

A peer-reviewed version of this preprint was published in PeerJ on 13 January 2016.

[View the peer-reviewed version](https://doi.org/10.7717/peerj-cs.40) (peerj.com/articles/cs-40), which is the preferred citable publication unless you specifically need to cite this preprint.

Kazemzadeh A, Gibson J, Georgiou P, Lee S, Narayanan S. 2016. A Socratic epistemology for verbal emotional intelligence. PeerJ Computer Science 2:e40 <https://doi.org/10.7717/peerj-cs.40>

A Socratic epistemology for verbal emotional intelligence

Abe Kazemzadeh, James Gibson, Panayiotis Georgiou, Sungbok Lee, Shrikanth Narayanan

We describe and experimentally validate a question-asking framework for machine-learned linguistic knowledge about human emotions. Using the Socratic method as a theoretical inspiration, we develop an experimental method and computational model for computers to learn subjective information about emotions by playing emotion twenty questions (EMO20Q), a game of twenty questions limited to words denoting emotions. Using human-human EMO20Q data we bootstrap a sequential Bayesian model that drives a generalized pushdown automaton-based dialog agent that further learns from 300 human-computer dialogs collected on Amazon Mechanical Turk. The human-human EMO20Q dialogs show the capability of humans to use a large, rich, subjective vocabulary of emotion words. Training on successive batches of human-computer EMO20Q dialogs shows that the automated agent is able to learn from subsequent human-computer interactions. Our results show that the training procedure enables the agent to learn a large set of emotions words. The fully trained agent successfully completes EMO20Q at 67% of human performance and 30% better than the bootstrapped agent. Even when the agent fails to guess the human opponent's emotion word in the EMO20Q game, the agent's behavior of searching for knowledge makes it appear human-like, which enables the agent maintain user engagement and learn new, out-of-vocabulary words. These results lead us to conclude that the question-asking methodology and its implementation as a sequential Bayes pushdown automaton are a successful model for the cognitive abilities involved in learning, retrieving, and using emotion words by an automated agent in a dialog setting.

A Socratic Epistemology for Verbal Emotional Intelligence

Abe Kazemzadeh, James Gibson, Panayiotis Georgiou, Sungbok Lee, Shrikanth Narayanan

Signal Analysis and Interpretation Laboratory

University of Southern California, Los Angeles, USA

Abstract

We describe and experimentally validate a question-asking framework for machine-learned linguistic knowledge about human emotions. Using the Socratic method as a theoretical inspiration, we develop an experimental method and computational model for computers to learn subjective information about emotions by playing emotion twenty questions (EMO20Q), a game of twenty questions limited to words denoting emotions. Using human-human EMO20Q data we bootstrap a sequential Bayesian model that drives a generalized pushdown automaton-based dialog agent that further learns from 300 human-computer dialogs collected on Amazon Mechanical Turk. The human-human EMO20Q dialogs show the capability of humans to use a large, rich, subjective vocabulary of emotion words. Training on successive batches of human-computer EMO20Q dialogs shows that the automated agent is able to learn from subsequent human-computer interactions. Our results show that the training procedure enables the agent to learn a large set of emotions words. The fully trained agent successfully completes EMO20Q at 67% of human performance and 30% better than the bootstrapped agent. Even when the agent fails to guess the human opponent's emotion word in the EMO20Q game, the agent's behavior of searching for knowledge makes it appear human-like, which enables the agent maintain user engagement and learn new, out-of-vocabulary words. These results lead us to conclude that the question-asking methodology and its implementation as a sequential Bayes pushdown automaton are a successful model for the cognitive abilities involved in learning, retrieving, and using emotion words by an automated agent in a dialog setting.

1. Introduction

Epistemology is the branch of philosophy that deals with knowledge and belief. According to basic results in epistemology, knowledge is defined as true, justified belief. This paper was inspired by reflecting on how humans justify their beliefs about emotions. This reflection led to a experimental method for collecting human knowledge about emotions and a computational model that uses the collected knowledge in an automated dialog agent.

The logician Charles S. Peirce identified three types of thought processes by which a person can justify their beliefs and thereby acquire knowledge: induction, deduction, and hypothesis (Peirce, 1868). Whereas

Email addresses: abe.kazemzadeh@gmail.com (Abe Kazemzadeh), jjgibson@usc.edu (James Gibson), georgiou@isi.usc.edu (Panayiotis Georgiou), sungbok1@usc.edu (Sungbok Lee), shri@isi.usc.edu (Shrikanth Narayanan)

induction is primarily involved with observational data, deduction and hypothesis have a linguistic, propositional component. The third of these, hypothesis (also known as abduction (Eco and Sebeok, 1988)), has been compared with the Socratic method of questioning-asking dialogs (Hintikka, 2007). The Socratic method was named after the ancient Greek philosopher Socrates who applied his method of inquiry to
35 examine concepts that seem to lack any concrete definition, in particular some of the complex moral and psychological concepts of his time like “justice”, “knowledge”, “piety”, “temperance”, and “love”. We claim that this method of inquiry can shed light on how people justify beliefs about emotional concepts, which also seem to defy concrete definition.

Question-asking allows people to learn about things without directly experiencing them. Since a computer
40 agent cannot directly experience emotions as a human would, question-asking can be leveraged for the computer agent to learn about emotional concepts. Question-asking has also been proposed as a stage in child development responsible for rapid learning and language acquisition (Frazier et al., 2009). Likewise, a computer agent can use question-asking to acquire knowledge and vocabulary. We call the approach of using question-asking to interactively acquire linguistic knowledge about emotion by a computer dialog agent
45 a *Socratic epistemology for verbal emotional intelligence*.

The knowledge acquired by the Socratic epistemology for verbal emotional intelligence is an informal, social type of knowledge. This informal knowledge about emotions is important because although there has been much recent progress toward understanding the underlying biological basis for emotion, humans have been able to understand emotions informally since ancient times. We call this informal, language-based
50 understanding of emotions *natural language description of emotion* (Kazemzadeh, 2013). Natural language descriptions of emotion are utterances that refer to emotions, as opposed to utterances that express emotions. This phenomenon can be seen as a specific subset of the larger phenomenon of emotional language, which also includes emotion or sentiment expressed towards some object, vocal modulation due to emotion, and persuasion and pragmatics. Studying language that deals with referential statements about emotions is a
55 novel focus that is distinct from the prevailing trends of studying the expressive characteristics of emotional language.

The framework we present also differs from other computational theories of emotion in that it aims to study how people describe emotions, rather than how emotions *should be* described. As such, it can be seen as a descriptive, rather than prescriptive, theory, and hence has commonalities with sociological studies of
60 emotions (King, 1989; Shaver et al., 2001; Mihalcea and Liu, 2006; Kamvar and Harris, 2009; Kazemzadeh et al., 2012) and crowd-sourcing. Although there is clearly a biological basis for emotions, the way that people describe emotions in natural language varies over time from culture to culture and person to person. A descriptive theory of how people describe emotions in natural language should reflect this variation. Thus, a particular instantiation of a theory of natural language description of emotion should be relative to a

65 particular agent or set of agents at a given time who hold similar beliefs and communicate using the same language. Rather than trying to establish a single, particular definition for emotion words, we study how shared meaning can be established through natural language dialog behavior even when people and agents may have different ways of experiencing and their describing emotions.

The key goals of this paper are to use question-asking to observe natural language descriptions of emotion in the context of interactive dialogs and to computationally model the communicative behavior that supports the referential link between language and emotions. To do this, we use an experimental game called emotion twenty questions (EMO20Q). EMO20Q is played like the traditional game of twenty questions, except that one player, the *answerer*, must choose a word that names an emotion. The other player, the *questioner*, must try to guess the emotion that the answerer picked by posing a series of twenty or fewer questions. 75 Importantly, we do not restrict the set of emotion words that the players can choose nor the questions they can ask: the set of emotion words and questions is not limited except by the players' judgment and cooperation. An example of a human-computer EMO20Q game is given below:

Agent: *Let me know when you have picked an emotion word.*

User: *ok*

80 **Agent:** *Is it very intense?*

User: *no*

Agent: *is it positive?*

User: *yes*

Agent: *is it an emotion that just pops up spontaneously (vs being triggered by something)?*

85 **User:** *depends*

Agent: *is it caused by anything in particular*

User: *well, yes, you need the right conditions*

Agent: *Is it relaxation?*

User: *Yes*

90 **Agent:** *So did I get it right?*

User: *Yes, you got it.*

Agent: *Awesome! Would you like to play again?*

User: *No, thanks.*

Agent: *Thank you for playing*

95 The data from the human-human EMO20Q games allows us to observe what words people consider to denote emotions, what questions people ask about emotions, what factors determine the length of question/answer exchanges, and what is the probability that the emotion words will be correctly guessed. These results are used to create the computer agent and test its performance with respect to human performance. For the

automated agent presented in this paper, we focused on the case where the computer agent plays the role of the questioner while future work for an answerer agent is discussed in Section 6.

The paper is organized as follows. Section 2 discusses the motivations and theory behind our work. Section 3 describes the computational model and algorithm we used to create an EMO20Q questioner agent. Section 4 discusses experiments we conducted of humans and computers playing EMO20Q. Section 5 describes the results of testing the agent. Finally Section 6 and Section 7 propose future work and provide discussion and links to open source software implementations.

2. Background

2.1. Natural Language Descriptions of Emotions

Just as memory addresses, variables, and URLs refer to electronic resources for computers, so do words and descriptions identify objects, both physical and conceptual, for humans. When processing natural language by computer, it can help to draw upon these similarities. This is especially helpful in the case of affective computing, when the objects we wish to refer to, emotions, are abstract and subjective.

In this paper we make a distinction between the emotion *expressed by* the speaker and the emotion *referred to* by the speaker. Currently there has been a great degree of interest in automatically analyzing emotional expression in language. The goal of such analysis is to determine emotions *expressed by* the speaker or writer, i.e., the emotions that the speaker currently feels. The language used as input to this kind of analysis can be a speech recording or textual representation of language. However, automatically analyzing the emotions expressed in an utterance or document is problematic when a speaker *refers to* emotions that are not his or her own current emotions. Some examples of this include quotations, storytelling/gossip, counterfactual reasoning, *post facto* emotional self-report, and abstract references to emotions.

He said that he was mad. (quotation)

Did you see how mad John was? (gossip)

If you eat my ice cream, I will get mad. (counterfactual)

I was mad when my car got stolen last year. (self-report)

Anger is one of the seven sins. (abstract reference).

In these examples, a naïve automated analysis would detect anger, but in fact the writer of these sentences is not actually feeling anger at the current time. In many cases, such as task-driven dialogs like ordering airline tickets from an automated call center, this distinction might not be pertinent. However, for open-ended dialog systems the distinction between expression and reference of emotions could be relevant, for example an automated agent for post-traumatic stress disorder therapy. The study of natural language descriptions of emotions brings the distinction between emotion expression and reference into focus.

The ability to talk about things beyond the here-and-now has been termed *displacement* (Hockett and Altmann, 1968). Displacement is an important characteristic that distinguishes human language from animal communication. In the context of this research, the ability to talk about an emotion without it being physically present is a key component of natural language description of emotion. Natural language description of emotion has been examined in ethnography, comparative linguistics, and cognitive science and it is beginning to be studied in the domain of natural language processing (King, 1989; Zoltán Kövecses, 2000; Rolls, 2005; Kazemzadeh et al., 2012).

At the most basic level, natural language description of emotion includes words that name emotions, e.g. *angry*, *happiness*, etc. However, due to the productive, generative nature of natural language, it is possible to refine and generalize emotion descriptions with longer natural language phrases. In order to communicate using natural language descriptions of emotions, people must be able to come to a shared understanding about the meaning of these descriptions. Russell (1905) introduced the notion of *definite descriptions*, a logical device used to model unique reference in the semantics of languages, both formal and natural. In this paper, we focus on the natural language definite descriptions. Common examples of natural language definite descriptions are proper names and noun phrases with the definite article “the”. Indefinite descriptions, on the contrary, are prefaced with indefinite articles, such as “a”, “some”, or “every”.

We maintain that natural language descriptions of emotions are definite descriptions when they are used in natural language interaction that terminates in mutual agreement. By considering terms that refer to emotions as definite descriptions, we are trying to capture the intuition that different people mean the same things when they use the same emotion terms. In Barrett (2006), the question is posed of whether emotions are natural kind terms, to which the paper answered no, i.e., that emotion words in general represent non-unique classes of human behavior rather than fundamentally distinct biological classes. The question of whether emotion terms are definite descriptions can be seen as a less stringent criterion than that of whether they are natural kinds. In this paper, we apply the notion of definite descriptions to capture the experimental data which indicates that there is a high degree of consensus about how emotions are described when measured by successful outcomes in human-human EMO20Q.

2.2. EMO20Q, Crowd-Sourcing, and Experimental Design

The game of EMO20Q was designed as a way to elicit natural language descriptions of emotion. Posing the experiment as a game leverages past results in crowd-sourcing and games with a purpose. From the perspective of natural language processing, the EMO20Q game can be seen as a *Wizard of Oz* experiment that collects human behavior to train the behavior of an automated agent. Games like EMO20Q can be seen as *games with a purpose* (von Ahn and Dabbish, 2004) whose purpose is *crowd-sourcing* (Howe, 2006) the collective knowledge and beliefs of the players (Kazemzadeh et al., 2011). The phenomenon of crowd-sourcing is closely tied to the emergent properties of online social communities (Zhong et al., 2000).

By relying on the wisdom of the masses, we venture a simple answer to the difficult question, “what is emotion?”. The answer, according to crowd-sourcing, is that emotion is what people say it is. Although this answer side-steps many important issues, such as physiological and psychological descriptions of emotions, it does bring other issues into sharper focus. There has been a trend toward studying *non-prototypical* emotional data (Mower et al., 2009). Non-prototypical emotional data is exemplified by disagreement among annotators when assigning emotional labels to data. We argue that our methodology provides a crowd-sourced description of emotions that can effectively deal with non-prototypical emotions. To avoid falling into the *ad populum* logical fallacy, we formulate the answer to the question “what is emotion?” not as a question of truth, but a question of knowledge and belief, i.e., an issue of epistemology as described in Section 1, in effect skirting the question of ground truth, but asking other interesting questions: “what do people believe about emotions, how do they express these beliefs in language, and how do they justify their beliefs through question-asking behavior?”

Annotation tasks can be seen as a type of crowd-sourcing to find consensus about assigning emotional labels to data. Elicitation of subjects also has aspects of crowd-sourcing to experimentally observe a diversity of emotional behavior in response to experimentally controlled stimuli. It can be argued that compared with annotation and elicitation of emotional data EMO20Q provides higher experimental validity and sensitivity and less experimental bias at the expense of experimental control and reliability.

In terms of experimental design, the human-human EMO20Q is a *quasi-experiment* or *natural experiment*, as opposed to a *controlled experiment*, which means that there is not a manipulation of variables made by the experimenters, but rather that these variables are observed as they vary naturally within the system. With annotation and elicitation tasks, experimenters can control the vocabulary of annotation labels and with elicitation tasks experimenters can control the stimuli that are presented. With this control, experiments are more easily repeated. In EMO20Q, we did not control what emotion words or questions the subjects picked so for another population the results could vary, leading to less experimental reliability. However, trading off control and reliability leads to more experimental sensitivity and validity and less experimental bias. In EMO20Q subjects can choose any words or questions they want and they communicate in a natural dialog setting. This way of characterizing emotion is closer to natural communication and more sensitive to nuances of meaning. When forced to annotate using a fixed vocabulary of emotion words, subjects are experimentally biased toward using that vocabulary.

The automated dialog agent is one way to enforce more experimental control for EMO20Q. Because the agent’s behavior is programmed we can use this as a way to better control and replicate experiments. Another way we aimed to improve experimental reliability is by prompting users to pick emotion words from three different difficulty classes. Sections 3 and 4 further describe the computational model for the agent’s behavior and our experimental design.

3. Model

Bayesian models have been successfully applied to a wide range of human cognitive abilities (Griffiths et al., 2008), including inductive inference of word meaning from corpora (Steyvers et al., 2006) and experimental stimuli (Xu and Tenenbaum, 2005) and powering affective dialog agents (Carofiglio et al., 2009). To our knowledge, this work is the first application of Bayesian cognitive models for learning emotion words from dialog interaction.

The model we use for the EMO20Q questioner agent is a sequential Bayesian belief update algorithm. This model fits the framework of Socratic epistemology, as described in the introduction, because it combines the notion of belief and question-asking. Intuitively, this algorithm instantiates an agent whose *semantic knowledge* is based on data from previous EMO20Q matches. The agent begins a new match of EMO20Q with a uniform *belief* about the emotion word to be guessed. Based on the previous semantic knowledge, the agent asks questions and updates its belief based on each *observation* of the user's answers to the questions. While the EMO20Q match is played, the observations are stored in the agent's *episodic buffer* (Baddeley, 2000), also known as working memory. After the match, the agent updates its semantic knowledge using the results of the match, clears its episodic buffer, and is then ready to play again. The words in italics are high-level abstractions used to create a cognitive model for the agent, which is underlyingly implemented as a sequential Bayesian statistical model. We ask that the reader keep this abstraction in his or her episodic buffer when reading the following description of the model's technical implementation.

The *semantic knowledge* described above is the conditional probability of observing a set of question-answer pairs given a hidden variable ranging over emotion words. This conditional probability distribution is estimated from the corpus of past human-human and human-computer EMO20Q matches as follows. Let E be the set of emotion words and let $\varepsilon \in E$ be this categorical, Bayesian (i.e., unobserved) random variable distributed over the set E . The probability of ε , $P(\varepsilon)$ is the *belief* about the emotion word to be guessed. Each question-answer pair from the match of EMO20Q is considered as an *observation* or feature of the emotion being predicted. Thus if Q is the set of questions and A is the set of answers, then a question $q \in Q$ and an answer $a \in A$ together compose the feature $f = (q, a)$, i.e. $f \in Q \times A$. The conditional probability distribution, $P(f|\varepsilon)$, which represents semantic knowledge, is estimated from the training data using a smoothing factor of 0.5 to deal with sparsity.

In this model we stipulate that the set of answers A are four discrete cases: "yes", "no", "other", and "none". When the answer either contains "yes" or "no", it is labeled accordingly. Otherwise it is labeled "other". The feature value "none" is assigned to all the questions that were not asked in a given dialog. "None" can be seen as a missing feature when the absence of a feature may be important. For example, the fact that a certain question was not asked about a particular emotion may be due to the fact that that question was not relevant at a given point in a dialog.

Similarly, we stipulate that the questions can be classified into some discrete class that is specified through a semantic expression derived from the annotation of questions, as described in Section 4.1. For example, the question “is it a positive emotion?” is represented as the semantic expression “e.valence==positive”. If the answer to this question was “maybe”, the resulting feature would be represented as (‘e.valence==positive’, ‘other’).

Using Bayes rule and the independence assumption of the naïve Bayes model, we can formulate the agent’s belief about the emotion vector ε after observing features $f_1 \dots f_t$, in one single batch, as opposed to sequentially (which will be formulated next):

$$P(\varepsilon|f_1, \dots, f_t) = \frac{\prod_{i=1}^t [P(f_i|\varepsilon)] P(\varepsilon)}{\prod_{i=1}^t P(f_i)}. \quad (1)$$

This is simply the formulation of naïve Bayes, where in this case $P(\varepsilon)$ is the prior probability of a player choosing a specific emotion word, $\prod_{i=1}^t [P(f_i|\varepsilon)]$ is the likelihood of seeing question-answer pairs given specific emotion words, and $\prod_{i=1}^t P(f_i)$ is the probability of observing question-answer pairs in general.

In terms of the high-level cognitive model, the set of observational feature vector $f_1 \dots f_t$ is what was described as the agent’s *episodic buffer*. $P(f|\varepsilon)$ is the agent’s *semantic knowledge* that relates question-answer features to emotion words. $p(\varepsilon)$ and $P(\varepsilon|f_1, \dots, f_t)$ are the agent’s initial/prior and final/posterior beliefs, respectively.

In Equation 1, the posterior belief of the agent of emotion e_k at time t , $P(\varepsilon = e_k|f_1, \dots, f_t)$ is computed only after the agent has asked all t questions. This model is known as naïve Bayes. In contrast the sequential Bayes model that we use is dynamic: the agent updates its belief at each time point based on the posterior probability of the previous step, i.e., at time t

$$P(\varepsilon|f_1, \dots, f_t) = \frac{P(f_t|\varepsilon)P(\varepsilon|f_1, \dots, f_{t-1})}{P(f_t)}$$

When the game begins, the agent can start with a uniform prior on its belief of which emotion is likely or it can use information obtained in previously played games. In the experiments of this paper, we use a uniform prior, $P(\varepsilon = e_k) = 1/|E|$, $\forall k = 1 \dots |E|$. We chose to use the uniform prior to initialize the agent because our training data contains many single count training instances and because we want to examine how the system performs with less constraints.

We introduce a new variable $\beta_{t,k} = P(\varepsilon = e_k|f_1, \dots, f_t)$ for the agent’s belief about emotion k at time t and postulate that the agent’s prior belief at a given time is the posterior belief of the previous step. Then, the agent’s belief unfolds according to the formula:

$$\begin{aligned} \beta_{0,k} &= P(\varepsilon = e_k) = 1/|E| \\ \beta_{1,k} &= \frac{P(f_1|\varepsilon = e_k)}{P(f_1)} \beta_{0,k} \\ \beta_{t,k} &= \frac{P(f_t|\varepsilon = e_k)}{P(f_t)} \beta_{t-1,k} \end{aligned} \quad (2)$$

Decomposing the computation of the posterior belief allows the agent to choose the best question to ask the user at each turn, rather than having a fixed battery of questions. We define “the best question” at time t to be the question that is most likely to have a “yes” answer given the posterior belief at time $t - 1$, $P(\varepsilon|f_1, \dots, f_{t-1})$:

$$\operatorname{argmax}_{q \in Q} P((q, \text{'yes'}) | \varepsilon) P(\varepsilon | f_1, \dots, f_{t-1})$$

This next-question criterion is a heuristic motivated by considering “yes” answers to be positive feedback that the agent is on the right track. While this heuristic worked well in practice, other next-question criteria are certainly possible and this is an area for future research.

At time t the agent asks the best question and takes the user’s response as input. It then parses the input to classify it into one of {“yes”, “no”, “other”}. This information is then used to update the agent’s posterior belief $\beta_{t+1,k}$ about each emotion $e_k \in E$, which will then be used as the prior in the following step. The unfolding of variable β in Equation 2 models the update of belief as it is justified by the agent’s question-asking and the user’s answers. It is this computational model of question-asking and belief update that represents the Socratic epistemology for verbal emotional intelligence in a software agent. Table 1 shows an example interaction between the automated EMO20Q questioner agent and a human user, along with a trace of the agent’s belief state that shows the justification of beliefs by question-asking.

Identity questions are a special type of question where the agent makes a guess about the emotion. Identity questions are chosen with the same best question criteria as other questions but trigger a transition to a different dialog state. An affirmative answer to an identity question (e.g., “is it happy?”) means that the agent successfully identified the user’s chosen emotion. Any other answer to an identity question will set the posterior probability of that emotion to zero because the agent can be sure it is not the emotion of interest.

The pseudo-code for the main loop of the adaptive Bayesian agent is shown in Algorithm 1. This automated, data-driven component was framed within a manually designed dialog graph, as shown in Figure 1. The dialog graph is implemented as a generalized pushdown transducer. Recall that a pushdown transducer is an transducer that can determines it output symbol and next state based on its current state, the input symbol, and the top of its stack (Allauzen and Riley, 2012). A generalized pushdown transducer is a pushdown transducer that is not limited to only the top of the stack when determining the output and next state. This aspect is important in the question asking loop because the stack represents the episodic memory, which stores the question-answer observations. Otherwise, the agent could be implemented as a plain pushdown transducer.

Algorithm 1 Sequential Bayesian EMO20Q agent. F is the observed question-answer features, E is the set of previously seen emotion words, $P(f|\varepsilon)$ is the semantic knowledge relating the observed question-answer pairs to emotion words, and $\beta_{t,k}$ is the belief about the emotion word indexed by k at time t . Because the agent is playing a twenty questions game, d is set to 20, but this could be changed for the agent to generalize to different question-asking tasks.

Input: $F = Q \times A$, E , and $P(f|\varepsilon)$

$\beta_{0,k} \leftarrow 1/|E|, \forall k = 1 \dots |E|$

for $i = 1$ **to** d **do**

$q^{(i)} = \underset{q \in Q}{\operatorname{argmax}} P((q, \text{'yes'})|\varepsilon)P(\varepsilon|f_1, \dots, f_{i-1})$

Print $q^{(i)}$

$a^{(i)} \leftarrow$ user's input answer

$f_i \leftarrow (q^{(i)}, a^{(i)})$

$\beta_{i,k} \leftarrow \beta_{i-1,k} \cdot P(f_i|\varepsilon = e_k)/P(f_1, \dots, f_i), \forall k = 1 \dots |E|$

if $(q^{(i)}$ is identity question for $e_k \wedge a^{(i)} = \text{'yes'}$) **then**

Return: $e^* = e_k$

end if

if $(q^{(i)}$ is identity question for $e_k \wedge a^{(i)} = \text{'no'}$) **then**

$\beta_{i,k} \leftarrow 0$

end if

end for

$k^* \leftarrow \underset{k \in 1 \dots |E|}{\operatorname{argmax}} [\beta_{i,k}]$

$e^* \leftarrow e_{k^*}$

Return: most likely emotion given observations: e^*

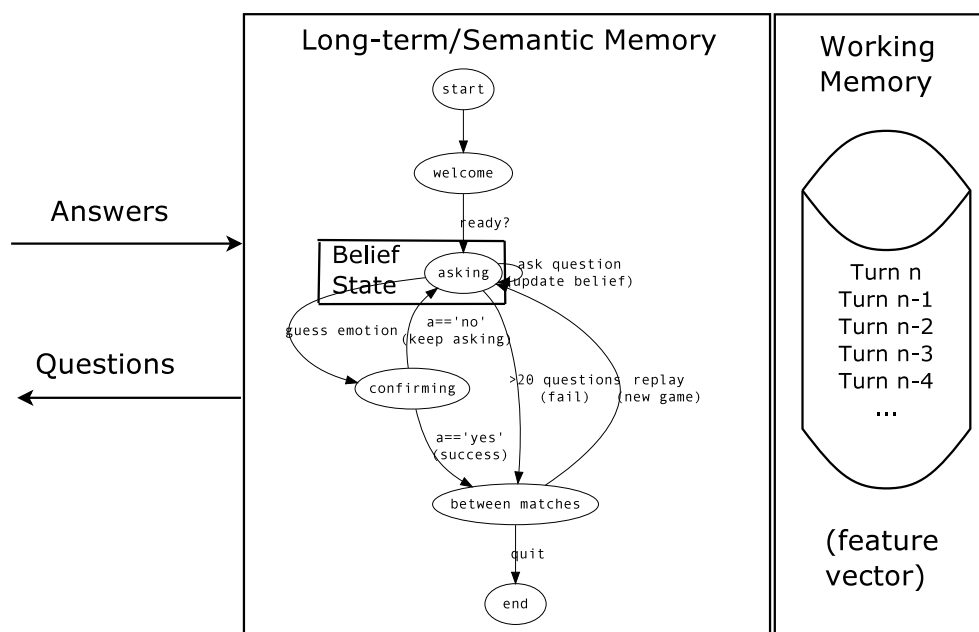


Figure 1: Dialog graph for the EMO20Q questioner agent. The loop labelled “asking” represents the functionality described by the sequential Bayesian belief model of Equation 2 and Algorithm 1. The dialog graph is implemented as a generalized pushdown automaton, where the stack represents the agent’s working memory of question-answer turns.

4. Experiments

The EMO20Q experiments we conducted can be partitioned into human-human and human-computer experiments. Section 4.1 will examine the data from human-human experiments, which was the initial corpus used to train the EMO20Q question-asking agent. Section 4.2 will focus on experiments with the question-asking agent described in Section 3.

4.1. Human-Human EMO20Q

The human-human EMO20Q results are described in an earlier conference paper (Kazemzadeh et al., 2011) but we include a brief description because it is important for understanding the development of the automated agent.

We collected a total of 110 matches from 25 players in the human-human experiments in which EMO20Q was played over text chat. The EMO20Q experiment was implemented as an online chat application using the Extensible Messaging and Presence Protocol (XMPP) and logged so that the games can be easily recorded and studied.

Early in our pilot studies, we realized that it was difficult to successfully terminate the game when the questioner guessed words that were synonyms of the that word the answerer picked. This led us to treat the phenomenon of synonyms with an additional rule that allowed the game to terminate if the answerer

could not verbally explain any difference between the two words. In this case, we considered the game to terminate successfully, but we flagged these matches and kept track of both words.

Of the 110 matches played between the 25 human players, 94 – approximately 85% – terminated successfully with the questioner correctly identifying the emotion that the answerer picked or a word that the answerer felt was a synonym. The mean and median number of questions asked per game was 12.0 and 10, respectively, when failures to correctly guess the emotion were averaged in as 20 questions.

Of the 94 successfully terminated matches, 22 terminated with synonyms. The 16 unsuccessfully terminated matches that were considered failures consisted of several distinct cases. The questioner player could give up early if they had no clue (5/16), they could give up at twenty questions (1/16), or they could pass twenty questions due to losing count or as a matter of pride (6/16). The four remaining cases were considered failures because the answerer inadvertently gave away the answer due to a typing error or giving an unduly generous hint.

There were 71 unique words that players chose in the human-human games, 61 of which were correctly identified. These are listed in Table 2.

There was a total of 1228 question-asking events. Of the questions, 1102 were unique (1054 after normalizing the questions for punctuation and case). In Table 3 we list some of the questions that occurred more than once. Since the surface forms of the questions vary widely, we used manual preprocessing to standardize the questions to a logical form that is invariant to wording. This logical form converted the surface forms to a pseudo-code language by converting the emotion names to nouns if possible, standardizing attributes of emotions and the relations of emotions to situations and events. Examples of the standardized questions are shown in Table 4. After this semantic standardization, there were a total of 727 question types.

4.2. Human-Computer EMO20Q

Using the human-human data described earlier in Section 4.1 and the computational model and algorithm described in Section 3, we built a computer agent to play the questioner role in EMO20Q games. The EMO20Q dialog agent was implemented using a server-side web application that maintained the belief state and episodic buffer for each open connection. The belief state was serialized to EmotionML (Schröder et al., 2012; Burkhardt et al., 2014) and saved in a session database between each question-answer turn.

To test the proposed model of Socratic epistemology for verbal emotional intelligence, we conducted two experiments to assess the performance of the agent. The first experiment was a small pilot study of 15 subjects who played three matches against the agent (Kazemzadeh et al., 2012). In the pilot study, the subjects were recruited locally. Subjects were asked to pick three emotion words, one that they thought was “easy”, one that was “medium”, and a third that was “difficult”. These difficulty ratings were described in terms of a person’s maturity and vocabulary: an “easy” emotion word was one that a child could guess, whereas a “difficult” word was one that would require maturity and a sophisticated vocabulary to guess. The

pilot study was designed to assess the feasibility of the agent design but did not use training beyond the original human-human data.

The second experiment was a larger experiment that forms the key experimental contribution reported by this paper. It followed the same methodology as the pilot study, but with 101 subjects recruited from Amazon Mechanical Turk. These subjects were selected to come from the United States, speak English fluently, and have high past acceptance rates as Mechanical Turkers.

In the second experiment, the parameters of the model were updated every ten subjects. Thus, there were ten waves of ten subjects, each playing 3 matches against the automated agent, which yielded 300 matches. After each ten subjects, the model described in Section 3 was updated based on the total counts of the corpus to that point. In addition to updating the probabilities of the models semantic knowledge (likelihoods), new vocabulary items were added if encountered.

5. Results

The results of our pilot experiments on fifteen subjects are summarized in Table 5. To compare the agent's performance with human performance, we used two objective measures and one subjective measure. The success rate, shown in column two of Table 5, is an objective measure of how often the EMO20Q matches ended with the agent successfully guessing the user's emotion. The number of turns it took for the agent to guess the emotion is the other objective measure. The last column, naturalness, is a subjective measure where users rated how human-like the agent was, on a 0-10 scale.

In the pilot study, the agent obtained a performance of 44% successful outcomes (where the emotion word was correctly guessed). This performance was much less than in the human-human experiments, where successful outcomes occurred in 85% of EMO20Q matches. However, the results indicated that this performance was due to sparsity of data. The emotion words chosen by the subjects as "easy" were recognized by the agent with similar success rate and number of required turns as human-human matches. Some examples of "easy" emotions are anger, happiness, and sadness. However, successful outcomes were fewer in emotions chosen as "medium" and "difficult". Some examples of "medium" emotions are contentment, curiosity, love, and tiredness. Pride, frustration, vindication, and jealousy are examples of "difficult" emotions. Overall, 28 new emotion words were encountered in the pilot study.

The results in terms of successful outcomes and number of turns required to guess the emotion word are roughly reflected in the percent of words that are in-vocabulary. Despite the low performance on emotion words rated "medium" and "difficult", there was not a corresponding decrease in the perceived naturalness of the questioner agent. This led us to believe that the model could reproduce somewhat natural behavior, but that the data we had was insufficient due to the amount of out-of-vocabulary words in the medium and difficult classes, which motivated us to perform the second, larger-scale experiment with 100 players from

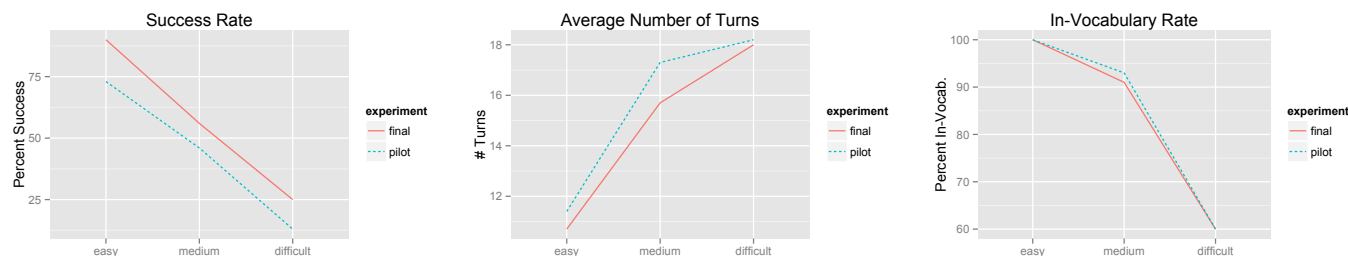


Figure 2: Results of initial automated agent pilot compared to the final experiment of 300 matches on Mechanical Turk, in which the agent was retrained every 30 matches.

Mechanical Turk.

In the larger scale Mechanical Turk experiment, we aimed to improve performance by retraining the model after each batch of 10 subjects. This strategy did in fact increase the successful outcome rate and reduced the length of the EMO20Q dialogs (number of questions), as can be seen from comparing Tables 5 and 6, which are visualized in Figure 2. Across all three difficulty classes, the successful outcome rate improved. The “difficult” class had the largest relative improvement in successful outcomes, increasing from 13% to 25%, and the overall successful outcome increased from 44% to 57%. The lengths of the EMO20Q dialogs decreased most for the medium difficulty class, resulting in an average of 1.6 less turns for this class. Overall, the decrease in dialog length decreased from 15.6 to 14.8 turns.

One surprising result was that even after collecting data from 300 EMO20Q dialogs (more than doubling the earlier human-human data), the out-of-vocabulary rate stayed nearly the same. We had expected out-of-vocabulary-words to become fewer as more data had been seen. However, with each round of the Mechanical Turk experiment, we continued to receive new emotion words rather than converging to a closed vocabulary. For the Mechanical Turk experiment, we did not ask subjects about the perceived naturalness of the agent in order to save on time, and hence costs to pay the Turkers, so unfortunately we cannot say whether the perceived naturalness increased.

Of the 101 subjects, only one was rejected, due to misunderstanding the task by choosing the words “easy”, “medium”, and “difficult” instead of emotion words. This level of acceptance, approximately 99% is rather high for Mechanical Turk, showing a high degree of cooperation. Several users commented that we could have paid less because the task was fun.

A complete listing of the words chosen by the subjects of the experiment is given in Table 7. It can be seen that there are a wide variety of words. A few (those marked by “?”) were questionable in the authors’ intuitions, but otherwise the words showed a high level of understanding and cooperation by the Mechanical Turkers. The three difficulty classes of words were not disjoint: some words like anger, disgust, love, and confusion spanned several categories. It can be concluded that these three difficulty levels do not form a

precise, natural classes of emotion words, but the levels do show a trend toward a smaller basic vocabulary and a wider open vocabulary. The difficulty levels also served as a method to elicit diverse words. The original human-human dialogs identified 71 unique emotion words, after the pilot study there were unique 99
400 emotion words, and after the large-scale mechanical Turk experiment there were 180 unique emotion words.

6. Discussion

The human-human EMO20Q data abounds in highly nuanced natural language descriptions of emotion. For example, one human-human EMO20Q game ended with a discussion of whether “pride” and “proud” refer to the same emotion:

405 [regarding “proud” vs. “pride”] because my intuition was that they’re different... you know
pride sometimes has a negative connotation

In another human-human EMO20Q dialog, a player had difficulty answering whether “anger” was a negative emotion:

[questioner:] so is it a negative emotion?

410 [answerer:] sort of, but it can be righteous

In one human-computer game, one player differentiated the emotion of loving from the emotion of being loved and another player picked the emotion “maudlin”, which the authors needed to look up in a dictionary.

Given the highly nuanced, idiosyncratic descriptions in the human-human data, we were surprised at the amount of successful game outcomes in the human-human EMO20Q games and we were initially unsure
415 whether devising an automated agent would be feasible. Although analyzing this level of detail is beyond the scope of many current systems, we saw that it is a task that humans can do with high success rates. In fact the successful outcome rates in the human-human EMO20Q games are comparable to agreement rates on emotional annotations at a much coarser level, such as labeling data with nine basic emotion labels (Busso et al., 2008).

420 The human-computer results showed us that it possible for computer agents to perform well at the questioner role of EMO20Q and moreover that the agent can learn new vocabulary items and improve its performance past the human-human bootstrap data. The fully trained agent successfully completed 57% of the EMO20Q games, which is 67% of human-human performance and 30% better than the bootstrapped agent. The agent’s emotion word vocabulary nearly doubled after the mechanical Turk experiment. Normally
425 larger emotion vocabularies results in less agreement in annotation tasks but this showed that in the EMO20Q dialog task, vocabulary size is not a weakness but rather a strength. Even when the agent fails to guess the human opponent’s emotion word in the EMO20Q game, the agent’s behavior of searching for knowledge

makes it appear human-like, which enables the agent maintain user engagement and learn from new, out-of-vocabulary words.

The ground truth issue involved in annotating recorded data with descriptive labels is a challenge that the Socratic epistemology can shed light on. The traditional annotation task seeks to have human annotators assign one of a number of labels to data. In the case of emotion research, usually the labels are a controlled vocabulary of several emotion descriptors, like “angry”, “happy”, “sad”, “disgusted”, “fearful”, “surprised”, and “neutral”. The problem with this approach is that these labels often do not fit realistic emotional data. Theoretically, our approach addresses the issue of ground truth in the annotation task with the notion of epistemology, which frames the issue as justification of belief rather than ground truth. Practically, our approach addresses the issue of non-prototypical emotions by enabling a more nuanced representation where the description is not a small, closed set of alternatives but rather an interactive process of communication over a large, open set of natural language descriptions. Though this more nuanced view brings with it new challenges, we have shown the design of an intelligent dialog agent is a feasible way of dealing with these challenges.

We plan to further continue this research in several ways. First, we hope to see the effect of modality on how people describe emotions in natural language. The current work was limited to text-based chat, so the paralinguistic data that may help to convey emotional information was minimized. Including audio and video data may allow greater convergence of the players to agree upon the unknown emotion in EMO20Q. Another area of future research will be to model the answerer role. The current research focused on the questioner role, but the answerer role will offer additional challenges and insights. In particular, automating the answerer role will require more robust natural language understanding because it will need to process to new, unseen questions from users, whereas the questioner used a fixed set of questions and only had to process answers to yes/no questions. The answerer would also likely require a different model than the Socratic, question-asking model presented in this paper. A successful answerer agent would allow a pleasing closed-loop simulation where both roles of EMO20Q are played by computer. There are also further areas to explore for the questioner agent, in particular, the criterion for choosing each question. Finally, we think that this approach can improve emotion annotation and other annotation tasks, such as coding behavioral data for psychological assessment. In these tasks human annotators are asked to label data using a controlled vocabulary of words and agreement is established statistically between isolated annotators. However, we have shown that humans are able to communicate with high accuracy using a large, subjective vocabulary and we feel that allowing natural language descriptions in an interactive, question-asking setting will allow for more accurate and less constrained annotations.

7. Conclusion

The main goals of this paper were to formulate a theoretical and computational model for a subset of human emotional language. We called this model the Socratic epistemology for verbal emotional intelligence because uses question-asking to justify beliefs about emotions in a natural language dialog context. We presented the emotion twenty questions (EMO20Q) game and showed that the level of human performance was high despite not limiting the players to any predefined emotion vocabulary. We also presented an automated agent that can play the question-asking role of EMO20Q. This agent uses a sequential Bayesian belief update algorithm to simulate a cognitive processing by which the agent updates its belief state of candidate emotion words over time. This framework was inspired by a method of question-asking that was proposed by the ancient philosopher Socrates and the field of epistemology:

[Gorgias:] Just as different drugs draw forth different humors from the body – some putting a stop to disease, others to life – so too with words: some cause pain, others joy, some strike fear, some stir the audience to boldness, some benumb and bewitch the soul with evil persuasion” (Gorgias, *Encomium of Helen*, c.415 B.C.).

Socrates: You, Gorgias, like myself, have had great experience of disputations, and you must have observed, I think, that they do not always terminate in mutual edification, or in the definition by either party of the subjects which they are discussing;. . . Now if you are one of my sort, I should like to cross-examine you, but if not I will let you alone. And what is my sort? you will ask. I am one of those who are very willing to be refuted if I say anything which is not true, and very willing to refute any one else who says what is not true, and quite as ready to be refuted as to refute. (Plato, *Gorgias*, 380 B.C.)

In the first quote above, Gorgias, a Sophist rhetorician, describes the effects of words on a person’s emotions. Gorgias describes emotions by making reference to the theory of physiological humors. Humankind’s conception of emotions has changed since the time of the ancients, who believed that emotions were generated from bodily “humors”, which in turn were derived from alchemical elements, but our conception of emotion is still largely expressible through language.

In the second quote, Socrates (as quoted by Plato) cross-examines Gorgias to determine Gorgias’ beliefs. Socrates applied his method of question-asking to understand beliefs about complex abstract concepts that were disputed in ancient times. Two millenia later we have used a computational implementation of this method to make a dialog agent better understand human beliefs about emotional concepts.

We have provided an anonymized version of data we gathered from EMO20Q, source code for the experiments, demos, and other resources at <http://sail.usc.edu/emo20q> .

8. References

Cyril Allauzen and Michael Riley. A pushdown transducer extension for the openfst library. In *Proceedings of the Conference on Implementation and Application of Automata*, 2012.

Allan D. Baddeley. The episodic buffer: A new component of working memory? *Trends in Cognitive Science*, 4(11):417–423, 2000.

Lisa Feldman Barrett. Are emotions natural kinds? *Perspectives on Psychological Science*, 1(1):28–58, March 2006.

Felix Burkhardt, Christian Becker-Asano, Edmon Begoli, Roddy Cowie, Gerhard Fobe, Patrick Gebhard, Abe Kazemzadeh, Igmarr Steiner, and Tim Llewellyn. Application of emotionml. In *Emotion, Social Signals, Sentiment, and Linked Open Data (ES3LOD) 2014*, 2014.

Carlos Busso, Murtaza Bulut, Chi-Chun Lee, Abe Kazemzadeh, Emily Mower, Samuel Kim, Jeannette Chang, Sungbok Lee, and Shrikanth Narayanan. IEMOCAP: Interactive emotional dyadic motion capture database. *Journal of Language Resources and Evaluation*, 42(4):335–359, 2008.

Valeria Carofiglio, Fiorella de Rosis, and Nicole Novielli. *Cognitive Emotion Modeling in Natural Language Communication*, pages 23–44. Springer, 2009.

Umberto Eco and Thomas A. Sebeok, editors. *The Sign of Three: Dupin, Holmes, Peirce*. Advances in Semiotics. Indiana University Press, 1988.

Brandy N. Frazier, Susan A. Gelman, and Henry M. Wellman. Preschoolers’ search for explanatory information within adult-child conversation. *Child Development*, 80(6):1592–1611, November/December 2009.

Thomas L Griffiths, Charles Kemp, and Joshua B Tenenbaum. Bayesian models of cognition. *Cambridge University Press*, 2008.

Jaakko Hintikka. *Socratic Epistemology: Explorations of Knowledge-Seeking by Questioning*. Cambridge University Press, 2007.

Charles F. Hockett and Stuart Altmann. *A note on design features*, pages 61–72. Indiana University Press, 1968.

Jeff Howe. The rise of crowdsourcing. *Wired Magazine*, 14.06, June 2006.

Sep Kamvar and Jonathan Harris. *We Feel Fine: An Almanac of Human Emotion*. Scribner, 2009.

Abe Kazemzadeh. *Natural Language Description of Emotion*. PhD thesis, University of Southern California, 2013.

Abe Kazemzadeh, Panayiotis G. Georgiou, Sungbok Lee, and Shrikanth Narayanan. Emotion twenty questions: Toward a crowd-sourced theory of emotions. In *Proceedings of ACHI'11*, 2011.

Abe Kazemzadeh, James Gibson, Juanchen Li, Sungbok Lee, Panayiotis G. Georgiou, and Shrikanth Narayanan. A sequential Bayesian agent for computational ethnography. In *Proceedings of Interspeech*, Portland, OR, October 2012.

Brian King. *The Conceptual Structure of Emotional Experience in Chinese*. PhD thesis, Ohio State University, 1989.

Rada Mihalcea and Hugo Liu. A corpus-based approach to finding happiness. In *AAAI Spring Symposium on Computational Approaches to Weblogs*, March 2006. URL <http://www.cse.unt.edu/~rada/papers/mihalcea.aaais06.pdf>.

Emily Mower, Angeliki Metallinou, Chi-Chun Lee, Abe Kazemzadeh, Carlos Busso, Sungbok Lee, and Shrikanth Narayanan. Interpreting ambiguous emotional expressions. In *ACII Special Session: Recognition of Non-Prototypical Emotion from Speech- The Final Frontier?*, Amsterdam, Netherlands, 2009.

Charles Sanders Peirce. Some consequences of four incapacities. *Journal of Speculative Philosophy*, 2:140–157, 1868. URL http://www.iupui.edu/~peirce/writings/v2/w2/w2_22/v2_22.htm.

Edmond T. Rolls. *What Are Emotions, Why Do We Have Emotions, and What Is Their Computational Basis in the Brain*, chapter 5, pages 117–146. Oxford University Press, 2005.

Bertrand Russell. On denoting. *Mind*, 14:479–493, 1905.

Marc Schröder, Paolo Baggia, Felix Burkhardt, Catherine Pelachaud, Christian Peter, and Enrico Zovato. W3C candidate recommendation: Emotion markup language (EmotionML) 1.0, May 2012. URL <http://www.w3.org/TR/emotionml/>. <http://www.w3.org/TR/emotionml/>.

Phillip R. Shaver, Upekkha Murdaya, and R. Chris Fraley. Structure of the Indonesian emotion lexicon. *Asian Journal of Social Psychology*, 4:201–224, 2001.

Mark Steyvers, Thomas L Griffiths, and Simon Dennis. Probabilistic inference in human semantic memory. *Trends in Cognitive Sciences*, 10(7):327–334, 2006.

Luis von Ahn and Laura Dabbish. Labeling images with a computer game. In *Proceedings of the SIGCHI conference on Human factors in computing systems*, 2004.

Fei Xu and Joshua B Tenenbaum. Word learning as bayesian inference: Evidence from preschoolers. In *Proceedings of the Twenty-Seventh Annual Conference of the Cognitive Science Society*, 2005.

Ning Zhong, Jiming Liu, Y.Y. Yao, and S. Ohsuga. Web intelligence. In *Computer Software and Applications Conference*, 2000.

Zoltán Kövecses. *Metaphor and Emotion: Language, Culture, and Body in Human Feeling*. Cambridge University Press, 2000.

Table 1: Example dialog from EMO20Q questioner agent where the human interlocutor has chosen “relaxation”.

Speaker/Event	Utterance/Content
Agent	Let me know when you are ready...
User	ok
Belief update	(Uniform distribution)
Agent	Is it very intense?
User	no
Belief update	relaxation:0.027, hope:0.027, boredom:0.027, ..., anger:0.0024
Agent	is it positive?
User	yes
Belief update	hope:0.15, relaxation:0.12, cheerfulness:0.091, ...
Agent	is it an emotion that just pops up spontaneously (vs being triggered by something)?
User	depends
Belief update	hope:0.33, relaxation:0.23, cheerfulness:0.15, ...
Agent	is it caused by something in particular?
User	well, yes, you do need the right conditions
Belief update	relaxation:0.50, hope:0.29, cheerfulness:0.097, ...
Agent	Is it relaxation?
User	Yes
Agent	so did I get it right?
User	Yes, you got it
Agent	Awesome! Would you like to play again?
User	No, thanks.
Agent	Thank you for playing

Table 2: Emotion words from human-human EMO20Q matches. Matches that terminated with words that the answerers felt to be synonyms are denoted by parentheses.

emotions (synonyms)	count	# correct	...	emotions (synonyms)	count	# correct
admiration	1	1		guilt	4	4
adoration	1	0		happiness	1	1
affection (love)	2	2		helplessness	1	1
amusement	1	1		hope (feeling lucky)	3	3
anger	2	1		insecurity (shyness)	1	1
annoyance (irritated)	2	2		jealousy (envy)	3	3
anxiety	3	3		joy	1	0
apathy (uninterested)	1	1		loneliness	1	1
awe	1	0		love	2	2
boredom	2	2		madness (anger)	1	1
bravery	1	1		melancholy	1	1
calm	2	2		pity (sympathy)	1	1
cheerfulness	1	1		pride	2	2
confidence	1	1		proud	1	1
confusion	2	1		regret	2	2
contempt	1	1		relief	5	5
contentment (calm)	2	1		sadness	2	2
depression (misery)	2	2		satisfaction	1	0
devastation	1	0		serenity	1	1
disappointment	1	1		shame	1	1
disgust	2	2		shock	1	1
dread (hopelessness)	1	1		shyness	1	1
eagerness (determination)	1	1		silly	1	1
embarrassment	2	2		soberness	1	0
enthusiasm (eagerness)	3	1		sorrow (sadness)	1	1
envy (jealousy)	3	3		stress	1	1
exasperation	1	1		suffering	1	0
excitement	1	1		surprise	3	3
exhilaration (thrill)	1	1		tense (uncomfortable)	1	0
exhaustion	1	1		terror	1	1
fear (distress, scared)	2	2		thankful	1	0
frustration	2	2		thrill (entrancement)	2	1
fury	1	1		tiredness	2	2
glee	1	0		wariness	1	0
gratefulness	1	1		worry (anxiety, scared)	3	3
grumpiness	1	1				
				total	110	94

Table 3: Examples of some of the questions that occurred multiple times (disregarding case and punctuation).

question	count
is it positive?	16
ok is it a positive emotion?	15
is it a positive emotion?	14
is it intense?	13
ok is it positive?	10
is it a strong emotion?	7
is it like sadness?	6
is it sadness?	5
is it pride?	5
is it neutral?	5
is it like anger?	5
is it surprise?	4
is it an emotion that makes you feel good?	4
thrilled?	3
regret?	3
pleased?	3
is it very intense?	3
is it love?	3
is it kinda like anger?	3
is it associated with sadness?	3
...	...
ok is it a negative emotion?	2
ok is it a good emotion?	2
okay is it a strong emotion?	2
is it highly activated?	2
is it directed towards another person?	2
is it directed at another person?	2
is it associated with satisfaction?	2
is it associated with optimism?	2
is it associated with disappointment?	2
is it an emotion that lasts a long time	2
does it vary in intensity?	2

Table 4: Examples of question standardization.

Standardized Question	Examples
cause(emptySet,e)	<i>can you feel the emotion without any external events that cause it?</i>
cause(otherPerson,e)	<i>is it an emotion that just pops up spontaneously (vs being triggered by something)?</i> <i>is it caused by the person that it's directed at?</i>
e.valence==negative	<i>Do you need someone to pull this emotion out of you or evoke it? if so, who is it?</i> <i>is it considered a negative thing to feel?</i> <i>2) so is it a negative emotion?</i>
situation(e,birthday)	<i>would you feel this if it was your birthday?</i> <i>is it a socially acceptable emotion, say, at a birthday party?</i>
e==frustration	<i>oh, is it frustrated?</i> <i>frustration?</i>

Table 5: Experimental results for 15 subject pilot study (45 EMO20Q games).

difficulty	% success	avg. turns	% in vocab.	naturalness
easy	73%	11.4	100%	6.9
medium	46%	17.3	93%	5.5
difficult	13%	18.2	60%	5.8
total	44%	15.6	84%	6.1

Table 6: Experimental results for 100 subject Mechanical Turk study (300 EMO20Q games).

difficulty	% success	avg. turns	% in vocab.
easy	90%	10.7	100%
medium	56%	15.7	91%
difficult	25%	18.0	60%
total	57%	14.8	83.7%

Table 7: Observed emotion words by difficulty. Words that were attested but which did not fit the authors' broad intuitions are marked with '?'. The same words in multiple categories indicate that different subjects had differing opinions about difficulty.

difficulty	examples
easy	happiness, anger, sadness, calm, confusion, love, mad, hate, joy
medium	anger, confusion, contentment, curiosity, depression, disgust, excitement, fear, hate, irritation, love, melancholy, sorrow, surprise, tiredness, envy, outrage, elation, suffering, jealousy, nervousness, sympathy, thrill, upset, joy, anxiety, frustration, flustered, enjoyment, exhaustion, fury, boredom, delight, cold, apathy, hostility, loved, annoyance, playfulness, downtrodden, stupor, despair, pissed, nostalgia, overjoyed, indifference, courage
difficult	devastation, disgust, ecstasy, ennui, frustration, guilt, hope, irritation, jealousy, morose, proud, remorse, vindication, zealousness, elation, mischievous, usure, angst, patience, despise, inspired, euphoria, exuberance, worrying, melancholy, ambivalence, love, loneliness, exacerbated(?), avarice, stress, envy, disillusionment, maudlin, depression, confusion, maniacal, ambiguity, concern, pleasure, shame, indifference, anger, suicidal, pessimism, annoyance, sense of failure, educated(?), manic, overwhelmed, astounded, discontent, energetic, introspective, appalled, serenity, dissatisfaction, anxiety, lust, conflicted, perplexed, jubilation, disappointment, satisfaction, remorse, embarrassment, downcast, guilty, enamored, alienation, exotic(?), hate, caring, resentment, pity, aversion, quixotic, infuriation