# A Taxonomy of Empathetic Response Intents in Human Social Conversations

**Anuradha Welivita and Pearl Pu**

School of Computer and Communication Sciences
École polytechnique fédérale de Lausanne
Switzerland
`{kalpani.welivita,pearl.pu}@epfl.ch`

## Abstract

Open-domain conversational agents or chatbots are becoming increasingly popular in the natural language processing community. One of the challenges is enabling them to converse in an empathetic manner. Current neural response generation methods rely solely on end-to-end learning from large scale conversation data to generate dialogues. This approach can produce socially unacceptable responses due to the lack of large-scale quality data used to train the neural models. However, recent work has shown the promise of combining dialogue act/intent modelling and neural response generation. This hybrid method improves the response quality of chatbots and makes them more controllable and interpretable. A key element in dialog intent modelling is the development of a taxonomy. Inspired by this idea, we have manually labeled 500 response intents using a subset of a sizeable empathetic dialogue dataset (25K dialogues). Our goal is to produce a large-scale taxonomy for empathetic response intents. Furthermore, using lexical and machine learning methods, we automatically analysed both speaker and listener utterances of the entire dataset with identified response intents and 32 emotion categories. Finally, we use information visualization methods to summarize emotional dialogue exchange patterns and their temporal progression. These results reveal novel and important empathy patterns in human-human open-domain conversations and can serve as heuristics for hybrid approaches.

## 1 Introduction

Inspired by the recent success of deep neural networks for natural language processing (NLP) tasks such as language modeling (Mikolov et al., 2010) and machine translation (Sutskever et al., 2014), neural response generation is currently at the forefront of research in the NLP community. Recent advances in this field have proven the efficacy of deep neural networks in modelling both task-oriented and open-domain dialogue systems (Wen et al., 2015; Sutskever et al., 2014; Vinyals and Le, 2015). Most of the existing neutral conversation models are capable of generating syntactically and contextually well-formed responses. Some of the work also focuses on enabling chatbots to generate emotionally colored and affect-rich responses (Asghar et al., 2018; Zhou et al., 2018; Xie et al., 2019). Despite the efforts in modeling affect in natural language, work that focuses specifically on modeling empathy in chatbots is relatively limited and remains an open research question (Spring et al., 2019).

Empathy plays a vital role in human psychological processes for smooth social interaction (Decety, 2010). Empathy-related responding includes caring and sympathetic concerns for other people. Humans are born with core *affect* neural circuitry, and they gradually develop the ability to apprehend the emotional states of others and respond in an empathetic manner. Empathy motivates pro-social behavior and increases the sense of social bonding (Eisenberg and Eggum, 2009). Therefore, in the context of social interaction, a chatbot needs to be empathetic to maintain healthy interaction with humans and develop trust. The task of augmenting social chatbots with empathy is challenging because the generated responses have to be appropriate in terms of both content and emotion information (Spring et al., 2019).

Several neural response generation models have attempted to address this challenge in a fully data-driven manner. For example, Rashkin et al. (2019), use the full transformer architecture (Vinyals and Le, 2015) pre-trained on 1.7 billion Reddit conversations and fine-tuned on the EmpatheticDialogues dataset (Rashkin et al., 2019) to generate empathetic responses. Lin et al. (2019) adapt the Generative Pretrained Transformer (GPT) (Radford et al., 2018) to empathetic response generation task by fine-tuning it on the PersonaChat (Zhang et al., 2018) and EmpatheticDialogues datasets. Even though these models are capable of mimicking human empathetic conversation patterns in some ways, it is often unpredictable what the chatbots might generate, for example, they may generate inconsiderate remarks, redundant responses, asking the same questions repeatedly, or any combinations of them. Since it is really important to respond to humans' emotions appropriately, we believe controllability of response generation is essential.

Several other neural response generation approaches attempt to gain control over the generated response by conditioning it on a manually specified emotion label (Zhou et al., 2018; Zhou and Wang, 2018; Hu et al., 2018; Song et al., 2019) or using affective loss functions based on heuristics such as minimizing or maximizing affective dissonance between prompts and responses (Asghar et al., 2018). These models claim to generate emotionally more appropriate responses than those generated from purely data-driven models. However, the primary concern of these handcrafted rules is its practicality. No prior work has shown normative associations between the speaker's emotions and the corresponding listener's emotions. As our work would reveal, listeners are much more likely to respond to sad or angry emotions with questioning than expressing similar or opposite emotions in the first turn. Xu et al. (2018), however, has shown the benefit of incorporating dialogue acts as policies in designing a social chatbot. They were able to avoid the need to manually condition the next response with a label by jointly modeling dialogue act selection and response generation. Their framework first selects a dialogue act from a policy network according to the dialogue history. The generation network then generates a response based on both dialogue history and the selected dialogue act. It is thus possible to explicitly learn human-human conversational patterns in social chitchat and generate more controlled and interpretable responses. Unfortunately, they did not study empathetic response generation.

To fill this gap, we have developed a taxonomy of empathetic listener intents by manually annotating around 500 utterances of the EmpatheticDialogues dataset (Rashkin et al., 2019), covering 32 types of emotion categories. In the following, we first describe in detail how this taxonomy was derived (Figure 1) and how we chose the dataset to support this annotation work. To extend this subset, we employ automatic techniques to label all speaker and listener utterances, covering 25k empathetic human-human conversations. To be able to explain the patterns and trends of the conversation flow, we employ visualization methods to illustrate the most frequent exchanges and reveal how they temporally vary as dialogues proceed. Finally, we discuss how these results can be used to derive more informed heuristics for controlling the neural response generation process.[1]

## 2   Related Work

To provide a background for this research, we begin by describing some of the existing theories related to empathy in other fields such as psychology and neuroscience and their limitations in incorporating them into the design of social chatbots. Then we describe seminal work on existing neural-based open-domain response generation systems and means by which they control the generated response. These studies serve as the motivation and inspiration for our work. Next, we discuss some existing dialogue-act/intent taxonomies and their limitations in modeling empathy in human social conversations.

### 2.1   Theories of Empathy in Psychology and Neuroscience

Zillmann (2008) defines empathy as a social emotion in response to the emotions of others. Further, he states that the evoked empathetic reaction itself constitutes an emotional experience, primarily because it is associated with increased excitement and awareness. It tends to be a *feeling with* or *feeling for* the observed party. Two of the most famous theories in the psychological literature to explain the phenomenon

---

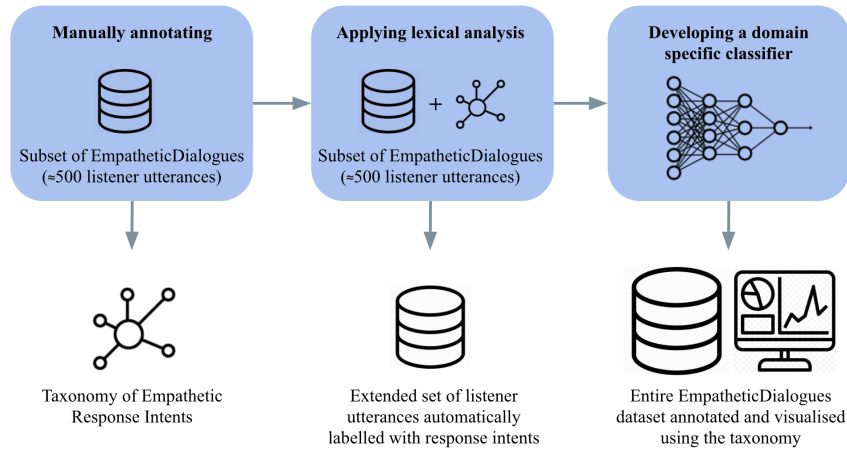[1]Our source code and results are available at `https://github.com/anuradha1992/EmpatheticIntents`.

Figure 1: Three development steps for constructing the taxonomy of empathetic response intents.

of empathy are the "simulation theory" (Gordon, 1995) and the "theory-theory" (Gopnik and Wellman, 1995). The "simulation theory" states that a person understands another or empathizes by imagining himself in the other's situation and seeing it from his perspective. The "theory-theory" states that the ability to understand what another person is feeling is based on rules for how one should think and feel.

A recent work by Singer and Klimecki (2014) in the field of Neuroscience states that empathy refers to the general capacity of humans to resonate with others' emotional states irrespective of their valence. However, when empathizing, they suggest that one should not confuse oneself with the other; i.e., one should still know that the emotion he resonates with is the emotion of another. The failure to separate that can lead to empathetic distress. According to them, the desirable way of empathizing with others is having compassion, which is a feeling of concern for another person's suffering accompanied by the motivation to help. However, all this work does not describe in detail specific means through which humans show empathy, especially via natural language dialogues. Also, most of these studies on empathic states focus on reactions to negative rather than positive events (Buechel et al., 2018). Hence, empathy for positive events remains less understood. In our study, we explore the means through which humans empathize with others both in positive and negative scenarios.

## 2.2 Neural Response Generation

Xie et al. (2019), describe an end-to-end Multi-turn Emotionally Engaging Dialog model (MEED), capable of recognizing emotions and generating emotionally appropriate and human-like responses. Their GRU based Seq2Seq dialogue model consists of a hierarchical mechanism to track the conversation history in multi-turn dialogues combined with an additional emotion RNN to process the emotional information in each history utterance. They model affect exchanges in human dialogues using a dedicated embedding layer. This emotion recognition step enables the model to produce more emotionally appropriate responses for a given context. But their approach is entirely data-driven and lacks control and interpretability over generated responses.

Chen et al. (2019), propose a model that can generate comments to posts in social media so that they are not only relevant in topic but also in emotion. To fully understand how emotions are expressed in conversations, they first analyse NTCIR-12 STC-1 collection (Shang et al., 2016), a social-media conversation dataset. The results show that for posts with different emotions, the distributions of comment emotions are very different from each other, and only several emotions are appropriate for responding to a given post. Inspired by the findings, they extend the basic encoder-decoder neural network architecture (Vinyals and Le, 2015) with an RNN-based response emotion estimator, which takes in a post and estimates how relevant an emotion is for responding to the post. This information is fed into the decoder when generating the response. In this, the classifier automatically determines the emotion of the response. Xu et al. (2018) incorporate dialogue acts as policies in their open-domain neural response generation model by performing learning with human-human conversations tagged with a dialogue act

classifier. They jointly model dialogue act selection and response generation using a GRU based neural network consisting of a policy network and a generation network. The policy network first selects a dialogue act according to the conversation history, and then the generation network generates a response based on the conversation history and the selected dialogue act. They claim that with dialogue acts, they not only achieve significant improvement over response quality for a given context but also can explain why such achievements are possible. The above work motivated us to develop explicit empathetic response intents from the dataset. We believe they can inform the development of empathetic social chatbots by providing more control to and interpretation of the responses generated and render human-machine conversations more natural and engaging.

## 2.3 Dialogue-Act/Intent Taxonomies

Work has been conducted to establish dialogue act/intent taxonomies both by analysing human-human and human-machine conversation datasets. Stolcke et al. (2000) propose a taxonomy of 42 mutually exclusive dialogue acts with the intention of enabling computational dialogue act modeling for conversational speech. They follow the standard Dialog Act Markup in Several Layers (DAMSL) tag set (Core and Allen, 1997) and modify it in several ways so they can easily distinguish utterances in conversational speech. Using this taxonomy, they produce a large hand-labeled database of 1,155 conversations from the Switchboard corpus of spontaneous human-to-human telephone conversations, which is widely used to train and test dialogue act classifiers. Montenegro et al. (2019) propose a dialogue act taxonomy for a task-oriented virtual coach designed to improve the lives of the elderly. It is a multi-dimensional hierarchical taxonomy comprising of topic, intent, polarity, and entity labels at the top, in which the intent label classifies the utterance in classes related to the user's communicative intentions such as 'question', 'inform', and 'agree'. They use the taxonomy to manually annotate the user turns in 384 human-machine dialogues collected from a group of elderly. It aims to help the dialogue agent to detect goals, realities, obstacles, and ways forward of the particular topics the agent is designed to deal with.

Existing dialogue-act/intent taxonomies are either too general as they were constructed for open-domain conversations or too specific as they were constructed for specific task-oriented scenarios. These taxonomies do not necessarily model empathy in human social conversations. Also, the above approaches do not use automatic approaches to extend the manual annotations, which make their datasets comparatively smaller. In our study, we present how a smaller set of human labeled sentences can be extended using lexical methods and use it to train a classifier to automatically annotate a larger corpus.

## 3 Dataset

Many open-domain conversation datasets are publicly available mainly to assist tasks such as neural dialogue generation. Out of them, some datasets are multi-modal (e.g. IEMOCAP (Busso et al., 2008), SEMAINE (McKeown et al., 2011), MELD (Poria et al., 2018)) containing visual, acoustic and textual signals. Since they contain a lot of back-channel communication through facial expressions and speech tones, the text may not fully represent the contextual expression of intent. Datasets containing dialogues extracted from social media platforms such as Twitter (e.g., the Twitter Dialog Corpus (Serban et al., 2017)) are often noisy, short, and different from real-world conversations and may contain a lot of toxic responses rather than compassionate ones. Also, datasets containing TV or movie transcripts (e.g., Emotionlines (Chen et al., 2018), OpenSubtitles (Lison et al., 2019)) and telephone recordings (e.g. Switchboard corpus (Stolcke et al., 2000)) are a translation of voice into text, which does not fully model interactions that happen only through text. Even purely text-based daily conversation datasets such as DailyDialog (Li et al., 2017) are not guaranteed to contain empathetic responses.

Rashkin et al. (2019) introduced the EmpatheticDialogues dataset consisting of 24,856 open-domain, human-human conversations as a benchmark dataset to train and evaluate dialogue systems that can converse in an empathetic manner. Each conversation in this dataset is based on a situation associated with one of 32 emotions, which are selected from multiple annotation schemes, ranging from basic emotions derived from biological responses (Ekman, 1992; Plutchik, 1984) to larger sets of subtle emotions derived from contextual situations (Skerry and Saxe, 2015). The dialogues are collected using ParlAI

| | |
|---|---|
| Label: | Afraid |
| Situation: | Speaker felt this when... *"I've been hearing noises around the house at night"* |
| Conversation: | Speaker: *I've been hearing some strange noises around the house at night.* |
| | Listener: *oh no! That's scary! What do you think it is?* |
| | Speaker: *I don't know, that's what's making me anxious.* |
| | Listener: *I'm sorry to hear that. I wish I could help you figure it out* |

Table 1: Example conversation taken from the EmpatheticDialogues dataset.

(Miller et al., 2017), integrated with Amazon Mechanical Turk (MTurk), recruiting 810 US workers. During construction, the workers were instructed to show empathy when responding to conversations initiated by their speaker counterparts. Since almost all the dialogues in this dataset are empathetic, purely text-based, and most of which do not contain any toxic responses, we chose it to derive our taxonomy. An example conversation from this dataset is given in Table 1. Table 2 shows the basic statistics of the dataset. The average number of turns per dialogue is close to 4. The maximum number of dialogue turns in the dataset is 8. However, not many dialogues exceed 4 turns. Close to 77% of the total number of dialogues contain only up to 4 turns, and only 1.4% of the dialogues contain up to 8 turns.

| Criteria | Statistics |
|---|---|
| Total no. of dialogues | 24,856 |
| Total no. of dialogue turns | 107,247 |
| Average no. of turns per dialogue | 4.31 |
| Maximum no. of turns per dialogue | 8 (345 dialogues) |
| Minimum no. of turns per dialogue | 1 (3 dialogues) |
| Total no. of speaker turns | 55,984 |
| Total no. of listener turns | 51,263 |
| Average no. of speaker tokens per dialogue turn | 17.88 |
| Average no. of listener tokens per dialogue turn | 13.69 |

Table 2: Statistics of the EmpatheticDialogues dataset used for analysis.

## 4 Taxonomy of Empathetic Response Intents

To develop the taxonomy, we investigated which intents are frequently associated with listeners when responding to different emotional situations in EmpatheticDialogues. We took a subset of the dataset with situations associated with the Plutchik's 8 basic emotions (Plutchik, 1984) (joyful, anticipating, trusting, surprised, angry, afraid, sad, and disgusted), and manually analysed it to derive the listener intents associated with each type of emotions. In this process, 20 dialogues belonging to each emotion were randomly selected and each sentence in all listener utterances were manually annotated by an expert evaluator with a label that best describes their intent. This resulted in 521 sentences manually annotated with intent labels. Because an utterance can have multiple sentences, we decided to annotate each sentence in a listener's utterance with a unique intent label. For example, the two sentences comprising the utterance *"Those symptoms are scary! Do you think it's Corona?"* would be annotated with separate intent labels "Acknowledging" and "Questioning", respectively. After analysing their occurrences and whether some of the intents can be grouped into a common intent, we were able to come up with a taxonomy of 15 empathetic response intents. Table 3 presents this taxonomy with corresponding examples and occurrence frequencies. Words and phrases that were most helpful in annotating these examples with their corresponding intents are underlined. Manual annotation of empathetic response intents was carried out with reference to the context preceding an utterance. This way, we were able to distinguish utterances using similar words in the same order depending on their context. For example, sentences such as *"I hope they find a vaccine soon."* can be categorised into two different intents, "Encouraging" and "Consoling" depending on whether the sentence follows a positive or negative emotional context, respectively.

| Category | Examples | Frequency |
|---|---|---|
| 1. Questioning (to know further details or clarify) | - *What are you looking forward to?* | 24.38% |
| 2. Acknowledging (Admitting as being fact) | - *That sounds like double good news. It was probably fun having your hard work rewarded.* | 22.46% |
| 3. Agreeing (Thinking/Saying the same) | - *That's a great feeling, I agree!* | 9.60% |
| 4. Consoling | - *I hope he gets the help he needs.* | 7.87% |
| 5. Encouraging | - *Hopefully you will catch those great deals!* | 5.37% |
| 6. Sympathizing (Express feeling pity or sorrow for the person in trouble) | - *So sorry to hear that.* | 5.37% |
| 7. Wishing | - *Hey... congratulations to you!* | 4.41% |
| 8. Suggesting | - *Maybe you two should go to the pet store to try and find a new dog for him!* | 4.03% |
| 9. Sharing own thoughts/opinion | *I would love to have a boy too, but I'm not sure if I want another one or not.* | 4.03% |
| 10. Sharing or relating to own experience | *I had a friend who went through the same thing.* | 3.84% |
| 11. Advising | *Don't take too much money with you.* | 2.69% |
| 12. Expressing care or concern | *I hope the surgery went successfully and with no hassle.* | 2.30% |
| 13. Expressing relief | *Phew.. That's a relief., I am glad you were okay.* | 1.53% |
| 14. Disapproving | *But America is so great now! look at all the great things that are happening.* | 1.15% |
| 15. Appreciating | *You are very trusting. It's nice to have a friend like you.* | 0.95% |

Table 3: Taxonomy of empathetic response intents with corresponding examples and occurrence frequencies based on the manually annotated 521 listener utterances in the EmpatheticDialogues dataset.

## 5 Automatic Labelling of EmpatheticDialogues Using the Taxonomy

### 5.1 Annotation Procedure

To annotate all the speaker and listener utterances in the EmpatheticDialogues dataset with emotion labels and response intents, we trained a BERT transformer-based classifier, as suggested by Devlin et al. (2019). Prior to selecting BERT as the classifier, we trained and tested a FastText classifier on the annotation task, but its accuracy was lower compared to BERT. We proceeded with the 8 most frequent intents (questioning, acknowledging, agreeing, consoling, encouraging, sympathizing, wishing, and suggesting) in our taxonomy of empathetic listener intents and the 32 types of emotion categories given in the EmpatheticDialogues dataset. The rest of the listener intents were classified as 'neutral' since the emotion behind those intents were more on the neutral side. To expand the training data collected by manual annotation, we searched through the rest of the dataset using n-grams that are most indicative of the intent categories. For example, n-grams such as '100 %', 'absolutely', 'definitely', 'i agree', 'me neither', 'me too', and 'i completely understand' are indicative of the intent 'agreeing' and were used to collect more example utterances corresponding to that category. The most indicative n-grams used to collect more utterances for each of the intents are listed in Appendix A.

During training, we initialized the representation network with weights from the pre-trained language model, RoBERTA (Liu et al., 2019), and fine-tuned the model on situation descriptions given in the EmpatheticDialogues dataset tagged with 32 emotions and listener utterances tagged with 8 out of 15 intents

from our taxonomy of empathetic response intents. The training, validation, and test sets comprised of 25023, 3544 and 3225 sentences respectively, which spanned equally across all emotion and intent categories. We trained the model with a peak learning rate of $2e^{-5}$ and a batch size of 32 for 10 epochs and obtained the classifier giving the lowest validation loss. The top-1 accuracy of our classifier with 41 labels over the test set was 65.88%, which is significantly higher than the accuracy of FastText (Joulin et al., 2016) and DeepMoji (Felbo et al., 2017) classifiers trained on 32 emotion labels in the Empathetic-Dialogues dataset. The latter two were considered as the state-of-the-art at that time, and achieved 43% and 48% accuracy on the EmpatheticDialogues test set, respectively (Rashkin et al., 2019).

## 5.2 Analysis of Emotion-Intent Exchange Patterns

Based on the above annotations, we analysed the most frequent response intents corresponding to different emotions expressed by speakers. In Figure 2, we visualize the emotion-intent exchanges taking place between speakers and listeners in the EmpatheticDialogues dataset. In this, each chord connects emotion-intent pairs that co-occur together. The chord leaving a particular arc represents the speaker's emotion or intent and gets connected to the arc representing the listener's emotion or intent that immediately follows. It can be seen that a significant proportion of speaker utterances contain a particular emotion in the 32 different emotion categories defined in EmpatheticDialogues, while most of the listener utterances contain a particular intent defined in our taxonomy. Instead of conveying a particular emotion, the listeners show their empathy via specific means described in our taxonomy. And the proportions of the arcs for each intent resembles the frequencies in the manually annotated subset. It serves as a validation that our taxonomy is indeed true for listener responses as it was applied to the entire dataset.
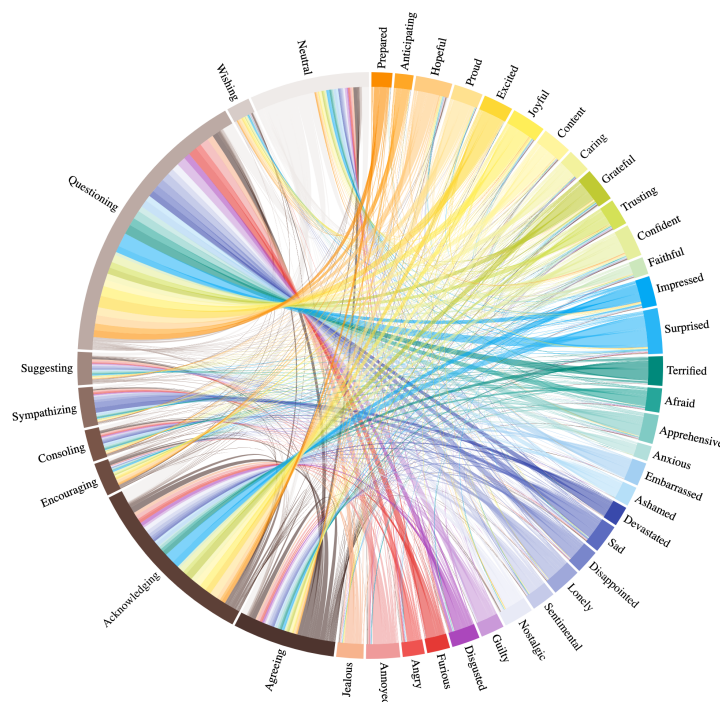


Figure 2: Visualization of emotion-intent exchanges between speakers and listeners in the Empathetic-Dialogues dataset irrespective of the dialogue turn. Each chord connects co-occurring emotion-intent pairs. The chord leaving a particular arc represents the speaker's emotion or intent and gets connected to the arc representing the listener's emotion or intent that immediately follows in a conversation.

It can be seen that 'questioning' and 'acknowledging' play a significant role in empathetic responses irrespective of the speaker's emotion—whether it is subtle or intense or has a positive or negative valence. Questioning enables the listener to sound more attentive and show interest in what the speaker describes. It prevents listeners from arriving at early conclusions, without knowing the situation in detail. It is also important to let speakers know that they have the right to feel the way they feel, even though listeners

may not completely agree with their choices. Expressions of 'acknowledgment' serve this purpose. This type of emotional interaction allows the speaker to elaborate on his feelings and what he is going through and feel validated at the same time. It can also be seen that some listener intents such as 'encouraging' and 'wishing' are frequently associated with positive speaker emotions, and some intents such as 'sympathizing' and 'consoling' are frequently associated with negative speaker emotions. A list of example utterance-response pairs corresponding to some of the most frequent emotion-intent exchanges ($\geq 100$ times out of $\approx 50k$ utterance-response pairs in EmpatheticDialogues) are given in Appendix B.

Next, we analysed how emotions and response intents shift over different turns in the dialogue as the dialogues progress in time. In this analysis, we discovered the most frequent emotion-intent flows that occur between speakers and listeners from the start to the end of conversations. To visualize the shift in emotions and intents over different dialogue turns, we computed the frequency of emotion-intent flow patterns up to 4 dialogue turns and plotted the ones having a frequency $\geq 5$. The reason for selecting only the first 4 turns in the dialogues is the fact that close to 77% of the dialogues in the EmpatheticDialogues dataset contain up to 4 turns and from this only 1.8% of the dialogues go up to 8 turns. Since there is comparatively much fewer data over dialogue turns from 5 to 8, we omitted these turns in our analysis.
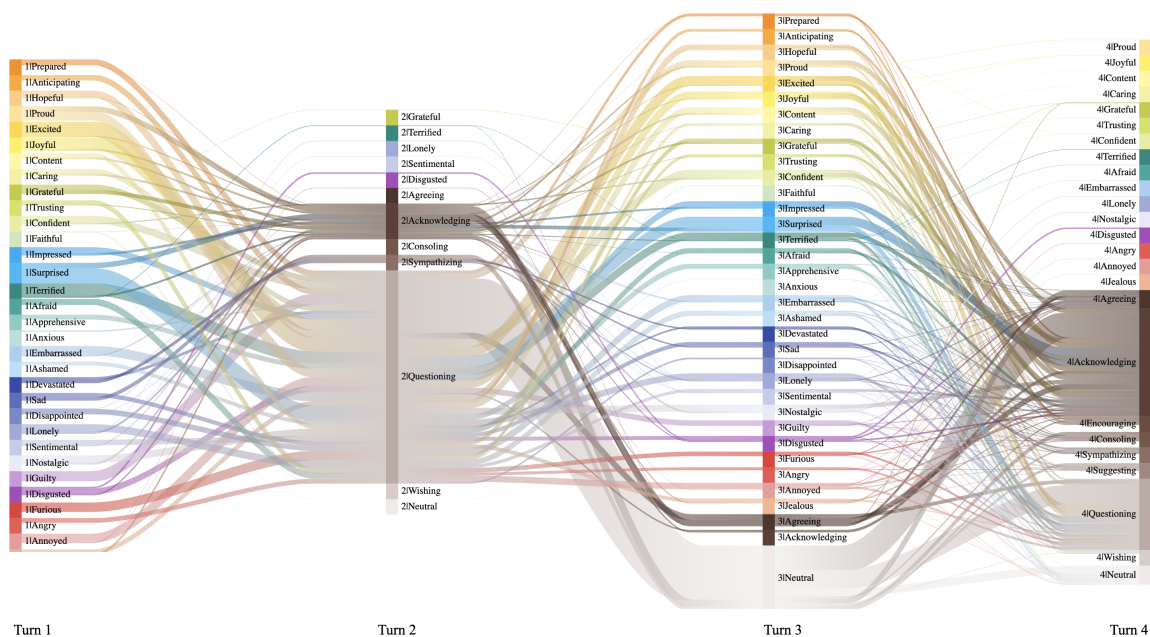


Figure 3: Visualization of the most common emotion-intent flow patterns (having a frequency $\geq 5$) throughout the first four dialogue turns in the EmpatheticDialogues dataset.

Figure 3 plots the most frequent emotion-intent flow patterns up to 4 turns in the dataset. Turns 1 and 3 correspond to speaker turns and turns 2 and 4 correspond to listener turns. According to the visualization, most emotions experienced by speakers are immediately followed by 'questions' as well as expressions of 'acknowledgment'. Expressions of 'sympathy' immediately follow more negatively intense emotions such as 'devastated' and 'sad'. Towards the end of dialogues, we can see more expressions of 'acknowledgment', 'agreement' and 'suggesting'. Expressions of 'encouragement' and 'wishing' can be seen in the case of positive emotional situations and 'sympathizing' and 'consoling' in the case of negative emotional situations. Another important observation is that towards the end of dialogues, listener utterances become more emotional compared to the beginning, as the speakers elaborate on their emotions. Such situations also reflect scenarios of personal distress—the phenomenon where one is unable to distinguish the emotion of their own from the emotion of another. Dialogue in Table 4 illustrates this phenomenon.

Still, they are not as frequent as how listeners choose to empathize healthily instead of making it a distress. But this sheds light on the fact that when the speaker goes on elaborating on his situation, sometimes the listener's ability to distinguish between the speaker's emotion and the emotion of his own may decrease, leading to consequences of personal distress. In the case of intense negative emotions,

| |
|---|
| S: *Bleh, I just had the worst food ever.* (Disgusted) |
| L: *What did you eat?* (Questioning) |
| S: *I was at Mcdonalds and was given a rotten cheese burger. I almost puked after I ate it.* (Disgusted) |
| L: *Oh gross, makes me never want McDonalds again.* (Disgusted) |

Table 4: Example conversation that illustrates personal distress towards the end of the dialogue.

it can lead to avoidance or a deliberate change of conversation topic. This is a scenario commonly experienced by people engaged in therapeutic and health professions and is described in researches by Singer and Klimecki (2014) and Buechel et al. (2018). However, in order to verify this observation more solidly, we need a corpus with a larger number of turns per dialogue.

The taxonomy we have developed can be incorporated into the design of social chatbots to gain more controllability and interpretability of the responses generated. It can be achieved by feeding in the conversation history, in which each utterance is tagged with an emotion or an intent label into a neural network that jointly models dialogue intent selection and response generation. Dialogue intent selection module will select the most appropriate intent based on the conversation history we feed in, and the response generation module will generate an appropriate response conditioned on the selected intent label. To help ensure more robustness, it is also possible to repeatedly sample plausible intent labels during training and feed them into the response generation module. The overall goal of modeling chatbots in this manner is to lead the conversation in a healthy and desirable direction with the controllability and interpretability provided by the taxonomy. Moreover, the taxonomy can be used as an annotation scheme to label utterances in other datasets and analyse them in terms of their empathetic quality in the same way described here. It also has implications in distinguishing between multiple forms of empathy—compassion and personal distress, as recognized in psychology and neuroscience fields.

One limitation of this study is the analysis results are highly dependent on the EmpatheticDialogues dataset. For future studies, we intend to curate a much larger empathetic dialogue dataset using a subset of the OpenSubtitles (8 million dialogues) (Lison et al., 2019), which will help us develop a more accurate emotion classifier and establish a more general taxonomy. Another limitation is the emotion classifier trained to automatically label utterances in the EmpatheticDialogues dataset is a sentence level classifier, which is unable to accurately distinguish similar utterances whose empathetic label can differ according to the context. We intend to improve our sentence level classifier into a classifier based on dialogue history that will be able to more accurately distinguish such cases. Automated labeling of intents using lexical methods also have the possibility to injure the robustness of the model due to considering only the most indicative n-grams in individual sentences without accounting for the surrounding context.

## 6 Conclusion

In this paper, we introduced a taxonomy of empathetic response intents capable of supporting automatic empathetic communication in social chitchat. The strategies relying on this taxonomy are essential for a chatbot to engage in prosocial conversations, expressing empathetic concern for its users, and keeping the users engaged. Another significant contribution from our work is to provide analysis on the EmpatheticDialogues corpus after automatically annotating it based on the most frequent intents from our taxonomy and 32 types of emotion categories defined in EmpatheticDialogues. We illustrate the most frequent emotion-intent exchange patterns in the dataset and how they vary temporally over the course of interaction. These results further validate the taxonomy of empathetic listener intents we derived and shed light on the frequent empathetic conversation patterns seen among humans when engaged in social chitchat. We explained how our taxonomy can be utilized in the development of an empathetic chatbot to achieve more controllability and interpretability in the responses generation process. The method described here can also be used as an annotation scheme to label utterances from other datasets and analyse them in terms of their empathetic quality. As future work, we plan on using these findings to develop a social chatbot capable of effectively engaging in empathetic conversations.

# References

Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. 2018. Improving language understanding by generative pre-training.

Alexander Miller, Will Feng, Dhruv Batra, Antoine Bordes, Adam Fisch, Jiasen Lu, Devi Parikh, and Jason Weston. 2017. A dialog research software platform. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 79—84.

Amy Skerry and Rebecca Saxe. 2015. Neural representations of emotion are organized around abstract event features. *Current Biology* 25:1945–1954.

Andreas Stolcke, Klaus Ries, Noah Coccaro, Elizabeth Shriberg, Rebecca Bates, Daniel Jurafsky, Paul Taylor, Rachel Martin, Carol Van Ess-Dykema, and Marie Meteer. 2000. Dialogue act modeling for automatic tagging and recognition of conversational speech. *Computational linguistics* 26(3):339–373.

Can Xu, Wei Wu, and Yu Wu. 2018. Towards explainable and controllable open domain dialogue generation with dialogue acts. *arXiv*, abs/1807.07255.

Carlos Busso, Murtaza Bulut, Chi-Chun Lee, Abe Kazemzadeh, Emily Mower, Samuel Kim, Jeannette N. Chang, Sungbok Lee, and Shrikanth S. Narayanan. 2008. IEMOCAP: Interactive emotional dyadic motion capture database. *Language resources and evaluation*, 42(4):335.

César Montenegro Portillo, Asier López Zorrilla, Javier Mikel Olaso Fernández, Roberto Santana Hermida, Raquel Justo Blanco, José Antonio Lozano Alonso, and María Inés Torres Barañano. 2019. A Dialogue-Act Taxonomy for a Virtual Coach Designed to Improve the Life of Elderly. *Multimodal Technologies Interact*, 3(52).

David M. Blei, Andrew Y. Ng, and Michael I. Jordan. 2017. Latent dirichlet allocation. *Journal of machine Learning research*, 3(Jan):993–1022.

Dolf Zillmann 2008. Empathy Theory. *W. Donsbach (Ed.), The International Encyclopedia of Communication.* pages 1530–1534, Malden, MA: Blackwell.

Edward Loper and Steven Bird. 2002. NLTK: The Natural Language Toolkit. In *Proceedings of the ACL Workshop on Effective Tools and Methodologies for Teaching Natural Language Processing and Computational Linguistics.*

Gary McKeown, Michel Valstar, Roddy Cowie, Maja Pantic, and Marc Schroder. 2011. The semaine database: Annotated multimodal records of emotionally colored conversations between a person and a limited agent. *IEEE transactions on affective computing*, 3(1):5–17.

Gery W. Ryan and H. Russell Bernard. 2003. Techniques to Identify Themes. *Field Methods*, 15(1):85–109.

Hannah Rashkin, Eric Michael Smith, Margaret Li and Y-Lan Boureau. 2019. Towards Empathetic Open-domain Conversation Models: A New Benchmark and Dataset. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 5370–5381, Florence, Italy.

Hao Zhou, Minlie Huang, Tianyang Zhang, Xiaoyan Zhu, and Bing Liu. 2018. Emotional chatting machine: Emotional conversation generation with internal and external memory. In *Thirty-Second AAAI Conference on Artificial Intelligence*.

Ilya Sutskever, Oriol Vinyals, and Quoc V. Le. 2014. Sequence to sequence learning with neural networks. In *Advances in neural information processing systems*, pages 3104–3112.

Iulian Vlad Serban, Alessandro Sordoni, Ryan Lowe, Laurent Charlin, Joelle Pineau, Aaron Courville, and Yoshua Bengio. 2017. A hierarchical latent variable encoder-decoder model for generating dialogues. In *Thirty-First AAAI Conference on Artificial Intelligence*.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 4171–4186, Minneapolis, Minnesota.

Jean Decety. 2010. The neurodevelopment of empathy in humans. *Developmental neuroscience*, 32(4):257–267.

Jianhua Yin, and Jianyong Wang. 2014. A dirichlet multinomial mixture model-based approach for short text clustering. In *Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 233–242.

Judith A. Holton. 2007. The coding process and its challenges. *The Sage handbook of grounded theory Part III*.

Lifeng Shang, Tetsuya Sakai, Zhengdong Lu, Hang Li, Ryuichiro Higashinaka, and Yusuke Miyao. 2016. Overview of the NTCIR-12 Short Text Conversation Task. In *Proceedings of NTCIR-12*, pages 473—484.

Mark Core and James Allen. 1997. Coding dialogs with the DAMSL annotation scheme. In *Working Notes of the AAAI Fall Symposium on Communicative Action in Humans and Machines*, pages 28–35, Cambridge, MA.

Nabiha Asghar, Pascal Poupart, Jesse Hoey, Xin Jiang, and Lili Mou. 2018. Affective neural response generation. In *European Conference on Information Retrieval* pages 154–166.

Nancy Eisenberg, and Natalie D. Eggum. 2009. Empathic responding: Sympathy and personal distress. *The social neuroscience of empathy* 6(2009):71–830.

Oriol Vinyals, and Quoc Le. 2015. A neural conversational model. In *Proceedings of the 31st International Conference on Machine Learning, Deep Learning Workshop*.

Paul Ekman. 1992. An argument for basic emotions. *Cognition & emotion* 6(3-4):169—200.

Pierre Lison, Jörg Tiedemann, and Milen Kouylekov. 2019. Open subtitles 2018: Statistical rescoring of sentence alignments in large, noisy parallel corpora. In *Eleventh International Conference on Language Resources and Evaluation (LREC)*, European Language Resources Association (ELRA).

Raymond Li, Samira Ebrahimi Kahou, Hannes Schulz, Vincent Michalski, Laurent Charlin, and Chris Pal. 2018. Towards deep conversational recommendations. In *Advances in neural information processing systems*, pages 9725–9735.

Robert M. Gordon 1995. The Simulation Theory: Objections and Misconceptions. *Mind & language* 7(1-2):11–34.

Robert M. Gordon 1995. Why the child's theory of mind really is a theory. *Folk Psychology: The Theory of Mind Debate* , edited by M. Davies and T. Stone, 232—258, Oxford: Blackwell.

Robert Plutchik. 1984. Emotions: A general psychoevolutionary theory. *Approaches to emotion* 197—219.

Ryuichiro Higashinaka, Kotaro Funakoshi, Masahiro Araki, Hiroshi Tsukahara, Yuka Kobayashi, and Masahiro Mizukami. 2015. Towards taxonomy of errors in chat-oriented dialogue systems. In *Proceedings of the 16th annual meeting of the special interest group on discourse and dialogue*, pages 87–95.

Saizheng Zhang, Emily Dinan, Jack Urbanek, Arthur Szlam, Douwe Kiela, and Jason Weston. 2018. Personalizing Dialogue Agents: I have a dog, do you have pets too?. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics*, pages 2204–2213, Melbourne, Australia.

Sheng-Yeh Chen, Chao-Chun Hsu, Chuan-Chun Kuo, and Lun-Wei Ku. 2018. Emotionlines: An emotion corpus of multi-party conversations. In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation*, Miyazaki, Japan.

Soujanya Poria, Devamanyu Hazarika, Navonil Majumder, Gautam Naik, Erik Cambria, and Rada Mihalcea. 2018. MELD: A multimodal multi-party dataset for emotion recognition in conversations. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 527–536, Florence, Italy.

Sven Buechel, Anneke Buffone, Barry Slaff, Lyle Ungar, and João Sedoc. 2018. Modeling empathy and distress in reaction to news stories. *arXiv*, abs/1808.10399.

Tania Singer, and Olga M. Klimecki. 2004. Empathy and compassion. *Current Biology* 24(18):R875–R878.

Tianran Hu, Anbang Xu, Zhe Liu, Quanzeng You, Yufan Guo, Vibha Sinha, Jiebo Luo, and Rama Akkiraju. 2018. Touch your heart: a tone-aware chatbot for customer care on social media. In *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems*, pages 1–12.

Timo Spring, Jacky Casas, Karl Daher, Elena Mugellini, and Omar Abou Khaled. 2019. Empathic Response Generation in Chatbots. *SwissText*.

Tsung-Hsien Wen, Milica Gasic, Nikola Mrksic, Pei-Hao Su, David Vandyke, and Steve Young. 2015. Semantically conditioned lstm-based natural language generation for spoken dialogue systems. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 1711–1721, Lisbon, Portugal.

Tomáš Mikolov, Martin Karafiát, Lukáš Burget, Jan Černocký, and Sanjeev Khudanpur. 2010. Recurrent neural network based language model. In *Eleventh annual conference of the international speech communication association*.

Vaswani, Ashish, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008.

Wanling Cai and Li Chen. 2019. Towards a Taxonomy of User Feedback Intents for Conversational Recommendations. In *Proceedings of ACM RecSys 2019 Late-breaking Results*, Copenhagen, Denmark.

Xianda Zhou and William Yang Wang. 2018. MojiTalk: Generating Emotional Responses at Scale. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics*, pages 1128–1137, Melbourne, Australia.

Yanran Li, Hui Su, Xiaoyu Shen, Wenjie Li, Ziqiang Cao, and Shuzi Niu. 2017. Dailydialog: A manually labelled multi-turn dialogue dataset. In *Proceedings of the Eighth International Joint Conference on Natural Language Processing*, pages 986–995, Taipei, Taiwan.

Yubo Xie, Ekaterina Svikhnushina, and Pearl Pu. 2019. A Multi-Turn Emotionally Engaging Dialog Model. *arXiv*, abs/1908.07816.

Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. *arXiv*, abs/1907.11692.

Armand Joulin, Edouard Grave, Piotr Bojanowski, and Tomas Mikolov. 2016. Bag of tricks for efficient text classification. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers* pages 427–431.

Bjarke Felbo, Alan Mislove, Anders Søgaard, Iyad Rahwan and Sune Lehmann. 2017. Using millions of emoji occurrences to learn any-domain representations for detecting sentiment, emotion and sarcasm. In *Proceedings of Empirical Methods in Natural Language Processing (EMNLP)*.

Zhaojiang Lin, Peng Xu, Genta Indra Winata, Zihan Liu, and Pascale Fung. 2019. Caire: An end-to-end empathetic chatbot. *arXiv*, abs/1907.12108.

Zhenqiao Song, Xiaoqing Zheng, Lu Liu, Mu Xu, and Xuan-Jing Huang. 2019. Generating Responses with a Specific Emotion in Dialog. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics* pages 3685–3695.

Zhongxia Chen, Ruihua Song, Xing Xie, Jian-Yun Nie, Xiting Wang, Fuzheng Zhang, and Enhong Chen. 2019. Neural Response Generation with Relevant Emotions for Short Text Conversation. In *CCF International Conference on Natural Language Processing and Chinese Computing*, pages 117–129.

**Appendix A. Words and phrases most indicative of the empathetic response intents that were used to extract more example listener utterances from the EmpatheticDialogues dataset for training the BERT transformer-based classifier**

| Response intent | Words and phrases most indicative of the intent |
|---|---|
| Agreeing | *100%, exactly, absolutely, definitely, agree, i know, me either, me neither, i understand, i completely understand, me too, that's right, you're right, correct* |
| Acknowledging | *it sucks, that sucks, i'd ... too, i would ... too, i feel you, that's splendid, i bet ... was, that's great", that's a good idea, i bet ... can't, that's pretty, i see, it's pretty, can understand, sounds, that would, i would have, must've, cool, nice, awesome* |
| Encouraging | *hopefully ... will, i hope ... will, works out for you, i bet ... will, i bet ... 'll, i bet ... can* |
| Consoling | *there you go, hopefully ... will, i hope ... will, cheer up, get better, will pass quickly* |
| Sympathizing | *i'm sorry, sorry to hear, oh no, bless you, deepest sympathy* |
| Suggesting | *maybe, i think ... should, perhaps, why don't you, you could always, what if* |

| Questioning | *what ... ?, why ... ?, when ... ?, where ... ?, how ... ?, are ... ?, is ... ?, did ... ?, do ... ?, does ... ?, have ... ?, has ... ?, had ... ?* |
|---|---|
| Wishing | *congratulations, happy birthday, happy anniversary, i wish you, wish you ... !, all the best, good luck* |

Table 5: Words and phrases that are most indicative of the empathetic response intents.


## Appendix B. Example speaker-listener utterance pairs corresponding to the taxonomy of emotion/intent exchanges

| Speaker's emotion | Listener's response emotion/intent | Example utterance-response pairs |
|---|---|---|
| Anticipa-ting | Questioning | S: *When tax season came I was in a hurry to get mine done. I was looking forward to a big refund.* (Anticipating)<br>L: *really? why is that?* (Questioning) |
| | Acknowledging | S: *I cannot wait for the newest Pokemon game, it looks amazing to me!* (Anticipating)<br>L: *Those games do seem fun* (Acknowledging) |
| Joyful | Questioning | S: *i was happy to see that i was able to get a new pet the other day* (Joyful)<br>L: *What pet did you get?* (Questioning) |
| | Acknowledging | S: *I jumped for joy when my baby was born.* (Joyful)<br>L: *wow that must have been a huge moment for you* (Acknowledging) |
| Trusting | Questioning | S: *Man, I let one of my friends take my Benz one day to run some errands. I really thought she would be careful with it.* (Trusting)<br>L: *Oh, no! Did she damage your car?* (Questioning) |
| | Acknowledging | S: *My therapist was so kind to me, I had to tell her a lot.* (Trusting)<br>L: *That's good you have someone that you can talk to about your problems and feelings. I'm sure it helps!* (Acknowledging) |
| Surprised | Questioning | S: *I was shocked when i got invited on a random trip* (Surprised)<br>L: *Was a happy shocked feeling or a bad one?* (Questioning) |
| | Acknowledging | S: *The other day I found out that my sister is having twins!* (Surprised)<br>L: *Oh that's wonderful twins seem really cool.* (Acknowledging) |
| | Neutral | S: *No one even knew she was dating anyone until the announcement, so I was very surprised.* (Surprised)<br>L: *I guess she wanted to keep it a secret for some reason.* (Neutral) |
| Afraid | Questioning | S: *It's so dark and creepy down there.* (Afraid)<br>L: *lol. Do you think there are monsters down there?* (Questioning) |
| | Acknowledging | S: *It was only off for a little over 2 hours, but I could not find a flashlight and it was so scary.* (Afraid)<br>L: *That sounds awful!* (Acknowledging) |

| Sad | Questioning | S: *I feel bad I don't always get to go through bad things and full get healed.* (Sad)<br>L: *Do you mean you feel bad that you don't get to go through bad things or that you don't get to be healed?* (Questioning) |
|------|-------------|---------------------------------------------------------------------------------------|
|  | Sympathizing | S: *I was extremely emotional when my dog passed away* (Sad)<br>L: *Aww man sorry for your loss, those are the worst.* (Sympathizing) |
|  | Acknowledging | S: *My favorite donut shop went out of business.* (Sad)<br>L: *Ah that's a pity. It really sucks to lose favorite shops.* (Acknowledging) |
|  | Agreeing | S: *I'm sad. My youngest son starts kindergarten tomorrow!* (Sad)<br>L: *I am sure it is a bittersweet moment. I can relate myself.* (Agreeing) |
| Disgusted | Questioning | S: *I am disgusted that so many people voted in favour of Brexit in the UK.* (Disgusted)<br>L: *Why is that?* (Questioning) |
|  | Acknowledging | S: *It was a brand new box of Rice crispies. When I opened it and poured it in my bowl, there were several live bugs.* (Disgusted)<br>L: *Well that sounds disgusting* (Acknowledging) |
|  | Disgusted | S: *I was at Mcdonalds and was given a rotten cheese burger. I almost puked after I ate it.* (Disgusted)<br>L: *Oh gross, makes me never want McDonalds again.* (Disgusted) |
|  | Agreeing | S: *Everytime I see my cat vomit on floor it makes me sick.* (Disgusted)<br>L: *i think you have the same attitude like me.* (Agreeing) |
| Angry | Questioning | S: *i was upset when i saw someone put a dent in my door* (Angry)<br>L: *Was this a parking lot?* (Questioning) |
|  | Acknowledging | S: *My grandma didn't make my oatmeal right yesterday. I was so mad.* (Angry)<br>L: *Oh wow! You were pretty angry* (Acknowledging) |

Table 6: Example speaker and listener utterances corresponding to the most common emotion exchanges between speakers and listeners, when the speaker's emotion is one of the Plutchik's 8 basic emotions.