<b>T-stat</b>	<b>Trust</b>	<b>T-stat</b>	<b>Share</b>	<b>T-stat</b>
<b>in</b>	0.106***	(4.116)	0.030	(0.426)	0.063***	(3.549)	0.018	(0.370)
<b>left</b>	0.124***	(4.350)	0.178*	(2.255)	0.100***	(5.023)	0.078	(1.383)
<b>will</b>	0.062*	(2.251)	0.108	(1.406)	0.005	(0.221)	0.103	(1.770)
<b>both</b>	-0.071*	(-2.000)	0.094	(0.960)	0.018	(0.699)	-0.006	(-0.081)
<b>we</b>	0.088*	(2.300)	0.127	(1.187)	0.019	(0.691)	0.154*	(2.024)
<b>go</b>	0.063*	(2.115)	0.001	(0.009)	-0.029	(-1.264)	-0.175**	(-2.705)
<b>get</b>	0.069*	(2.204)	-0.025	(-0.287)	0.038	(1.436)	-0.011	(-0.150)
<b>u</b>	0.076*	(1.972)	-0.076	(-0.712)	0.006	(0.198)	0.081	(0.868)
<b>'s</b>	0.036	(0.859)	0.312*	(2.560)	-0.020	(-0.654)	0.129	(1.494)
<b>that</b>	-0.007	(-0.198)	-0.047	(-0.462)	0.056*	(2.114)	0.078	(1.030)
<b>make</b>	0.084	(1.902)	-0.016	(-0.128)	0.080**	(2.936)	0.001	(0.013)
<b>out</b>	-0.084*	(-2.212)	0.093	(0.879)	-0.148***	(-5.827)	0.083	(1.155)
<b>fair</b>	0.015	(0.398)	0.298**	(2.839)	0.080**	(2.679)	-0.022	(-0.259)
<b>'ll</b>	-0.114*	(-2.210)	-0.087	(-0.610)	0.028	(0.977)	0.118	(1.438)
<b>best</b>	0.102*	(2.077)	0.215	(1.563)	0.004	(0.147)	-0.075	(-0.894)
<b>maximize</b>	0.155**	(3.161)	-0.025	(-0.164)	0.037	(0.861)	-0.255*	(-2.119)
<b>our</b>	-0.014	(-0.235)	-0.394*	(-2.131)	-0.004	(-0.095)	0.160	(1.276)
<b>each</b>	-0.157**	(-3.053)	0.118	(0.831)	0.064	(1.690)	-0.050	(-0.468)
<b>always</b>	-0.010	(-0.227)	0.095	(0.787)	0.110**	(3.273)	0.215*	(2.261)
<b>pick</b>	0.098*	(2.277)	-0.029	(-0.240)	0.003	(0.071)	0.055	(0.458)
<b>have</b>	0.094*	(2.075)	0.009	(0.071)	0.043	(0.981)	0.207	(1.630)
<b>right</b>	-0.169***	(-3.474)	-0.220	(-1.628)	-0.134***	(-3.653)	-0.217*	(-2.026)

Table 5: New token regressions controlling for the four message types, modeling promises and non-promises separately

is obviously a strong assumption and although we have demonstrated the usefulness of this approach, it does not capture the full complexity of our subjects' free-form text messages. To examine some of this additional complexity, we again incorporate our promise labels. This time we use them to partition our pool into two groups: promises and all other messages. The results of these regressions are presented in table 5.

This new specification, adding a degree of context to the tokens, highlights previously unseen message features. Of particular interest are "always" and "fair." Both have strong, positive effects when used in non-promise messages. Further, promises which include the "fair" token are much highly significantly correlated with trustworthy promises. Examining the messages, we find that non-promises which use the "always" token nearly universally use it to emphasize their intention to be fair. The effectiveness of these messages in gaining the recipient's trust suggests a previously unknown class of messages, soft promises. Although not coded as promises, as they don't explicitly state an intention to play Left, Senders stating that they are fair has a clear implication in the context of the game.

Further, receivers are suspicious of messages which highlight both players: those which say outcomes which they 'both' or 'each' may get.

## 4.4 Cluster Analysis

What the above analysis shows, in short, is that there are features in our corpus of messages that have strong correlations with both trust and trustworthiness choices even after controlling for literal meaning. With this approach — running OLS regressions on dummy variables for common tokens — we establish an initial result that literal meaning does not capture all the useful information contained in this free-form text data. These significant tokens suggest new classes of messages worth further investigation. However, this method of treating each message as linear combinations of independent words leaves out much of the complexity present in natural language.

To address this shortcoming, we employ cluster analysis to identify groups of similar messages which we may then aggregate and examine further. Cluster analysis produces groups of messages such that the within-group similarity is greater than the inter-group similarity. The specifics of this task vary greatly based on the chosen measure of similarity and objective function comparing intra- and inter-group similarity. This permits us a great deal of freedom to choose a similarity measure which captures more of the complexity of natural language than in our previous approach.

Previous authors have employed cluster analysis on laboratory data including subject choices (Ellingsen et al., 2018; Grimm et al., 2021; Woods et al., 2022) and eye tracking data (Polonio and Coricelli, 2019; Devetag et al., 2016).

Our approach involves two steps. First, we employ Affinity Propagation, (Frey and Dueck, 2007) to perform our cluster analysis and group our messages. Affinity Propagation allows for both the number of clusters and the number of elements assigned to each to vary, which we require in the absence of *ex ante* knowledge of the shape of the message space. Flexible cluster size also allows nonsequitur messages to separate themselves into smaller clusters instead of forcing themselves into unsuitable clusters to achieve an arbitrary number of groups. Each cluster also forms around an exemplar message, giving us a natural starting point when we manually interpret and label each cluster.

Following this aggregation, we then employ standard OLS, fitting a linear probability model to test which areas of our message space have significantly higher rates of trust (greater likelihood Principals choose In) and trustworthiness (greater likelihood Agents choose Left). As with our bag-of-words token regressions, we add a manually coded control variable for literal meanings.

Conducting cluster analysis with affinity propagation requires us to define a similarity matrix  $\mathbf{S}$  for the preference of each element for potential exemplars. Given this degree of flexibility, we conduct cluster analysis with two different similarity measures, one using solely the information in our corpus of text (corpus-only approach) and the other informed by outside natural language knowledge (knowledge-based approach). By using multiple similarity measure, we test the robustness of our conclusions against changes in our clustering procedures. If cluster analysis is a useful tool for us to process free-form laboratory text data, we expect the results to be relatively consistent across these choices.

In the Corpus-only approach, we define our similarity matrix  $\mathbf{S}_J$  by taking each message, again, as an unordered collection of tokens and calculating the Jaccard index of overlap between sets. For messages  $m_i$  and  $m_k$ , both being sets of tokens:

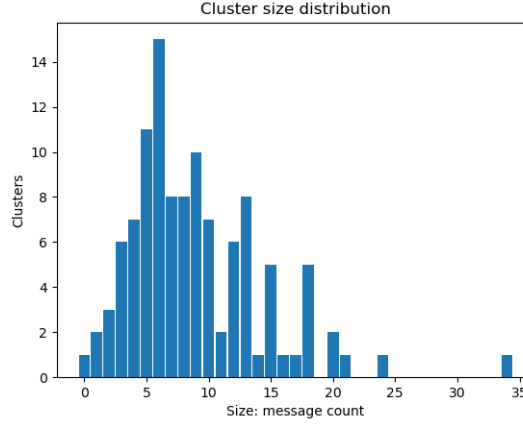
$$s(i, k)_{i \neq k} = j(m_i, m_k) = \frac{|m_i \cap m_k|}{|m_i \cup m_k|}; s(k, k) = \text{Median}(\mathbf{J})$$

Applying this measure to each pair in our corpus yields a complete graph in matrix form; each element  $j(i, j)$  representing the similarity of message  $i$  to message  $j$  and vice versa; a higher value implying greater similarity.

## 4.5 Results: Affinity Propagation using corpus-only similarity measure

Affinity Propagation identifies 112 clusters of similar messages in our corpus. The median cluster size is 8 messages with a mean of 9.1. Figure 6 presents the distribution of the number of messages assigned to each cluster.

Figure 5: Distribution of cluster sizes



We conduct cluster analysis with the goal of identifying key features in our corpus, without a priori knowledge of the number of expected clusters or their contents. Instead, we must investigate each cluster and assess whether or not the clusters share a common meaning to a sufficient degree. We must also manually assign a label to each cluster, describing the common theme shared among the messages assigned to it, using the exemplar element for each cluster as our starting point. In Table 6 below, we present two example clusters and their exemplars, demonstrating the similarity between the messages grouped by affinity propagation.



Table 6: Example Clusters

Cluster	Message, ( <b>exemplar</b> )
55	<b>I believe in fair play</b>
55	I believe in equality
55	I am playing in
55	I believe in playing fair and having an equal outcome.
55	I believe in fair, even-steven financial deals. I like a win-win situation. ...
55	I believe in being fair with people.
58	<b>Hey, I want to make sure we both get the highest amount of bonus!...</b>
58	Hopefully we both want to get the highest points possible
58	I will make sure we get the same amount of tokens
58	Hello. I hope we both make a good amount of bonus money today! Good luck!
58	Hi there! I always split these 50-50 because I would want someone to do the ...
58	Hello. Lets make this to where we both can get the largest bonuses out of th...
58	To make sure you guys get a higher bonus, please choose IN in each step so w...

In both example clusters, the elements are sufficiently similar to lead to a clear substantive interpretation. Cluster 55 messages state that the sender believes in fairness, which is captured well by the exemplar message, “I believe in fair play.” Messages in cluster 58 state the author’s hopes for the game outcome, speaking in terms of the amount of tokens. Again, the exemplar, “Hey, I want to make sure we both get the highest amount of bonus...” represents the typical message in the cluster well.

#### 4.5.1 Regressing Trust choices on clusters

With these cluster assignments, we again fit two Ordinary Least Squares models, regressing the frequency of In choices by the receiver (Trust) and the frequency of Left choices by the author (Share) on dummy variables for each message’s assigned cluster. As before, we include literal meaning as a control. For the sake of clarity, we present only clusters with the most statistically significant coefficient estimates.

We first examine the most effective messages in gaining the recipient’s trust. For each cluster, we include both the exemplar message and a brief description of the contents of the common themes. The full list of messages in each of the clusters we discuss is in the appendix.

Table 7: Most Trusted Corpus-only Clusters

cluster	message	Cluster size	trust (t-stat)	share_coef (t-stat)
78	<b>Requests In with 'please', promises Left, and specifies '12 token' outcome</b> Please choose "IN" so that we can both get 12 points. I will choose "LEFT."		0.2*** (4.19)	-0.02 (-0.17)
84	<b>Promises 'what's best for everyone'</b> Hi, I will make the choice that is beneficial for all of us :)		0.18*** (3.89)	0.01 (0.1)
96	<b>Requests In so both can 'maximize'</b> If you choose In we can both maximize our tokens.		0.19*** (3.29)	-0.16 (-0.98)
40	<b>"Let's try and be fair"</b> lets try and be fair		0.22** (3.20)	0.05 (0.25)
37	<b>Promises Left, specifies outcome for 'both'</b> I'm going to choose Left, please choose In so we both get 12 tokens		0.11** (2.88)	-0.02 (-0.19)
58	<b>Suggests that both 'get the highest bonus'</b> Hey, I want to make sure we both get the highest amount of bonus! Hopefully y...		0.17** (2.73)	-0.02 (-0.13)
44	<b>Requests In, promises Left, and specifies '12 tokens'</b> If you choose In I will choose Left and we will both get 12 tokens		0.08** (2.54)	0.1 (1.11)
45	<b>Request In, promise Left, and specifies outcome for 'both of us'</b> Hello, I will chose left if you decide to go in. I want this to be a win/win ...		0.10** (2.36)	0.16 (1.28)
81	<b>Promises 'most for both of us'</b> Trying to make the most for both of us		0.11** (2.29)	0.32** (2.36)
55	<b>"I believe in fair play"</b> I believe in fair play		0.15* (2.22)	0.46** (2.38)

The top ten most significant, positive coefficients when regressing Trust on cluster dummies, controlling for promise labels, along with the corresponding coefficient estimates in our Share model. For each cluster, we include both the exemplar and a brief description of the common themes among the messages \* :  $p < .05$ , \*\* :  $p < .01$ , \*\*\* :  $p < .001$

Again, we are cautious of over-inference, and avoid delving too deeply into any individual cluster. Instead, we focus on the common features across these groups of highly effective messages. The most effective messages largely fall into two broad categories. First, four of the top ten clusters, (numbers 78, 37, 45, and 44,) follow a common three part structure — explicitly requesting the receiver plays In, promising to choose Left, and stating the positive result should both these choices be made. These messages conform with our definition of a

promise, but go a step further in describing the outcome and requesting specific action from the Principal. An additional two clusters, (96 and 58) do not follow this three-part formula but are also detailed promises.

Second, four of the ten clusters are ‘soft promises,’ (clusters 55, 40, 81, 84). These messages fail to meet our definition of a promise, requiring the author to make a clear statement of a future intended action. Instead, they make softer claims: “Let’s try to be fair,” for example. Although these are not strictly promises in the grammatical sense, the implication is clear in the context of the game. Recipients of these messages interpret them as statements of intent to choose Left even without an explicit action indicated.

None of the increased rates of trust in these clusters have a corresponding significant increase in the author’s trustworthiness. However, we cannot reliably interpret this as evidence of effective deception. Our design focuses on the In/Out decision and thus has lower statistical power to estimate Left/Right coefficients; each message has ten readers but only one author.

Examining the least effective clusters, (those with the most significant negative coefficients in the Trust regression,) we can again identify two broad themes. As with the most effective clusters, we also present each cluster’s coefficient in the Share regression, modeling the agents’ choices given the message they sent.

Table 8: Least Trusted Corpus-only Clusters

cluster	message	cluster_size	trust (t-stat)	share_coef (t-stat)
97	choose out	6	-0.29*** (-4.27)	0.03 (0.17)
32	choose right!	10	-0.21*** (-4.01)	-0.35* (-2.32)
41	I chose out.	12	-0.18*** (-3.5)	-0.09 (-0.65)
34	hello..let’s play	10	-0.14** (-2.55)	0.12 (0.75)
102	I AM PLAYER A	18	-0.09* (-2.08)	-0.19’ (-1.66)
111	yes goodsss	4	-0.16* (-1.96)	-0.05 (-0.21)
92	OUT PLAYER 8	5	-0.14’ (-1.93)	0.25 (1.20)

The top ten most significant, negative coefficients when regressing Trust on cluster dummies, controlling for promise labels, along with the corresponding coefficient estimates in our Share model. For each cluster, we include both the exemplar and a brief description of the common themes among the messages \* :  $p < .05$ , \*\* :  $p < .01$ , \*\*\* :  $p < .001$

As suggested by the token regressions, Principals distrust messages that seem to go against the Agents’ best interests (clusters 32, 97, and 41). There is no immediately apparent

explanation for sending a message stating your intent to play Right or asking the reader to choose Out, although these messages do not violate the rules of the game. Also in this group are messages demonstrating apparent confusion about the game (clusters 41, 102, 31, 92, and 56).

The second theme among the most distrusted messages is newly revealed by our Cluster Analysis; brief messages containing greetings but no other content (cluster 34) have highly significant, negative coefficients even after controlling for their lack of a promise to play Left.

#### **4.5.2 Regressing Share choices on clusters**

Moving to the Agents' Left/Right choices, we fit an OLS model, using dummy variables for each message's cluster assignment and human-coded promise label as independent variables to estimate the likelihood that its author chooses to share the tokens (choose Left in step 2). As a consequence of our experiment's design, this model has significantly less statistical power than our Trust model; there are ten Trust choices for each message but only one associated Share choice. We focus primarily on relating the Share coefficient estimates to their Trust model counterparts, asking whether Principals' choices are rational responses to accurate predictions about Agents' behavior. Table 9 presents clusters with the most significant positive and most significant negative Share model coefficient estimates.

Table 9: Significant Corpus-only Clusters: Share

cluster	message	cluster_size	share t-stat	trust t-stat
79	<b>"I like this game"</b> I like this game.	8	0.45** (2.71)	0.06 (0.98)
35	<b>Suggestion to 'work together'</b> Hiya. I hope we can work together.	11	0.37* (2.55)	0.09' (1.71)
108	<b>Request In, promise Left, specify '12 tokens'</b> Lets opt in and choose left so we both get 12 points.	7	0.44* (2.44)	0.01 (0.09)
55	<b>"I believe in fair play"</b> I believe in fair play	6	0.46* (2.38)	0.15* (2.22)
81	<b>Promises 'most money for both of us'</b> Trying to make the most for both of us	12	0.32* (2.36)	0.11 (2.29*)
89	<b>Promises left conditional on In</b> If you choose in, I will choose left.	21	0.23* (2.06)	0.07' (1.70)
9	<b>"Good player"</b> GOOD PLAYER	8	-0.30 (-1.78)	-0.03 (-0.44)
29	<b>"I think you should choose In"</b> I think you should choose In	8	-0.32 (-1.89)	-0.09 (-1.52)
61	<b>"Going left"</b> Going left	6	-0.45* (-2.29)	-0.02 (-0.34)
32	<b>"choose right!"</b> choose right!	10	-0.35* (-2.32)	-0.21*** (-4.01)
52	<b>Appeal to 'maximization'</b> let's maximize OUR bonuses	9	-0.45* (-2.82)	0.02 (0.32)

The top ten most significant, positive coefficients when regressing Share on cluster dummies, controlling for promise labels. \* :  $p < .05$ , \*\* :  $p < .01$ , \*\*\* :  $p < .001$

The clusters with the most significant coefficient estimates in the Share model largely follow the patterns seen in the Trust model. Detailed promises (clusters 89 108, and 81) and soft promises (clusters 55 and 35) are both included in the positive Share clusters.

Similar to the Trust model, the messages in these least trustworthy clusters are almost all noticeably lower quality than the typical message in our corpus. This includes both nonsensical messages like those in cluster 32 and the briefest possible promises in cluster 61. The cluster with the strongest negative coefficient estimate is also the exception to this pattern of poor quality, cluster 52, which contains explicit appeals to payoff maximization. Unlike the low-quality message clusters, there is no corresponding significant negative coefficient in our Trust model. Recipients of these appeals to maximization continue to choose In, entrusting the tokens to their partner, at the same rate as other non-promise messages.

Considering both the Trust and Share models formed from our corpus-only clusters, we find five broad classes of messages which warrant further study. Our primary goal is to study the determinants of trust and we find two themes associated with increased trust and two

associated with decreased likelihood of trust. Both detailed promises — explicitly promising an action, giving their partner a call to action, and detailing the promised outcome — and soft promises, (e.g. "I think fairness is important,") are associated with increased trust. These messages tend to have corresponding positive coefficient estimates in our Share model, however, we cannot make strong claims about this agreement due to the Share model's relative lack of statistical power. In contrast, Principals trust poor-quality messages and brief greetings significantly less often. Our coefficient estimates in the Share model generally support their skepticism, with estimates agreeing in direction if not also in significance.

The final notable cluster, appeals to maximization, provides our only evidence of possible effective deception. Despite our Share model estimates indicating that these messages are the least likely to be trustworthy, we find no significant effect in our Trust model, despite it having much greater statistical power.

## 4.6 Affinity Propagation using Knowledge-based similarity

In the preceding section, we show the potential of cluster analysis as a tool for making sense of unstructured text data. From over 1000 messages, Affinity Propagation successfully aggregates individual messages into clusters which have substantial within-cluster similarity. These clusters allow us to conduct careful, limited inference and identify previously unknown message features which impact their recipients' choices.

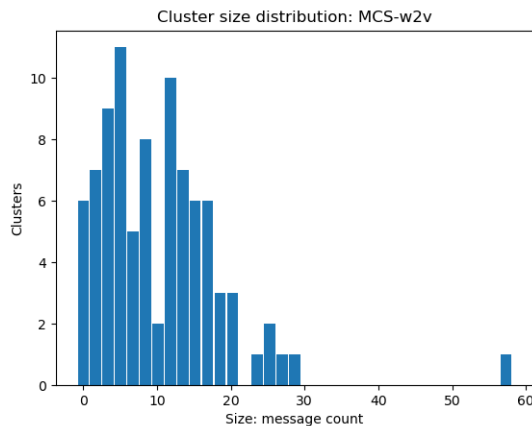
However, these results are not without caveats. The clusters themselves are sensitive to our choice of similarity measure and other tuning parameters in the affinity propagation model. The Jaccard index as a similarity measure is both easily interpretable and based solely on the text contained in our message corpus. This simplicity makes it an attractive baseline similarity measure but, in choosing this measure, we ignore the large body of existing natural language processing work. To illustrate this weakness, consider the brief greeting clusters, 34 and 4 in the above model, (exemplar messages "Hello, let's play" and "Hi, how are you?" respectively.) Although the within-cluster similarity is satisfactory, the division between the two is based solely on the use of "Hi" versus "Hello" to convey nearly identical information. We make no claim about this specific partition being the definitive description of the message space, but this over-clustering raises concerns about whether our results are robust to our specification of similarity measure.

We test the robustness of our corpus-only Affinity Propagation results by conducting the same clustering and regression methods but using an alternative similarity. This Knowledge-based approach is trained on existing information outside of our corpus of messages. In this second approach, we measure the similarity between two messages using an approach first introduced by (Mihalcea et al., 2006) to compare short text documents. For two sets of words  $T_1$  and  $T_2$ , the MCS similarity is:

$$S_{mcs} = \frac{\sum_{w \in T_1} \text{maxsim}(w, T_2) \cdot \text{idf}(w)}{2 \sum_{w \in T_1} \text{idf}(w)} + \frac{\sum_{w \in T_2} \text{maxsim}(w, T_1) \cdot \text{idf}(w)}{2 \sum_{w \in T_2} \text{idf}(w)}$$

$\text{maxsim}(w, T)$  is the maximum semantic similarity between a word  $w$  and any word in the sentence  $T$ , the precise details of which vary depending on the choice of semantic similarity measure. We measure word-to-word semantic similarity with the Google word2vec neural network, (Mikolov et al., 2013). After training, word2vec take individual words as inputs

Figure 6



and outputs a multi-dimensional vector representation reduction of latent meaning. The network we implement here is pre-trained on a corpus of articles available through Google News, using the Gensim python package (Rehurek and Sojka, 2011). The resultant neural network encodes individual words into a 300-dimension vector; the similarity between any two words is the cosine distance between their word2vec vector forms.

$\text{idf}(w)$  is the inverse document frequency for word  $w$ ,  $\text{idf}(w) = \ln(\frac{N}{n(w)})$  where  $N$  is the total number of messages in our corpus and  $n(w)$  is the number of messages containing  $w$  (Sparck Jones, 1972). A higher idf value for a word implies greater specificity and thus receives greater emphasis when calculating similarity.

After calculating the similarity matrix using the MCS and w2v, we again apply affinity propagation, resulting in 89 unique clusters. The mean size is 11.3 messages with a median size of 10. Figure 7 shows the distribution of cluster sizes.

Table 10 presents two example clusters formed by Affinity Propagation using the MCS-w2v measure of message similarity. As in the corpus-only approach, there is a high degree of within-cluster similarity around a common theme. For the two selected examples, cluster 3 and 59, we label them as greetings and emphatic promises to choose Left respectively. As before, the exemplar elements serve as reliable summaries of each cluster’s contents. We also observe that the greater flexibility afforded by the Knowledge-based similarity measure achieves the desired result of softening the boundaries between clusters while maintaining within-group similarity. Cluster 3 now includes short greetings regardless of whether they use "Hi" or "Hello."

Table 10: Example word2vec Clusters

Cluster	Message, ( <b>exemplar</b> )
3	<b>Hi, how are you?</b>
3	a bit confussed but attemp either way. hi
3	HELLO! HOW ARE YOU?
3	Are you ready?
3	Hello, what are you going to choose?
3	hi how are you
3	HI HOW ARE YOU
3	hi hoe are you
3	hi, how are you?
3	hello how are you doing
3	Hello, how are you doing?
59	<b>I will always choose LEFT.</b>
59	I will always choose Left. I want to split the bonus in a fair way.
59	Hello. Let's try to maximize our earnings. I will always choose LEFT as I be...
59	I will always choose LEFT. That means we are both assured of getting a fair ...
59	I always split evenly with my fellow turkers.
59	I will always play Left if you play In.
59	I will choose Left. I always try to split earnings with other workers when I...
59	I'll always choose Left.
59	I promise to always play left so we all benefit.

With the new, reduced-size clusters, we again fit Ordinary Least-Squares models, regressing both Trust and Share choices on dummy variables for each message's cluster assignment, controlling for literal meaning. Our goal is to assess the robustness of our previous approach and examine whether the same communication strategies are correlated with changes in behavior across the two partitions of the message space. Table 11 and table 12 present the most significant positive and most significant negative coefficient estimates in the Trust model.



Table 11: Most Trusted word2vec clusters

cluster	Cluster Description exemplar message	size	trust t-stat	share t-stat
31	<b>State intent to do 'best for both of us'</b> I'll do the best for both of us!		0.12*** (3.33)	0.12 (1.30)
43	<b>"Choose In to get 12 tokens"</b> choose IN so that we will both get 12 tokens		0.11** (2.49)	0.05 (0.42)
39	<b>Propose 'maximizing both' player's earnings</b> If you choose IN, I will choose LEFT to maximize both of our earnings		0.11** (2.40)	0.00 (0.02)
59	<b>Promises to choose Left, including 'always'</b> I will always choose LEFT.		0.11** (2.34)	0.25* (2.01)
66	<b>"Let's work together"</b> Lets work together		0.10** (2.16)	0.13 (1.08)
64	<b>Suggest 'both earning as much as possible</b> Hey! Let's try and make as much as we can for everyone!		0.17** (2.14)	0.07 (0.32)
15	<b>"I believe in equality"</b> I believe in equality		0.18** (2.04)	0.38 (1.62)
68	<b>Promise to 'always' choose Left</b> im going left		0.10** (2.00)	0.08 (0.64)

Significant coefficient estimates when regressing In/Out choices on cluster dummies using a knowledge-based measure of similarity, controlling for literal meaning. \*:  $p < .05$ , \*\*:  $p < .01$ , \*\*\*:  $p < .001$

Table 12: Least trusted word2vec clusters

clusterw2v	exemplar message	clustersize	trust (t-stat)	sharecoef )t-stat)
80	I choose out	16	-0.21*** (-4.72)	0.03 (0.23)
27	choose right!	15	-0.20*** (-4.46)	-0.39** (-3.05)
61	it really good and nice	6	-0.22** (-3.13)	-0.20 (-1.03)
42	out player a-8 player -b 8	7	-0.17** (-2.66)	0.32' (1.82')
85	I AM PLAYER A	8	-0.16** (-2.62)	-0.16 (-0.94)
3	Hi, how are you?	21	-0.09** (-2.34)	0.06 (0.56)
2	I am choosing out.	6	-0.17* (-2.16)	-0.13 (-0.63)
74	Please choose IN.	27	-0.07 (-1.88)	0.07 (0.67)

Significant negative coefficient estimates when regressing In/Out choices on cluster dummies using a knowledge-based measure of similarity, controlling for literal meaning. \* :  $p < .05$ , \*\* :  $p < .01$ , \*\*\* :  $p < .001$

Compared to the corpus-only model, we find fewer clusters with significant, positive coefficients the Knowledge-based Trust model. However, the same message features stand out as the most effective. The clusters significant at the 5% level follow the same two patterns we observed in the previous model. We see detailed promises represented in clusters 43, 68 and 59 and soft promises in cluster 31 and 15. Additionally, we see clusters featuring suggestions to cooperate without explicitly promising an action or requesting one from its recipient (clusters 64 and 39).

The least trusted messages in the knowledge-based model are effectively identical to those found in the corpus-based model. The bulk of these clusters contain messages that suggest the author does not understand the game. They state intent to either choose actions that are impossible within the bound of the game, e.g. "I choose player A," or state their intention to choose Right. Additionally, we again see that recipients of short, simple greetings are significantly more likely to distrust their sender, even relative to other non-promise messages.

In the Share model (Table 13) we see once again that the broad themes from our model using a knowledge-based measure of similarity match with the previous, corpus-only approach. The clusters associated with higher likelihoods of their senders choosing Left are detailed promises (clusters 73, 52, 59) and implied promises, specifying an outcome should the recipient choose In (cluster 69).

The clusters with significant, negative coefficients also support our conclusions from the Corpus-only approach. Agents who state their intention to play Right tend, strongly, to do so and recipients respond appropriately to these messages by choosing Out. Our only

significant evidence of effective deception, explicit appeals to payoff maximization, is also robust to the alternative measure of similarity. The authors of these messages choose Right 34 percentage points more often, however, there is no significant change in their recipients' In/Out choices.

Table 13: word2vec clusters: Share

cluster_w2v	Description exemplar message	cluster_size	share (t-stat)	trust (t-stat)
73	<b>Promises Left conditional on In &amp; Specifies the outcome</b> If you choose in, I will choose left.	59	0.20** (2.77)	0.03 (1.18)
52	<b>Promises Left conditional on In, no stated outcome</b> If you go IN I'll go left.	20	0.24* (2.16)	0.01 (0.28)
59	<b>Promise to choose left, including 'always'</b> I will always choose LEFT.	16	0.25* (2.01)	0.11* (2.34)
69	<b>Requests In and specifies the outcome</b> Please choose "IN" so that we can both get 12 ...	18	0.21' (1.86)	0.06 (1.49)
32	<b>Suggestion to 'maximize'</b> lets maximize	14	-0.34** (-2.65)	0.09' (1.89)
27	<b>Promise to choose Right</b> choose right!	15	-0.39** (-3.05)	-0.20*** (-4.46)

The most significant coefficient estimates when regressing Left/Right choices on cluster dummies using a Knowledge-based measure of similarity, controlling for literal meaning. \* :  $p < .05$ , \*\* :  $p < .01$ , \*\*\* :  $p < .001$

Overall, the conclusions from our knowledge-based approach agree with the previous, Corpus-only method. Although the resulting partitions are, indeed, sensitive to the choice of similarity measure, the same message features stand out in both the Trust and Share models. Principals who receive emphatic, detailed promises are significantly more likely to choose In, relative to the typical promise. Similarly, soft promises, which imply an intention to choose Left without an explicit statement of intent, yield greater trust than other non-promise messages.

With less statistical power, the Share model provides only limited insights into Agent behavior, with three effect estimates significant at the 5% level. This includes both promises to choose Right and appeals to payoff maximization, both identified in the Corpus-only model. Despite the strong, negative coefficient estimate in the Share model, Principals are no less likely to trust an appeal to payoff maximization than other non-promise messages.

## 5 Reducing Cluster Size

Given the size of our corpus, containing a little over 1000 unique messages, having nearly one cluster for every ten messages raises concerns of over-clustering. We again reiterate that our purpose requires only that the within cluster similarity is sufficiently high. Our purpose is to identify new, relevant message features, not to determine a definitive partitioning of the message space which captures all of the potential complexity and nuance of written communication. Still, with these cluster sizes being implications of largely arbitrary rule-of-thumb parameter choices, we want to examine whether our results are robust to different tuning parameter decisions.

To address this concern, we conduct the same cluster analysis followed by OLS regression procedure as above except that we adjust our tuning parameter, a message’s affinity towards itself, to reduce the number of resulting clusters. Specifically, we select this parameter with the objective of maximizing the BIC of the resulting linear model independently for both the corpus-only and word2vec similarity measures.

This optimization reduces the number of clusters in the corpus-only and word2vec model to 27 and 16 down from 112 and 89 respectively. As before, using cluster assignments as dummy variables, we fit a linear probability model, estimating the frequency recipients choose Left, trusting the sender, and the frequency with which senders choose Left, sharing the tokens if given the opportunity. We fit this model both with and without controlling for promise labels. Figures 14 and 15 present the coefficient estimates for these models.

### 5.1 Reduced Corpus-only clusters

Out of the 27 clusters grouped using only information found in our corpus, 7 clusters (6, 10, 11, 18, & 20) have positive, significant coefficient estimates, indicating that they are trusted at a higher frequency relative to the mean across all messages in our data-set. We again see the two common themes from the previous cluster regressions. Clusters 10, 11, and 18 include detailed promises, requesting In and specifying an outcome in addition to the stated intention to choose Left. The remaining two — clusters 6 and 20 — contain suggestions to work together, implying both a request for In and an intent to play Left without explicitly stating this.

In contrast, three clusters are associated with lower frequencies of trust among their recipients (clusters 1, 3, 9). These messages fall into two distinct themes of common content. The first group consists of errors — messages which are either entirely nonsensical or which are not in line with a basic understanding of the game structure (clusters 3 and 9). The remaining cluster (cluster 1) contains brief greetings with no other content. Both brief greetings and nonsensical messages are also found in the previous cluster specifications.

In the reduced corpus-only Share model, we largely find agreement between the two models, with the signs of significant coefficient estimates in the Share model agreeing, at least in direction with those in the Trust models. There are no clusters which have coefficient estimates significantly different from zero in the Share model with significant, positive estimates in the Trust model.

In the Share model using the reduced, corpus-only clusters, four have significant, positive coefficient estimates. Clusters 6 and 20 both consist of suggestions to ‘work together’ or

Table 14: Reduced size corpus-only clusters

cluster	size	Description exemplar message	size	trust t-stat	share t-stat
1	21	<b>"Hi, how are you?"</b> Hi, how are you?		-0.08* (-2.08)	-0.05 (-0.5)
3	18	<b>"I choose right"</b> I choose right.		-0.17*** (-4.36)	-0.29** (-2.62)
5	22	<b>Good luck / good morning</b> good luck!		-0.04 (-1.03)	-0.24* (-2.41)
6	27	<b>Suggestions to 'work together'</b> I hope we can work together to win.		0.1** (3.00)	0.22* (2.41)
8	45	<b>Brief request for In and promise of Left</b> I will choose in and go left		0.04' (1.68)	-0.02 (-0.24)
9	25	<b>"I choose out"</b> I chose out.		-0.19*** (-5.4)	0.05 (0.51)
10	44	<b>Request In, promise Left, and specify the 'maximum' or 'best' outcome</b> If you choose IN, I will choose LEFT to maximize both of our earnings.		0.06* (2.20)	0.06 (0.74)
11	103	<b>Request In, promise Left, and specify '12 tokens' as the outcome</b> If you choose In I will choose Left and we will both get 12 tokens		0.05* (2.48)	0.08 (1.37)
17	83	<b>Promise Left conditional on In</b> I will be choosing left if you choose in		0.01 (0.55)	0.17** (2.98)
18	64	Please choose "IN" so that we can both get 12 points. I will choose "LEFT."		0.07** (2.82)	-0.09 (-1.47)
19	18	<b>"I like this game"</b> I like this game.		0.01 ( 0.20)	0.35** (3.05)
20	42	<b>States intent to make the 'most for both of us'</b> Trying to make the most for both of us		0.07* (2.49)	0.26*** (3.59)
22	24	<b>Suggests 'best payoff' without a request or promise</b> ALL THE BEST		0.0 (0.06)	-0.22* (-2.24)
24	34	<b>"I am Player A"</b> I AM PLAYER A		-0.04 (-1.27)	-0.17* (-2.09)

Coefficient estimates when regressing Left/Right and In/Out choices on cluster dummies using the corpus-only measure of similarity, controlling for literal meaning. Only clusters with significant coefficient estimates in either model are shown for brevity. \*:  $p < .05$ , \*\*:  $p < .01$ , \*\*\*:  $p < .001$

statements of intent to do so. Cluster 17 includes promises to choose Left, conditional on the Receiver’s choice to choose In. This conforms with the previous results showing detailed promises outperforming simple ones. Finally, messages in which the Sender states that they “like this game” are significantly more likely to come from trustworthy senders.

Negative, significant clusters in the reduced, corpus-only share model also agree with the previous models’ results. Nonsensical messages, promising to choose Right or to choose Player A are significantly, negatively associated with their senders choosing Left. Senders of messages in cluster 22 state their intent to ‘maximize’ their result and are significantly more likely to choose Right.

## 5.2 Reduced word2vec clusters

Table 15: Reduced word2vec clusters

cluster	size	Description exemplar message	trust (t-stat)	share (t-stat)
0	43	<b>Request for In including 'please'</b> Please Choose In	-0.02 (-0.88)	0.07 (0.95)
1	46	<b>"I choose Player A"</b> I choose player A	-0.10*** (-3.83)	0.09 (1.25)
2	140	<b>Requests In, promises Left, and specifies the 'best' outcome</b> I will choose left if you choose IN	0.03' (1.68)	0.05 (1.09)
3	26	<b>Promises Right</b> choose right	-0.13*** (-4.06)	-0.36*** (-3.92)
4	115	<b>Promises Left only</b> i will choose left	0.03 (1.60)	-0.01 (-0.16)
5	113	<b>States intent for what's 'best for everyone'</b> Hello, I want to do what's best for the both of us so I'm ...	0.05** (2.76)	0.00 (0.06)
6	34	<b>Greeting including 'good luck'</b> good luck!	-0.02 (-0.63)	-0.03 (-0.41)
7	69	<b>Greeting including 'let's play'</b> Hello Let's play the game	-0.04' (-1.68)	-0.03 (-0.48)
8	70	<b>"Let's make the most money"</b> lets go in	0.03 (1.27)	-0.09 (1.51)
9	69	<b>Requests In, promises Left, and specifies '12 tokens' as the outcome</b> If you choose In I will choose Left and we will both get 1...	0.05* (2.19)	0.07 (1.12)
10	80	<b>Request for In only</b> Choose The In	0.02 (0.81)	-0.01 (-0.12)
11	29	<b>Suggestions to 'work together'</b> Lets work together	0.09** (2.93)	0.07 (0.81)
12	64	<b>Promises Left only</b> im going left	0.02 (1.00)	-0.1 (-1.61)
13	38	<b>Brief greetings</b> HELLO! HOW ARE YOU?	-0.05' (-1.7)	0.02 (0.28)
14	37	<b>Promises Out</b> I choose out	-0.15*** (-5.47)	0.05 (0.64)
15	37	<b>States intent to be 'fair'</b> I will be fair	0.06* (2.24)	0.14' (1.78)

Coefficient estimates when regressing Left/Right and In/Out choices on cluster dummies using the word2vec measure of similarity, controlling for literal meaning. Only clusters with significant coefficient estimates in either model are shown for brevity. \*:  $p < .05$ , \*\*:  $p < .01$ , \*\*\*:  $p < .001$

This alternative clustering model, using the word2vec similarity measure, yields fewer clusters compared to the corpus-only approach. As intended, this approach mitigates some of the potential over-clustering seen in the first approach and we have fewer groups of closely related messages divided into multiple categories. After this reduction, the same general themes remain common among the messages in clusters with significant coefficient estimates, both positive and negative, in the Trust model. Of the 16 clusters, 4 have significant positive coefficient estimates, controlling for literal meaning (clusters 5, 9, 11, and 15). These clusters capture the same patterns observed in the models fitted to previous similarity measures. Clusters 4 and 15 include 'soft promises' in which the Sender implies that they will choose Left without explicitly stating intended actions. Similarly, messages in cluster 11 suggest that the players 'work together.' This is also a soft promise, but also implies a request that the Receiver choose In as well. Finally, cluster 9 includes detailed promises, with requests for In and a specified outcome of '12 tokens.'

The clusters with significant, negative coefficient estimates in the reduced word2vec trust model are again familiar. All significant, negative clusters (clusters 1, 3, and 14) are either incoherent (cluster 1 messages promise to "choose Player A") or make little strategic sense in the context of the game. Clusters 3 and 14 promise to choose Right or request the Receiver chooses Out, respectively.

In the Share model, there is one cluster with a significant coefficient estimate, positive or negative. Messages in cluster 3 are promises to choose Right. Consistent with our previous results, Senders of these messages keep their word and choose Right at higher frequencies.

### 5.3 Sentiment

Using the bag-of-words and affinity propagation approaches allow us to gain valuable insights to the message topics and their association with trust behaviors. Although this has already allowed us to discover new potentially impactful features of persuasive messages, there remains a great deal of nuance left out by this reduction. To capture a portion of this additional nuance, we conduct further analysis into our messages by examining their sentiment.

To do this, we re-code our set of tokens into positive, negative, and neutral words based on the classifications provided by the OpinionFinder Subjectivity Lexicon (Wilson et al., 2005). As with previous steps, we must alter this procedure slightly to better fit our environment, removing 'right' from the list of words with positive sentiment as it has a neutral meaning as a potential action in our game; all usages of this word in our data refer to the game action. From this, we construct three variables for each message: a count of positive words, negative words, and a control for total word count. As before, we use these variables to fit three models with Trust, Trustworthiness, and Promise labels as our dependent variables. Additionally, we fit two additional models for Trust and Trustworthiness with promise labels as an additional control. Table 6 presents the results of these regressions.



Table 16: Regressions on Sentiment

	Trust	Truth	Promise	Trust w Control	Truth w Control
positive words	0.014* (0.006)	0.019 (0.015)	0.033* (0.014)	0.008 (0.006)	0.013 (0.014)
negative words	-0.032* (0.013)	-0.052 (0.031)	-0.106*** (0.030)	-0.014 (0.012)	-0.033 (0.031)
total words	0.006*** (0.001)	0.005* (0.002)	0.013*** (0.002)	0.004*** (0.001)	0.003 (0.002)
CONSTANT	0.489*** (0.010)	0.537*** (0.024)	0.201*** (0.024)	0.454*** (0.010)	0.502*** (0.025)
IS_PROMISE				0.173*** (0.013)	0.175*** (0.032)
N	1004	997	1004	1004	997
R <sup>2</sup>	0.097	0.015	0.086	0.238	0.042

\* :  $p < .05$ , \*\* :  $p < .01$ , \*\*\* :  $p < .001$

Initially, we find that messages with positive sentiment have a significant negative relationship with Trust, with each additional positive word being associated with a 2 percent decrease in Trust frequency. There is an additional, positive impact of increased message length on trust. However, we see matching, significant associations between both positive sentiment and total words with the likelihood a message is classified as a promise. There is no significant from positively or negatively charged words on Trustworthiness in our model.

In response to the similarities between the promise and Trust models, we include a control variable to test if our observed sentiment effect is driven by the correlation between promises and sentiment. After this control is added, the significance of this sentiment effect disappears. Neither positive nor negative words have an impact on the likelihood a receiver will trust them. As in the uncontrolled model, this is also true for Trustworthiness. In our environment, message sentiment is not a useful predictor of behavior.

## 6 Conclusion

This paper presents a novel approach to examining free-form communication data from laboratory experiments, avoiding human classification in favor of a machine learning approach, informed primarily by the data.

Our first approach, we consider each message as a series of independent tokens, reducing the message space into a manageable collection of binary variables that we can then approach with standard OLS regression methods. From this, we establish that our text data contain useful information even after controlling for manually-coded promise labels.

The negative impact of these messages on Trust suggests that the principals are unwilling to trust agents who demonstrate poor understanding of the game. This is consistent with the qualitative observation we made above in the selection of most and least trustworthy messages. Importantly, there is no corresponding significant effect when Trustworthiness is

our dependent variable. These seemingly confused agents are not significantly less likely to act in their principals' interests than the rest of the subject pool. This observation warrants further investigation in future experiments.

Like the 'out' token, 'right' has a highly significant negative coefficient. However, unlike 'out', 'right' has a corresponding significant impact on Trustworthiness as well. The impact on Trust is, perhaps, unsurprising; the impact on Trustworthiness, on the other hand, suggests something novel to the communication literature. These subjects tend to take the entire pool if entrusted with it but are upfront about their intentions in their pre-play communication. This indicates that there are subjects who, although intending to act against their partners' interests, retain a preference for honesty strong enough for them to reveal their intended action ahead of time. In our subject pool, there appear to be subjects unwilling to lie by omission.

In our second approach, we employ Affinity Propagation to group our data into clusters of rhetorically similar messages. We then use those cluster labels in OLS regressions as we did with individual tokens, again controlling for manually coded promise labels. We take a purposefully conservative approach to interpreting these regression results. Given that the cluster output is sensitive to our choices of calibration parameters and distance measurements, we cannot make strong statements about subtle differences between clusters. These clusters in and of themselves do not have substantial meaning to us as researchers; two messages appearing in the same cluster does not necessarily imply a definitive link between them. However, our clusters identify dense regions in our message space which we can then use to aggregate our unstructured data to yield enough power to perform meaningful analysis.

We conduct our cluster analysis using two separate approaches: one relying exclusively on information within the corpus of messages, and the other, knowledge-based, approach incorporates outside information from Google's word2vec natural language processing neural-network to link messages using similar, but not necessarily identical words.

In our first approach, comparing messages based solely on overlapping tokens, we find that Affinity propagation produces quality clusters, each with a clear substantive interpretation. These clusters allow us to identify several key patterns significantly associated both with higher rates of Trust and Trustworthiness, (Share.) First, highly effective promises consist of three parts: A specific promise to choose Left, a request for the recipient to choose In, and an explicit statement of the outcome if these things occur. Of our ten clusters with the most significant, positive coefficient for Trust, four followed this pattern, with an additional two having two of the three elements. The remaining four out of these ten messages fall the category of soft promises. These messages are not promises in the grammatical sense and were not labeled as such by our research assistants.

As for the least trusted messages, we again observe two broad patterns. First, as suggested by the highly significant coefficient estimates on "Right" and "Out" in our token regressions, most of the clusters significantly associated with recipients choosing Out are of poor quality and display potential lack of comprehension about the game's rules. These clusters include messages like "choose out" and "I AM PLAYER A" among other similarly incongruous statements. Eight of the ten least effective clusters fall into this group. More surprisingly, the remaining two clusters each contain messages where principals send brief greetings their partner: "hello... let's play" or "Hi, how are you?".

Despite these limitations, the findings we present make a strong case for the value of this

approach. We show that, even after controlling for human-coded promises, certain classes of messages are significantly more effective than others. The most effective messages contain three parts: a statement of intent, a call to action, and a promised benefit. Despite their positive impact on Trust, they are no more trustworthy than less detailed promises. In addition, I identify an under-trusted class of messages consisting of ‘soft promises.’ These messages are not as direct statements of intent or calls to action. Instead, they are vague personal statements such as “I believe in fairness” or “we should work together.” Removed from their context, readers are unlikely to classify these statements as promises. However, the implication within the game is clear. There is only one choice in the Agent’s action space reasonably interpret-able as fair; there is only one game outcome reasonably considered cooperation. Agents making these statements are significantly more trustworthy than the average message, demonstrating that we may want to expand our notion of a promise. Unfortunately for their authors, these soft promises are not met with statistically greater trust by their recipients. Senders making a promise should be clear about their intent as their meaning appears lost otherwise.

Finally, we measure each message’s positive or negative sentiment and find no significant correlations with behavior.

The above analyses have shown that the messages in our trust game contain useful information for predicting Trust beyond what is captured solely by identifying promises. By using natural language processing on our messages, we are able to successfully identify some of these features which we can then use to guide future work. Action labels; ‘left’, ‘right’, ‘in’, and ‘out’; remain significant predictors of principals’ behavior even after we control for promises using manually assigned labels. Their effects are in the expected directions, with In and Left positively impacting Trust and Out and Right negatively impacting Trust. When promises and non-promises are considered separately, we uncover further useful predictors for both Trust and Trustworthiness. Finally, a brief analysis of message sentiment does not find a significant effect on either Trust or Trustworthiness.

Out and Right are particularly interesting to us as it was initially apparent why anyone would choose to use them. ‘Out’ has a highly significant, negative impact on Trust, but no corresponding effect on Trustworthiness. Although ‘out’ is a valid action for the principals, a large portion of these messages are incoherent with the procedures of the game. Out of the 96 messages which contain the ‘out’ token, 25 contains statements that are clearly in violation of the rules of the game — mostly agents promising that they will play Out, which they cannot do in their current pairing. This number does not include a further 12 messages in which the agent instructs the principals to select Out themselves, which, although counter to the agents’ interests, is still a valid action the principal could take.

Currently, the underlying causes of these tokens’ effects are only conjectures. However, the ability to generate new, testable hypotheses from our messages is what makes natural language processing so useful. Having found a set of unexpected, significant tokens for predicting behavior, we can conduct further experiments to test potential reasons for their effect. The significant ‘out’ and ‘right’ tokens both would not have been uncovered without the use of natural language processing and do not conform well to existing theories of communication. For future work where we can test our conjectured explanations that principals trust agents who comprehend the game and that aversion to lying extends to lies of omission.

## References

- Andreoni, J. and Rao, J. M. (2011). The power of asking: How communication affects selfishness, empathy, and altruism. *Journal of public economics*, 95(7-8):513–520.
- Andres, M., Bruttel, L., and Friedrichsen, J. (2023). How communication makes the difference between a cartel and tacit collusion: A machine learning approach. *European Economic Review*, 152:104331.
- Ben-Ner, A., Putterman, L., and Ren, T. (2007). Lavish returns on cheap talk: Non-binding communication in a trust experiment. *Available at SSRN 1013582*.
- Buchan, N. R., Johnson, E. J., and Croson, R. T. (2006). Let’s get personal: An international examination of the influence of communication, culture and social distance on other regarding preferences. *Journal of Economic Behavior & Organization*, 60(3):373–398.
- Bursztyn, L., Egorov, G., Haaland, I., Rao, A., and Roth, C. (2023). Justifying dissent. *The Quarterly Journal of Economics*, 138(3):1403–1451.
- Camerer, C. F., Nave, G., and Smith, A. (2019). Dynamic unstructured bargaining with private information: theory, experiment, and outcome prediction via machine learning. *Management Science*, 65(4):1867–1890.
- Charness, G. and Dufwenberg, M. (2006). Promises and partnership. *Econometrica*, 74(6):1579–1601.
- Charness, G. and Dufwenberg, M. (2010). Bare promises: An experiment. *Economics letters*, 107(2):281–283.
- Charness, G. and Gneezy, U. (2008). What’s in a name? anonymity and social distance in dictator and ultimatum games. *Journal of Economic Behavior & Organization*, 68(1):29–35.
- Chen, D. L., Schonger, M., and Wickens, C. (2016). otree—an open-source platform for laboratory, online, and field experiments. *Journal of Behavioral and Experimental Finance*, 9:88–97.
- Chen, J. and Houser, D. (2017). Promises and lies: can observers detect deception in written messages. *Experimental Economics*, 20:396–419.
- Cooper, D. J. and Kühn, K.-U. (2016). Communication and cooperation: A methodological study. *Southern economic journal*, 82(4):1167–1185.
- Crawford, V. (1998). A survey of experiments on communication via cheap talk. *Journal of Economic theory*, 78(2):286–298.
- Devetag, G., Di Guida, S., and Polonio, L. (2016). An eye-tracking study of feature-based choice in one-shot games. *Experimental Economics*, 19:177–201.

- Ellingsen, T. and Johannesson, M. (2004). Promises, threats and fairness. *The Economic Journal*, 114(495):397–420.
- Ellingsen, T., Östling, R., and Wengström, E. (2018). How does communication affect beliefs in one-shot games with complete information? *Games and Economic Behavior*, 107:153–181.
- Fellbaum, C. (2010). Wordnet. In *Theory and applications of ontology: computer applications*, pages 231–243. Springer.
- Fiedler, M. and Haruvy, E. (2009). The lab versus the virtual lab and virtual field—an experimental investigation of trust games with communication. *Journal of Economic Behavior & Organization*, 72(2):716–724.
- Fonseca, M. A. and Normann, H.-T. (2012). Explicit vs. tacit collusion—the impact of communication in oligopoly experiments. *European Economic Review*, 56(8):1759–1772.
- Frey, B. J. and Dueck, D. (2007). Clustering by passing messages between data points. *science*, 315(5814):972–976.
- Fudenberg, D. and Liang, A. (2019). Predicting and understanding initial play. *American Economic Review*, 109(12):4112–4141.
- Grimm, S., Kocher, M. G., Krawczyk, M., and Le Lec, F. (2021). Sharing or gambling? on risk attitudes in social contexts. *Experimental Economics*, 24:1075–1104.
- Hanaki, N. and Ozkes, A. I. (2022). Strategic environment effect and communication. *Experimental Economics*, pages 1–34.
- Loper, E. and Bird, S. (2002). Nltk: The natural language toolkit. *arXiv preprint cs/0205028*.
- Lundquist, T., Ellingsen, T., Gribbe, E., and Johannesson, M. (2009). The aversion to lying. *Journal of Economic Behavior & Organization*, 70(1-2):81–92.
- Marcus, M., Santorini, B., and Marcinkiewicz, M. A. (1993). Building a large annotated corpus of english: The penn treebank.
- Mihalcea, R., Corley, C., Strapparava, C., et al. (2006). Corpus-based and knowledge-based measures of text semantic similarity. In *Aaai*, volume 6, pages 775–780.
- Mikolov, T., Chen, K., Corrado, G., and Dean, J. (2013). Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*.
- Mohlin, E. and Johannesson, M. (2008). Communication: Content or relationship? *Journal of Economic Behavior & Organization*, 65(3-4):409–419.
- Penczynski, S. P. (2019). Using machine learning for communication classification. *Experimental Economics*, 22(4):1002–1029.

- Polonio, L. and Coricelli, G. (2019). Testing the level of consistency between choices and beliefs in games using eye-tracking. *Games and Economic Behavior*, 113:566–586.
- Rehurek, R. and Sojka, P. (2011). Gensim–python framework for vector space modelling. *NLP Centre, Faculty of Informatics, Masaryk University, Brno, Czech Republic*, 3(2):2.
- Robinson, J., Rosenzweig, C., Moss, A. J., and Litman, L. (2019). Tapped out or barely tapped? recommendations for how to harness the vast and largely unused potential of the mechanical turk participant pool. *PloS one*, 14(12):e0226394.
- Snowberg, E. and Yariv, L. (2021). Testing the waters: Behavior across participant pools. *American Economic Review*, 111(2):687–719.
- Sparck Jones, K. (1972). A statistical interpretation of term specificity and its application in retrieval. *Journal of documentation*, 28(1):11–21.
- Tebbe, E. and Wegener, B. (2022). Is natural language processing the cheap charlie of analyzing cheap talk? a horse race between classifiers on experimental communication data. *Journal of Behavioral and Experimental Economics*, 96:101808.
- Tetlock, P. C. (2007). Giving content to investor sentiment: The role of media in the stock market. *The Journal of finance*, 62(3):1139–1168.
- Wilson, T., Wiebe, J., and Hoffmann, P. (2005). Recognizing contextual polarity in phrase-level sentiment analysis. In *Proceedings of human language technology conference and conference on empirical methods in natural language processing*, pages 347–354.
- Woods, D., Abdallah, M., Bagchi, S., Sundaram, S., and Cason, T. (2022). Network defense and behavioral biases: an experimental study. *Experimental Economics*, 25(1):254–286.

# Appendix

## 6.1 Affinity Propagation algorithm description

Affinity Propagation begins with a similarity matrix  $\mathbf{S}_{n \times n}$  and iteratively updates two matrices:  $\mathbf{R}_{n \times n}$ , the responsibility matrix and  $\mathbf{A}_{n \times n}$ , the availability matrix until stable clusters emerge. Both  $\mathbf{A}$  and  $\mathbf{R}$  initialize with all zero elements. The algorithm proceeds as follows:

1)  $\mathbf{R}$  updates following:

$$r(i, k) \leftarrow s(i, k) - \max_{\forall k \neq k'} \{a(i, k') + s(i, k')\}$$

2)  $\mathbf{A}$  updates following:

$$a(i, k) \leftarrow \min\{0, r(k, k) + \sum_{i' \notin i, k} \max(0, r(i', k))\}$$

$$a(k, k) \leftarrow \sum_{i' \neq k} \max(0, r(i', k))$$

Final clusters assignments are given by the criteria matrix  $\mathbf{C} = \mathbf{A} + \mathbf{S}$ ; element  $i$  chooses  $k$  as its exemplar if the highest value in row  $i$  in  $\mathbf{C}$  is in column  $k$ . This process repeats, updating  $\mathbf{A}$  and  $\mathbf{R}$  until cluster assignments remain unchanged between iterations.

Intuitively, the Affinity Propagation algorithm considers, first, each element's choice for its own exemplar out of the other elements in the set,  $\mathbf{R}$ . After this, these preferences are updated to account for the overall appropriateness of a given message to serve as an exemplar in the final cluster assignments. if a given message is chosen as a most preferred exemplar by relatively few other messages, it will incur a penalty reflecting its overall unsuitability to act as an exemplar.

Table 1 Expanded: Example Messages: Most and Least Effective

Type	Messages
Most Effective	<p>I'm going LEFT, so please select IN. This will maximize the earnings for both of us equally.</p> <p>Hello, I want to do what's best for the both of us so I'm going to be choosing left so that we can both get 12 tokens choose IN so that we will both get 12 tokens</p> <p>I'll choose LEFT if you choose IN.</p> <p>Choose "In" and I'll make sure we both get 12 tokens!</p> <p>I give you my word I will get each of us 12 if you go IN. That way we will both make more.</p> <p>Hi there. I would understand why you would choose OUT, but I promise to make the choice ...</p> <p>I promise you that I will always choose "left" so that we both get 12 points each.</p> <p>Hello!!! Let's make this work for both of us in the best possible way! Have a great day!</p> <p>Hey, I'm not a greedy person and in these situations, I always choose the most beneficial option for ...</p> <p>You can be sure that i will choose left, as it is the best option for both us. Please go ahead and select 'In'</p> <p>If pick in, I guarantee you will get 12 instead of 8 if you pick out</p> <p>I'm good for the 12. I believe that there are real people here, and I'm completely down with splitting the 24 and ...</p> <p>I'm picking left so if you want a few extra tokens, go for in. Totally understand if you don't but that's my choice.</p> <p>Let's work together to make things fair and so that we can both earn as many tokens as possible.</p> <p>I'm always going to choose "left" because I think it's dickish to screw another player over.</p> <p>So we can maximize our winnings, I will choose left on each turn if you choose In on each turn. That way we both...</p> <p>I will maximize both our outcomes ie LEFT if you stay in.</p> <p>Cooperation is better than being selfish. Love reigns.</p>
Least Effective	<p>I'm going to choose Right</p> <p>i get 24 you get 0</p> <p>U GO ON STRITE</p> <p>Hello Let's play the game</p> <p>I want choose Player A</p> <p>hi, how are you doing?</p> <p>Hey. let's work together, let's both choose out.</p> <p>I choose out</p> <p>choose out</p>



Table 17: New token regressions controlling for the four message types

<b>Tokens</b>	<b>Trust</b>	<b>T-stat</b>	<b>Trustworthiness</b>	<b>T-stat</b>
<b>i</b>	0.003	(0.212)	0.063	(1.411)
<b>choose</b>	-0.023	(-1.512)	-0.099*	(-2.259)
<b>in</b>	0.068***	(3.545)	0.036	(0.662)
<b>left</b>	0.060**	(3.278)	0.083	(1.608)
<b>you</b>	0.004	(0.217)	0.076	(1.618)
<b>to</b>	0.029*	(2.069)	0.088*	(2.187)
<b>will</b>	-0.005	(-0.282)	0.081	(1.773)
<b>be</b>	-0.012	(-0.847)	-0.077	(-1.942)
<b>the</b>	0.028	(1.761)	-0.069	(-1.517)
<b>both</b>	-0.013	(-0.635)	0.042	(0.764)
<b>we</b>	0.047*	(2.202)	0.138*	(2.306)
<b>and</b>	0.031*	(1.998)	0.046	(1.051)
<b>a</b>	-0.009	(-0.529)	-0.064	(-1.326)
<b>go</b>	-0.013	(-0.732)	-0.114*	(-2.326)
<b>if</b>	0.004	(0.188)	-0.001	(-0.022)
<b>get</b>	0.033	(1.689)	-0.023	(-0.430)
<b>of</b>	0.016	(0.754)	-0.020	(-0.333)
<b>for</b>	0.008	(0.398)	0.106*	(1.974)
<b>let</b>	0.011	(0.487)	-0.116	(-1.815)
<b>so</b>	-0.026	(-1.196)	-0.108	(-1.726)
<b>u</b>	0.048*	(2.007)	0.003	(0.043)
<b>player</b>	-0.023	(-0.863)	0.060	(0.794)
<b>'s</b>	0.001	(0.038)	0.159*	(2.363)
<b>that</b>	0.042*	(2.043)	0.043	(0.742)
<b>can</b>	0.046*	(2.084)	-0.143*	(-2.275)
<b>make</b>	0.066**	(3.015)	0.010	(0.165)
<b>out</b>	-0.101***	(-4.933)	0.087	(1.509)
<b>hi</b>	-0.001	(-0.055)	0.047	(0.866)
<b>hello</b>	0.029	(1.464)	0.069	(1.243)
<b>12</b>	0.041	(1.664)	-0.007	(-0.102)
<b>play</b>	-0.022	(-1.007)	0.028	(0.462)
<b>do</b>	-0.021	(-0.947)	-0.062	(-1.009)
<b>token</b>	0.018	(0.720)	-0.074	(-1.063)
<b>'m</b>	0.010	(0.411)	0.009	(0.137)
<b>fair</b>	0.052*	(2.316)	0.113	(1.793)
<b>this</b>	0.063**	(2.730)	0.099	(1.510)
<b>it</b>	-0.046	(-1.856)	-0.120	(-1.708)
<b>good</b>	0.004	(0.189)	-0.039	(-0.615)
<b>'ll</b>	-0.028	(-1.127)	0.046	(0.655)
<b>best</b>	0.028	(1.170)	0.005	(0.076)
<b>please</b>	0.013	(0.549)	0.070	(1.065)
<b>more</b>	0.030	(1.181)	-0.139	(-1.924)
<b>maximize</b>	0.067*	(2.247)	-0.148	(-1.692)
<b>point</b>	-0.031	(-1.285)	0.025	(0.357)
<b>our</b>	-0.007	(-0.198)	<sup>41</sup> -0.004	(-0.042)
<b>each</b>	-0.015	(-0.545)	0.008	(0.100)
<b>most</b>	-0.023	(-0.808)	0.121	(1.532)

Table 5 Expanded: Regressions Using Bag of Words Model: Promises and Non-Promises Separate

Coefficients	Promises Only				No Promises			
	Trust	T-stat	Trustworthiness	T-stat	Trust	T-stat	Trustworthiness	T-stat
i	-0.033	(-0.802)	0.060	(0.489)	0.021	(1.147)	0.063	(1.244)
choose	0.013	(0.554)	-0.011	(-0.163)	-0.024	(-1.119)	-0.070	(-1.183)
in	0.013	(0.604)	0.039	(0.602)	0.119***	(6.089)	0.017	(0.313)
left	0.004	(0.139)	-0.138	(-1.765)	0.080**	(3.019)	0.251***	(3.426)
you	0.018	(0.721)	-0.039	(-0.539)	-0.004	(-0.205)	0.090	(1.482)
to	-0.004	(-0.202)	0.098	(1.541)	0.057**	(3.016)	0.069	(1.321)
will	0.001	(0.068)	0.192**	(3.022)	-0.008	(-0.310)	0.060	(0.875)
be	-0.026	(-1.180)	-0.104	(-1.590)	-0.014	(-0.757)	-0.041	(-0.794)
the	0.045	(1.840)	-0.025	(-0.348)	0.007	(0.344)	-0.122*	(-2.074)
both	-0.020	(-0.699)	-0.006	(-0.072)	0.017	(0.605)	0.098	(1.295)
we	-0.005	(-0.141)	0.046	(0.480)	0.099***	(3.453)	0.098	(1.237)
and	0.001	(0.072)	-0.004	(-0.071)	0.034	(1.600)	0.035	(0.594)
a	-0.016	(-0.646)	-0.081	(-1.117)	-0.020	(-0.886)	-0.057	(-0.908)
go	0.004	(0.171)	0.005	(0.069)	-0.018	(-0.676)	-0.185*	(-2.509)
if	0.002	(0.081)	0.042	(0.530)	-0.004	(-0.104)	-0.003	(-0.030)
get	0.068**	(2.703)	-0.026	(-0.351)	0.036	(1.278)	0.003	(0.042)
of	-0.001	(-0.047)	-0.136	(-1.665)	-0.050	(-1.639)	0.039	(0.464)
for	0.032	(1.194)	0.145	(1.809)	-0.009	(-0.333)	0.100	(1.371)
let	-0.039	(-0.868)	-0.037	(-0.270)	0.041	(1.488)	-0.123	(-1.617)
so	-0.008	(-0.299)	0.048	(0.575)	0.024	(0.622)	-0.254*	(-2.385)
u	0.024	(0.778)	0.050	(0.537)	0.089*	(2.553)	-0.076	(-0.791)
player	-0.016	(-0.315)	0.308*	(2.042)	0.020	(0.657)	0.016	(0.197)
's	0.051	(1.313)	0.053	(0.446)	-0.036	(-1.155)	0.164	(1.933)
that	0.045	(1.865)	0.100	(1.410)	-0.000	(-0.002)	-0.102	(-1.114)
can	0.033	(1.209)	-0.253**	(-3.145)	-0.014	(-0.403)	-0.042	(-0.444)
make	0.085*	(2.581)	-0.194*	(-2.001)	0.043	(1.414)	0.136	(1.639)
out	-0.014	(-0.369)	0.080	(0.718)	-0.126***	(-5.112)	0.072	(1.070)
hi	-0.039	(-1.260)	0.238**	(2.604)	0.013	(0.529)	-0.056	(-0.812)
hello	0.094*	(2.535)	0.290**	(2.647)	0.027	(1.108)	-0.020	(-0.299)
12	0.025	(0.856)	0.089	(1.010)	0.064	(1.601)	-0.171	(-1.546)
play	-0.004	(-0.095)	0.198	(1.701)	-0.030	(-1.170)	-0.007	(-0.094)
do	0.010	(0.305)	0.019	(0.197)	-0.042	(-1.415)	-0.100	(-1.212)
token	-0.043	(-1.390)	-0.098	(-1.052)	0.075*	(1.968)	-0.059	(-0.563)
'm	0.054	(1.841)	0.056	(0.647)	-0.044	(-1.060)	-0.014	(-0.125)
fair	-0.003	(-0.101)	-0.052	(-0.602)	0.102**	(2.906)	0.222*	(2.288)
this	0.054	(1.363)	0.148	(1.264)	0.058*	(2.019)	0.134	(1.690)
it	-0.077*	(-2.400)	-0.124	(-1.309)	-0.027	(-0.756)	-0.038	(-0.386)
good	0.029	(0.647)	0.063	(0.478)	0.003	(0.112)	-0.108	(-1.454)
'll	-0.054*	(-1.993)	0.104	(1.296)	0.017	(0.321)	-0.109	(-0.702)
best	-0.014	(-0.344)	0.094	(0.807)	0.051	(1.623)	-0.077	(-0.897)
please	0.050	(1.600)	-0.017	(-0.179)	-0.006	(-0.185)	0.072	(0.784)
more	0.007	(0.207)	0.033	(0.304)	0.045	(1.307)	-0.158	(-1.669)
maximize	0.029	(0.752)	-0.250*	(-2.090)	0.125*	(2.384)	-0.190	(-1.322)
point	-0.025	(-0.788)	0.053	(0.565)	-0.023	(-0.615)	0.059	(0.583)
our	0.004	(0.100)	0.105	(0.839)	-0.021	(-0.376)	0.048	(0.306)
each	0.009	(0.254)	-0.002	(-0.017)	-0.005	(-0.107)	0.027	(0.225)
most	-0.009	(-0.216)	0.182	(1.478)	-0.009	(-0.234)	0.119	(1.143)
game	0.009	(0.154)	-0.163	(-0.973)	-0.028	(-0.817)	-0.012	(-0.123)
all	-0.053	(-1.209)	-0.171	(-1.312)	-0.043	(-1.218)	0.063	(0.642)
b	0.034	(0.485)	-0.461*	(-2.210)	0.014	(0.366)	0.004	(0.037)
always	0.069**	(2.734)	0.154*	(2.023)	0.012	(0.159)	0.168	(0.793)
money	-0.014	(-0.274)	-0.174	(-1.130)	0.032	(0.842)	-0.002	(-0.022)
pick	0.032	(0.970)	0.115	(1.167)	0.060	(1.145)	-0.010	(-0.069)
with	0.045	(1.145)	-0.066	(-0.558)	-0.030	(-0.767)	0.065	(0.595)
have	0.022	(0.561)	0.050	(0.427)	0.041	(1.069)	0.144	(1.342)
right	-0.036	(-0.565)	-0.023	(-0.121)	-0.130***	(-4.158)	-0.246**	(-2.827)
CONSTANT	0.653***	(14.807)	0.611***	(4.584)	0.400***	(24.460)	0.514***	(11.373)
N	381		379	42	623		618	
R <sup>2</sup>	0.24		0.22		0.38		0.16	

