

# Use of a Taxonomy of Empathetic Response Intents to Control and Interpret Empathy in Neural Chatbots

Anonymous ACL submission

## Abstract

A recent trend in the domain of open-domain conversational agents is enabling them to converse empathetically to emotional prompts. Current approaches either follow an end-to-end approach or condition the responses on similar emotion labels to generate empathetic responses. But empathy is a broad concept that refers to the cognitive and emotional reactions of an individual to the observed experiences of another and it is more complex than mere mimicry of emotion. Hence, it requires identifying complex human conversational strategies and dynamics in addition to generic emotions to control and interpret empathetic responding capabilities of chatbots. In this work, we make use of a taxonomy of eight empathetic response intents in addition to generic emotion categories in building a dialogue response generation model capable of generating empathetic responses in a controllable and interpretable manner. It consists of two modules: 1) a response emotion/intent prediction module; and 2) a response generation module. We propose several rule-based and neural approaches to predict the next response’s emotion/intent and generate responses conditioned on these predicted emotions/intents. Automatic and human evaluation results emphasize the importance of the use of the taxonomy of empathetic response intents in producing more diverse and empathetically more appropriate responses than end-to-end models.

## 1 Introduction

End-to-end neural dialogue response generation has revolutionized the design of open-domain conversational agents or chatbots due to requiring little or no manual intervention and its ability to largely generalize (Sordoni et al., 2015; Shang et al., 2015; Vinyals and Le, 2015). It overcomes many limitations of traditional rule-based response generation techniques such as the cost of domain expertise and predictability of responses. But due to the

Dialogue context:	
Speaker:	<i>I think that the girl of my dreams likes somebody else. I feel very sad about it.</i>
Listener:	<i>Ooh, am so sorry about that. <b>Have you tried to talk to her?</b></i>
Speaker:	<i>It's tough as she has been out of the country for a month, so I will likely discuss it when she returns.</i>
Possible responses:	
(No control)	<i><b>Have you talked to her about it yet?</b></i> (Repetitive)
(No control)	<i><b>I don't think that's a good idea.</b></i> (Not encouraging to the speaker)
(Conditioned on: Encouraging)	<i><b>I hope everything works out for you.</b></i> (Empathetically appropriate)

Table 1: An example dialogue showing how controllability affects response generation.

black-box nature of these end-to-end models, they offer very little controllability to the developer and generate responses that are difficult to interpret (Wu et al., 2018, 2020; Gupta et al., 2020), making these approaches less reliable and fail-safe (d’Avila Garcez and Lamb, 2020). A recent example is Microsoft’s Taybot that started producing unintended, and offensive tweets denying the Holocaust as a result of learning from racist and offensive information on Twitter (Lee, 2016). Having control over the generated responses would have enabled the chatbot to avoid malicious intentions and carefully choose how to converse. Thus, it is important to look at ways how developers can gain control over the responses generated by end-to-end neural response generation models and how they can be made interpretable.

Recent research has taken efforts to induce controllability and interpretability into end-to-end models. For example, Xu et al. (2018) explore how the flow of human-machine interactions can be managed by introducing dialogue acts as policies to the dialogue generation model. Sankar and Ravi (2019) show that conditioning the response generation process on interpretable dialogue attributes such as dialogue acts and sentiment helps to eliminate repetitive responses and makes the model more interesting and engaging.

In contrast to task-oriented dialogue systems designed to help people complete specific tasks, open-domain chatbots are designed to engage users in human-machine conversation for entertainment and emotional companionship (Wu and Yan, 2018). Hence, in open-domain conversations, controllability should also be studied with respect to aspects such as humor, personality, emotions, and empathy, which cannot be achieved using generic dialogue acts. In this study, our focus is on controlling empathy in open-domain chatbot responses, which requires understanding conversational strategies used in human-human empathetic conversations.

Earlier studies gain control in this aspect by conditioning the response on either manually specified (Zhou et al., 2018; Zhou and Wang, 2018; Hu et al., 2018; Song et al., 2019) or automatically predicted (Chen et al., 2019) sentiment or emotion labels. However, an analysis by Welivita and Pu (2020) on human-human conversations of the EmpatheticDialogues dataset (Rashkin et al., 2018) reveals, listeners are much more likely to respond to positive or negative emotions with specific empathetic intents such as *acknowledgment*, *consolation* and *encouragement*, rather than expressing similar or opposite emotions. They introduce a taxonomy of eight response intents that can better describe empathetic human responses to emotional dialogue prompts. In this paper, we explore how end-to-end response generation can be combined with more advanced control of empathy by utilizing the above taxonomy of empathetic response intents in addition to existing emotion categories. To provide a glimpse of what we aim to achieve, in Table 1 we show how conditioning the response on an empathetic response intent chosen based on the dialogue history can serve in producing a more empathetically appropriate response. It avoids repetitive or sub-optimal responses generated by end-to-end approaches without any control.

Our empathetic response generation model consists of two modules: 1) a response emotion or intent prediction module; and 2) a response generation module. We experiment with both rule-based and neural approaches for predicting the next response’s emotion or intent. For the rule-based approaches for predicting the response emotion/intent, we develop two decision tree-based response emotion and intent prediction methods. For the neural approach for predicting the response emotion/intent, we develop a classifier based on

the BERT transformer architecture (Vaswani et al., 2017; Devlin et al., 2019). The reason why we evaluate the performance of rule-based approaches is that they are much simpler than neural models and save a lot of training time and resources. Thus, if considerable performance can still be achieved through rule-based approaches compared to the baselines, it is worth considering the use of such simpler approaches over sophisticated neural approaches, especially in resource-limited environments. The emotions and intents predicted by these methods are then used to condition the responses generated by the response generation module. For training and evaluating these models, we use two state-of-the-art dialogue datasets containing empathetic conversations: 1) the EmpatheticDialogues dataset (Rashkin et al., 2018); and 2) the EDOS (Emotional Dialogues in OpenSubtitles) dataset (Welivita et al., 2021). The automatic and human evaluation results confirm the importance of the use of the taxonomy in generating more diverse and empathetically more appropriate responses than end-to-end models.

Our contributions in this paper are three folds. 1) We explore the ability of a taxonomy of empathetic response intents in controlling and interpreting the responses generated by open-domain conversational agents for emotional prompts. 2) We propose an empathetic response generation model consisting of a response emotion/intent prediction module and a response generation module to generate empathetic responses in a controllable and interpretable manner. 3) We experiment with both rule-based and neural approaches in predicting the next response’s emotion or intent and evaluate their performance in conditional generation of empathetic responses using automatic and human evaluation metrics against standard baselines.

## 2 Literature Review

Existing conversational agents are designed for either open-domain or specific task completion (Gao et al., 2018). Regarding the former, a common practice is to generate dialogue in an end-to-end fashion (Sordoni et al., 2015; Shang et al., 2015; Vinyals and Le, 2015). Often responses generated by these methods are unpredictable and not fail-safe (d’Avila Garcez and Lamb, 2020). Hence, recent research has focused on methods to control and interpret the responses generated by open-domain neural conversational agents. Mainly we

find three methods they use to control the generated response: 1) by a manually specified value (Zhou et al., 2018; Zhou and Wang, 2018; Hu et al., 2018; Song et al., 2019); 3) by rules that are predefined or derived from the training data (Hedayatnia et al., 2020); 3) by an automatically predicted value from a neural network model (Xie and Pu, 2021; Wu et al., 2018; Sankar and Ravi, 2019; Santhanam et al., 2020; Ke et al., 2018; Lee et al., 2020).

Specially in studies addressing emotional response generation, a manually specified sentiment, emotion (Zhou et al., 2018) or an emoji (Zhou and Wang, 2018) was used to control the sentiment or emotionality of the responses generated. Later, more and more research focused on automatically predicting values or deriving rules such that they could be used to control the generated response without manual intervention. For example, Sankar and Ravi (2019) used an RNN based policy network to predict the next dialogue act given previous dialogue turns and dialogue attributes. Hedayatnia et al. (2020) used rules designed as a set of dialogue act transitions from common examples in the Topical-Chat corpus (Gopalakrishnan et al., 2019) to plan the content and style of target responses.

But all the above work focused on achieving controllability using generic dialogue acts or generating controlled emotional responses conditioned on similar or opposite emotions, emojis, or sentiment tags. These labels do not suffice the controlled generation of meaningful empathetic responses because humans demonstrate a wide range of emotions and intents when regulating empathy (Welivita and Pu, 2020). Previous work also lacks comparisons between rule-based and automatic conditioning methods used to control response generation. In this work, we address the above gaps by investigating how empathy in neural responses can be controlled using a taxonomy of eight empathetic response intents (Welivita and Pu, 2020), in addition to 32 emotion categories, while evaluating the applicability of both rule-based and automatic control mechanisms for this task.

### 3 Methodology

Our controllable and interpretable empathetic response generation architecture consists of two modules: 1) the response emotion/intent prediction module; and 2) the response generation module. The emotion or intent predicted by the first module is input into the second to condition the response

generated by the second module. In the following sections we discuss the datasets used for our experiments, the different rule-based and automatic emotion/intent prediction methods we propose, how the emotions and intents predicted by these modules are used to generate responses that are both controllable and interpretable, and the different evaluation methods we utilize to compare the performance of these approaches on two state-of-the-art dialogue datasets containing emotional dialogue prompts.

#### 3.1 Datasets

We utilized the EmpatheticDialogues dataset proposed by Rashkin et al. (2018), and the OS (Open-Subtitles) and EDOS (Emotional Dialogues in OpenSubtitles) dialogue datasets proposed by Welivita et al. (2021) to train and evaluate our models. The EmpatheticDialogues (ED) dataset contains  $\approx 25$ K open-domain human-human conversations carried out between a speaker and a listener. Each conversation is conditioned on one of 32 emotions selected from multiple annotation schemes. The OS and EDOS datasets are curated by applying a series of preprocessing and turn segmentation steps on the movie and TV subtitles in the Open-Subtitles 2018 corpus (Lison et al., 2019). The EDOS dataset contains 1M highly emotional dialogues filtered from the rest of the OS dialogues. Even though the speaker and listener turns in the OS and EDOS datasets are not clearly defined, we assumed the odd-numbered turns (1, 3, 5, ...) as speaker turns and even-numbered turns (2, 4, 6, ...) as listener turns for our experiments. We used the OS dialogues dataset containing  $\approx 3$ M dialogues for pre-training and the ED and EDOS datasets to separately fine-tune the models. The statistics of these datasets are denoted in Table 2. From each dataset, 80% of the data was used for training, 10% for validation, and the remaining 10% for testing.

Dataset	Dialogues	Turns	Turns/dialogue
OS	2,989,774	11,511,060	3.85
ED	24,847	107,217	4.32
EDOS	1,000,000	2,940,629	2.94

Table 2: Statistics of the datasets used for training and evaluating the models.

We used a BERT (Devlin et al., 2019) transformer-based dialogue emotion classifier proposed by Welivita et al. (2020) to automatically annotate all dialogue turns in the above datasets. This classifier is trained on 25K situation descrip-

Empathetic intent	Example response
1. Questioning	<i>What's the matter? What's wrong?</i>
2. Agreeing	<i>Exactly, I get that entirely!</i>
3. Acknowledging	<i>Sounds awesome!</i>
4. Encouraging	<i>Just give it a trial.</i>
5. Consoling	<i>I hope everything works out for you.</i>
6. Sympathizing	<i>I am sorry to hear that.</i>
7. Wishing	<i>Congrats, that's a step forward.</i>
8. Suggesting	<i>Maybe you should talk to her.</i>

Table 3: The taxonomy of listener specific empathetic response intents used to achieve controllability and interpretability in the responses generated.

tions from EmpatheticDialogues labeled with 32 emotion classes, 7K EmpatheticDialogues listener turns labeled with eight empathetic response intents and *Neutral*, and 14K emotion and intent annotated dialogue turns from the OSED dataset. It has a final annotation accuracy of 65.88% over 41 labels, which is significant compared to the other state-of-the-art dialogue emotion classifiers (Welivita et al., 2020). We use the emotion and intent labels suggested by the above classifier as ground-truth labels for our experiments.

## 3.2 Response Emotion/Intent Prediction

To generate controlled and interpretable empathetic responses, we utilized 32 fine-grained emotions and a taxonomy of listener-specific empathetic response intents. The 32 emotions are emotion categories on which dialogues in the EmpatheticDialogues dataset are conditioned on (Rashkin et al., 2018). They range 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). We further utilized the taxonomy of empathetic response intents proposed by Welivita and Pu (2020), which is derived by analysing the listener responses in the EmpatheticDialogues dataset. These intents are denoted in Table 3 along with corresponding examples. To predict the emotion or intent of the next response, we propose several rule-based and neural response emotion/intent prediction methods, which are described in the following subsections.

### 3.2.1 Baselines

As a baseline, we sample a response emotion or intent from the set of eight empathetic response intents plus the most recent emotion encountered in the last  $k$  ( $k = 3$ ) dialogue turns. This is based on the observations by Welivita and Pu (2020) on the EmpatheticDialogues dataset (Rashkin et al., 2018). They state that in human empathetic conversations,

the listener’s response to emotional prompts mostly contain an empathetic response intent identified by their taxonomy or a statement with similar emotion. This baseline is inspired by the work of Hedayatnia et al. (2020), in which the response dialogue act is chosen among the most frequently seen dialogue acts based on an equal probability distribution.

The other baseline we used when generating responses is the plain end-to-end transformer model proposed by Vinyals et al. (2017), in which no conditioning is used when generating the response.

### 3.2.2 Rule-based Decision Tree Approaches

We propose two non-neural, decision tree-based response intent prediction methods that leverage the knowledge of the emotion-intent flow of the dialogues in the training dataset. The basic idea of a decision tree for this context is denoted along with an example in Figure 1. The probabilities of emotions and intents in the branches in the decision tree are learned from the training data itself by traversing through dialogues using a window of size  $k$ , where  $k$  is the maximum depth of the decision tree. The window is moved forward two dialogue turns at a time capturing the probability of speaker-listener emotion-intent exchanges in the training dataset.

Here, we used a window of size 4 mainly because most dialogues contained in the ED, OS, and EDOS datasets were limited to four dialogue turns. During inference, an emotion or an intent is sampled based on the sequence of emotions and intents in the previous  $(k - 1)$  dialogue turns. We used two different methods: 1) argmax; and 2) probabilistic sampling, to sample the response emotion or intent from the decision tree. In the argmax method, we chose the emotion or intent with the highest probability in the decision tree based on the sequence of emotions and intents in the previous  $(k - 1)$  dialogue turns. In the probabilistic sampling method, we sampled an emotion or an intent based on the distribution of probabilities in the decision tree given the sequence of emotions and intents in the previous  $(k - 1)$  dialogue turns. We refer to these two decision tree-based methods as *DT (argmax)*, and *DT (prob. sampled)*.

We have more control over the above methods than neural response intent prediction methods since we can foresee where the dialogue will be directed by visualizing the decision trees beforehand. For example, the decision trees generated using the EmpatheticDialogues and EDOS training

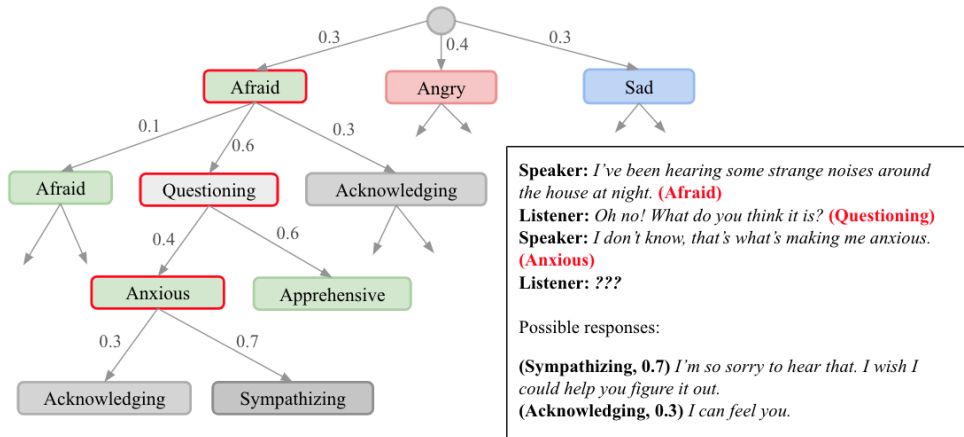


Figure 1: Visualization of a simpler version of our decision tree approach to predict the response emotion or intent.

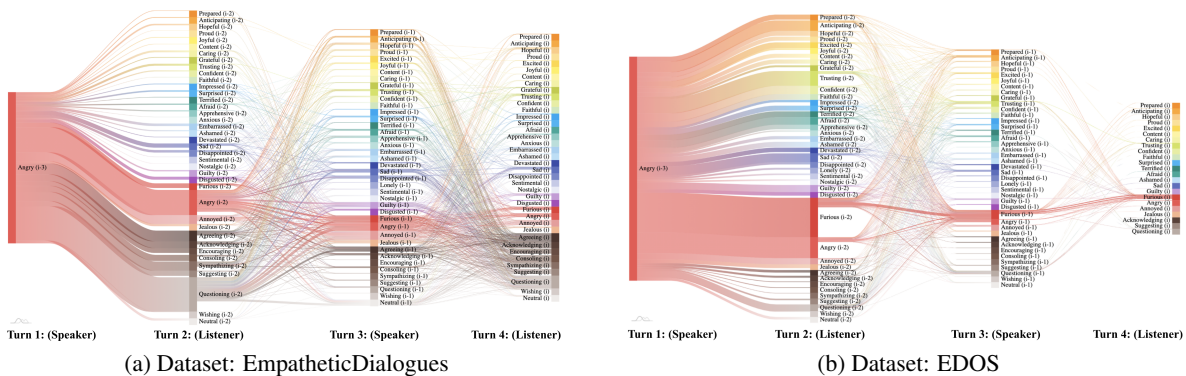


Figure 2: Decision trees generated using the EmpatheticDialogues and EDOS training datasets when the emotion of the beginning dialogue prompt is *Angry*.

355 datasets when the emotion of the beginning dia- 376  
 356 logue prompts is *Angry* are denoted in Figure 2. 377  
 357 As it could be observed, in the ED dataset, the lis- 378  
 358 teners mostly respond to speakers' emotions with 379  
 359 one of the intents from the taxonomy of empathetic 380  
 360 response intents. The EDOS dataset by nature is 381  
 361 more dramatic, in which both the speaker and the 382  
 362 listener become emotional. This phenomenon is 383  
 363 called "emotional contagion" in the psychological 384  
 364 literature (Hatfield et al., 1993). For example in 385  
 365 EDOS, if the speaker is angry, the listener also 386  
 366 tends to reply back with anger. These communica- 387  
 367 tion patterns could clearly be visualized with the 388  
 368 decision trees created and the developer can predict 389  
 369 beforehand how the chatbots whose responses are 390  
 370 conditioned on these emotion-intent patterns would 391  
 371 behave for a given emotional prompt. 392

### 3.2.3 Neural Response Emotion and Intent Predictor

372 An automatic method for predicting the next 393  
 373 response's emotion or intent is using a neural 394  
 374 395  
 375 396  
 397

network-based response emotion/intent predictor. 376  
 An advantage of using neural approaches to deter- 377  
 mine the emotion or intent of the next response 378  
 is that they can leverage clues from the semantic 379  
 content of the previous dialogue turns in addition 380  
 to the flow of emotions and intents when predict- 381  
 ing the response emotion or intent. Our neural re- 382  
 sponse emotion/intent predictor consists of a BERT 383  
 transformer-based encoder architecture (representa- 384  
 tion network) followed by an attention layer for 385  
 aggregating individual token representations, a hid- 386  
 den layer, and a softmax as depicted in Figure 3. 387  
 The BERT-base architecture with 12 layers, 768 388  
 dimensions, 12 heads, and 110M parameters is 389  
 used as the representation network. It is initialized 390  
 with weights from the pre-trained language model 391  
 RoBERTa (Liu et al., 2019). 392

We concatenate the previous  $k$  dialogue turns 393  
 as depicted in Figure 3 and they are input to the 394  
 encoder of the model. The emotions and intents 395  
 corresponding to these  $k$  dialogue turns are added 396  
 to the word embeddings and positional embeddings 397

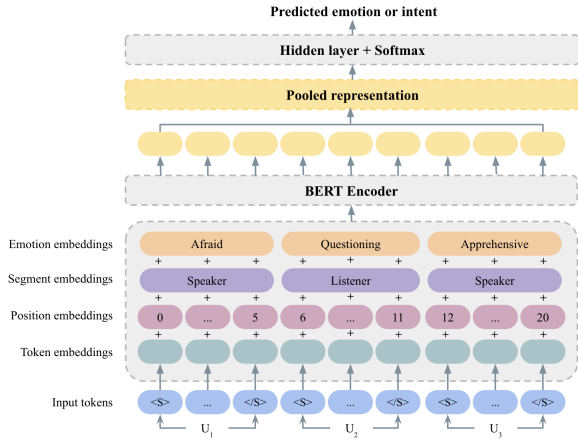


Figure 3: Architecture of the neural response emotion/intent predictor.

in the original transformer architecture. This additional knowledge helps the model to get a better understanding of the flow of emotions and intents in the previous dialogue turns. The emotions and intents are embedded into a vector space having the same dimensionality as the word and position embeddings so they can add up. In addition, we also incorporate segment embeddings that differentiate between speaker and listener turns. We pre-trained the model on the OS dialogues dataset and fine-tuned it separately on ED and EDOS datasets. The hyper-parameters used during training and other training details are described in the appendices.

### 3.3 Response Generation

For response generation, we used a plain transformer-based encoder-decoder architecture (end-to-end model) as a baseline (Vaswani et al., 2017). To generate controlled empathetic responses, we incorporated the different response emotion/intent prediction methods described above as input to the decoder. Figure 4 shows the overall architecture of our models.

The input representation for the encoder of the generation model is the same as the input representation used for the neural response emotion/intent predictor described in section 3.2.3. The vector representation generated by the encoder is input into the decoder along with the embedding of the emotion or intent predicted by the response emotion/intent predictor. During training, instead of the predicted emotion or intent, we used the ground-truth emotion or intent. The generation model is first pre-trained on OS dialogues and then fine-tuned on ED and EDOS datasets separately.

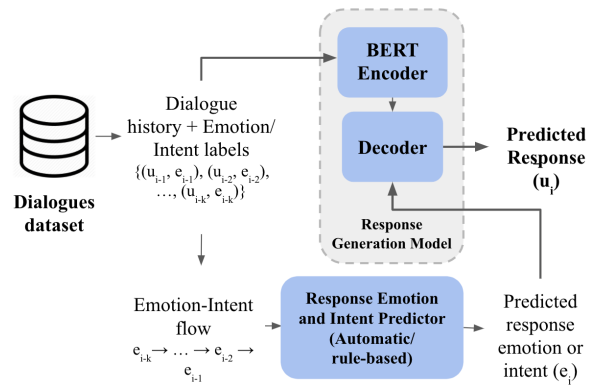


Figure 4: Overall architecture of the controllable and interpretable empathetic response generation model.

## 4 Evaluation and Results

### 4.1 Automatic Evaluation Results

Evaluation by means of automatic metrics was carried out separately for response emotion/intent prediction and conditional response generation. The following subsections describe the results obtained in these evaluations.

#### 4.1.1 Prediction Performance

The weighted precision, recall, F1, and balanced accuracy scores computed for different response emotion/intent prediction methods across ED and EDOS testing datasets are indicated in Table 4.

According to the weighted precision, recall, F1, and accuracy scores, the neural emotion/intent predictor performed the best compared to other prediction methods. Among rule-based approaches for response emotion/intent prediction, the DT (argmax) method performed the best. The DT (argmax) method had considerable improvement in recall, F1, and accuracy scores over the equally sampled baseline.

#### 4.1.2 Generation Performance

To evaluate the performance of response generation, we computed the perplexity, diversity metrics (distinct unigram and distinct bigram scores), and vector extrema cosine similarity on ED and EDOS testing datasets. They are denoted in Table 5. We also evaluated the responses generated by a model conditioned on the ground-truth emotion or intent of the next response to see how well the taxonomy of empathetic response intents alone contributes to better empathetic response generation performance.

According to the results, the models whose response was conditioned on the ground-truth re-

Model	Trained on: OS + ED Tested on: ED				Trained on: OS + EDOS Tested on: EDOS			
	Prec.	Recall	F1	Acc.	Prec.	Recall	F1	Acc.
Equally sampled	0.1138	0.0667	0.0638	0.0410	0.0981	0.0221	0.0232	0.0285
DT (argmax)	0.0959	0.0883	0.0883	0.0692	0.0755	0.1016	0.0799	0.0419
DT (prob. sampled)	0.0715	0.0663	0.0680	0.0480	0.0627	0.0616	0.0619	0.0345
Neural predictor	0.1634	0.1636	0.1472	0.1163	0.1306	0.1712	0.1181	0.0679

Table 4: Weighted precision, recall, F1 and accuracy scores computed for ED and EDOS test datasets. The cells in dark green indicate the best scores and the cells in light green indicate the second best scores.

Model	Trained on: OS + ED Tested on: ED				Trained on: OS + EDOS Tested on: EDOS			
	PPL	D-1	D-2	Embed. extrema	PPL	D-1	D-2	Embed. extrema
GT emotion/intent	11.74	0.0823	0.2812	0.5181	12.57	0.0846	0.2552	0.4539
End-to-end model	12.26	0.0544	0.1612	0.5015	13.13	0.0784	0.228	0.4365
Equally sampled	13.48	0.0761	0.2469	0.4824	14.20	0.0754	0.2229	0.433
DT (argmax)	13.23	0.0865	0.2977	0.4892	14.14	0.0727	0.2419	0.4458
DT (prob. sampled)	13.37	0.0795	0.2761	0.4828	14.23	0.0763	0.2418	0.436
Neural predictor	13.15	0.0835	0.2811	0.4851	13.97	0.0805	0.2415	0.4403

Table 5: Perplexity (PPL), diversity metrics (distinct unigrams: D-1; and distinct bigrams: D-2), and vector extrema cosine similarity (Embed. extrema) calculated on ED and EDOS testing datasets.

467 sponse emotion or intent performed the best in  
468 terms of perplexity and embedding extrema in both  
469 ED and EDOS datasets and in terms of diversity  
470 metrics in the EDOS dataset. These results em-  
471 phasize the usefulness of the taxonomy of em-  
472 pathetic response intents and the 32 fine-grained  
473 emotion categories in generating controlled em-  
474 pathetic responses. The models incorporating the  
475 DT (argmax) approach scored the best in terms of  
476 diversity metrics in the ED test dataset.

## 477 4.2 Human Evaluation

478 In addition to the automatic metrics, we carefully  
479 designed a human evaluation experiment in Ama-  
480 zon Mechanical Turk (AMT) to evaluate responses’  
481 empathetic appropriateness. We selected a total of  
482 1,000 dialogue cases: 500 ED and EDOS dialogues  
483 for testing. The AMT workers had to drag and  
484 drop responses generated by five models (end-to-  
485 end; models whose response was conditioned on  
486 the equally sampled baseline, DT argmax, DT prob.  
487 sampled and the neural predictor) into areas *Good*,  
488 *Okay*, and *Bad*, depending on their empathetic ap-  
489 propriateness. We neglected responses conditioned  
490 on the ground-truth emotion or intent since we are  
491 more interested in automatically predicted labels.  
492 We bundled 10 dialogues into a HIT (Human Intel-  
493 ligence Task) so that one worker works on at least  
494 10 cases to avoid too much bias between answers.  
495 To evaluate the quality of the work generated, we  
496 included three quiz questions equally spaced in  
497 a HIT. In these, we included the ground-truth re-  
498 sponse among the other responses generated by

499 the models. If a worker rated the ground-truth re-  
500 sponse either as *Good* or *Okay*, then a bonus point  
501 was added. To encourage attentiveness to the task,  
502 for those who obtained at least two out of three quiz  
503 questions correct, we gave a bonus of 0.1\$. Three  
504 workers were allowed to work on a HIT and only  
505 the ratings that were agreed by at least two work-  
506 ers, both who have obtained bonuses, were taken  
507 to compute the final scores. As a result, 8.33%  
508 of the answers were disqualified. The results of  
509 the experiment are denoted in Table 6. The ex-  
510 periment yielded an inter-rater agreement (Fleiss’  
511 kappa) score of 0.2294 indicating fair agreement.

512 According to the results, the neural predictor  
513 scored the highest percentage of *Good* ratings in  
514 both ED and EDOS testing datasets. The models  
515 that use the equally sampled approach performed  
516 the worst producing the highest percentage of re-  
517 sponses ranked *Bad*. An interesting observation  
518 is that the DT (argmax) method scored the most  
519 number of combined *Good* and *Okay* responses  
520 in ED and EDOS testing datasets confirming that  
521 rule-based approaches such as the decision tree ap-  
522 proach we propose could be used to control and  
523 interpret the responses without losing significant  
524 accuracy.

## 525 4.3 Case Study

526 In Table 7 we show some example responses gener-  
527 ated by different models for a given dialogue con-  
528 text. It could be noticed that having no response  
529 control mechanism and having a response condi-  
530 tioned on an inappropriate intent both can result

Model	Trained on: OS + ED Tested on: ED				Trained on: OS + EDOS Tested on: EDOS			
	Good	Okay	(Good + Okay)	Bad	Good	Okay	(Good + Okay)	Bad
End-to-end model	46.94	15.72	62.66	37.34	27.48	27.02	54.50	45.50
Equally sampled	25.22	25.22	50.44	49.57	19.35	29.84	49.18	50.82
DT (argmax)	37.61	33.70	71.31	28.69	21.70	36.72	58.42	41.58
DT (prob. sampled)	24.17	21.68	45.85	42.13	21.96	31.78	53.74	46.26
Neural predictor	50.00	17.39	67.39	32.61	24.82	33.01	57.83	42.17

Table 6: Human evaluation results (as a percentage) corresponding to ED and EDOS testing datasets.

Dialogue context:			
Speaker: <i>Last week when Toys R Us closed it really made me start thinking of the 90s and my childhood (Nostalgic)</i>			
Listener: <i>I've heard a lot of people say something similar. Did you go there a lot? (Questioning)</i>			
Speaker: <i>I did! I was talking to my husband about the different toys I had. Even started talking about the music and shows from that era. (Nostalgic)</i>			
Model	Predicted emotion/intent	Generated response	Human eval.
Ground-truth emotion/intent	Jealous	<i>That sounds like a lot of fun. I wish I had a husband like that.</i>	Good
End-to-end model	-	<i>I'm sorry to hear that. What kind of toys were you talking about?</i>	Bad
Equally sampled	Wishing	<i>I wish you the best of luck!</i>	Bad
DT (argmax)	Questioning	<i>What kind of toys were you talking about?</i>	Good
DT (prob. sampled)	Encouraging	<i>I bet that was fun</i>	Good
Neural predictor	Encouraging	<i>I bet that was fun</i>	Good

Table 7: An example dialogue showing that both lack of controllability and conditioning the response on an inappropriate emotion or intent can lead to responses that are empathetically inappropriate with the dialogue context.

in responses that are empathetically inappropriate with the dialogue context. The neural predictor, as well as the decision tree-based mechanisms, generate some emotion or intent that is appropriate to the dialogue context, enabling the generation model to generate responses that are more empathetically appropriate, guiding the conversation in a meaningful direction.

## 5 Discussion and Conclusion

This study investigated the use of a taxonomy of empathetic response intents along with 32 fine-grained emotions in controlling and interpreting the responses generated by open-domain conversational agents for emotional prompts. In this regard, several rule-based and automatic response control methods were proposed and were compared in terms of their prediction and generation performance on two state-of-the-art dialogue datasets containing emotional dialogues.

It was observed that the neural response emotion/intent predictor we proposed outperformed the rest including the end-to-end model in terms of evaluation metrics related to both prediction and generation. This implies the importance of leveraging semantic clues in addition to the flow of emotions and intents in the previous turns when predicting the next response's emotion or intent. However, there are some disadvantages to using this approach: 1) developers cannot foresee the la-

bel that the model would predict next; and 2) cost of time and resources spent for training the model. As a remedy, we proposed two decision tree-based response emotion/intent prediction approaches.

Across evaluation metrics for prediction and generation, the performance of the decision-tree methods was considerably better than the end-to-end approach and the equally sampled baseline. The decision tree (argmax) method performed the best in terms of diversity metrics related to response generation. In the human evaluation stage, we saw that the DT (argmax) method produced the most number of combined *Good* and *Okay* responses in ED and EDOS test datasets, pointing to the fact that the rule-based approaches we proposed can still be used without a significant degrade in performance in resource-limited environments.

On the whole, the results of this study inform developers about the utility of the taxonomy of empathetic response intents in controlling the responses generated by open-domain chatbots and which optimal methodology to use (rule-based or automatic conditioning) based on the operational environment.

## References

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 Pro-*



590		<i>cessing and Chinese Computing</i> , pages 117–129.	
591		Springer.	
592	Artur d’Avila Garcez and Luis C. Lamb. 2020.	<a href="#">Neurosymbolic ai: The 3rd wave.</a>	646
593			647
594	Jacob Devlin, Ming-Wei Chang, Kenton Lee, and		648
595	Kristina Toutanova. 2019. BERT: Pre-training of		649
596	deep bidirectional transformers for language under-		650
597	standing. In <i>Proceedings of the 2019 Conference</i>		
598	<i>of the North American Chapter of the Association</i>		
599	<i>for Computational Linguistics: Human Language</i>		
600	<i>Technologies, Volume 1 (Long and Short Papers)</i> ,		
601	pages 4171–4186, Minneapolis, Minnesota. Associ-		
602	ation for Computational Linguistics.		
603	Paul Ekman. 1992. An argument for basic emotions.		
604	<i>Cognition &amp; emotion</i> , 6(3-4):169–200.		
605	Jianfeng Gao, Michel Galley, and Lihong Li. 2018.		
606	Neural approaches to conversational ai. In <i>The</i>		
607	<i>41st International ACM SIGIR Conference on Re-</i>		
608	<i>search &amp; Development in Information Retrieval</i> ,		
609	pages 1371–1374.		
610	Karthik Gopalakrishnan, Behnam Hedayatnia, Qin-		
611	lang Chen, Anna Gottardi, Sanjeev Kwatra, Anu		
612	Venkatesh, Raefer Gabriel, and Dilek Hakkani-		
613	Tür. 2019. Topical-Chat: Towards Knowledge-		
614	Grounded Open-Domain Conversations. In <i>INTER-</i>		
615	<i>SPEECH</i> .		
616	Prakhar Gupta, Jeffrey P Bigham, Yulia Tsvetkov,		
617	and Amy Pavel. 2020. Controlling dialogue gen-		
618	eration with semantic exemplars. <i>arXiv preprint</i>		
619	<i>arXiv:2008.09075</i> .		
620	Elaine Hatfield, John T Cacioppo, and Richard L Rap-		
621	son. 1993. Emotional contagion. <i>Current directions</i>		
622	<i>in psychological science</i> , 2(3):96–100.		
623	Behnam Hedayatnia, Karthik Gopalakrishnan,		
624	Seokhwan Kim, Yang Liu, Mihail Eric, and		
625	Dilek Hakkani-Tur. 2020. <a href="#">Policy-driven neural</a>		
626	<a href="#">response generation for knowledge-grounded dialog</a>		
627	<a href="#">systems</a> . In <i>Proceedings of the 13th International</i>		
628	<i>Conference on Natural Language Generation</i> ,		
629	pages 412–421, Dublin, Ireland. Association for		
630	Computational Linguistics.		
631	Tianran Hu, Anbang Xu, Zhe Liu, Quanzeng You, Yu-		
632	fan Guo, Vibha Sinha, Jiebo Luo, and Rama Akki-		
633	raju. 2018. Touch your heart: A tone-aware chatbot		
634	for customer care on social media. In <i>Proceedings of</i>		
635	<i>the 2018 CHI conference on human factors in com-</i>		
636	<i>puting systems</i> , pages 1–12.		
637	Pei Ke, Jian Guan, Minlie Huang, and Xiaoyan Zhu.		
638	2018. <a href="#">Generating informative responses with con-</a>		
639	<a href="#">trolled sentence function</a> . In <i>Proceedings of the</i>		
640	<i>56th Annual Meeting of the Association for Compu-</i>		
641	<i>tational Linguistics (Volume 1: Long Papers)</i> , pages		
642	1499–1508, Melbourne, Australia. Association for		
643	Computational Linguistics.		
	Dave Lee. 2016. <a href="#">Tay: Microsoft issues apology over</a>		644
	<a href="#">racist chatbot fiasco</a> .		645
	Hung-yi Lee, Cheng-Hao Ho, Chien-Fu Lin, Chiung-		646
	Chih Chang, Chih-Wei Lee, Yau-Shian Wang,		647
	Tsung-Yuan Hsu, and Kuan-Yu Chen. 2020. Invest-		648
	igation of sentiment controllable chatbot. <i>arXiv</i>		649
	<i>preprint arXiv:2007.07196</i> .		650
	Pierre Lison, Jörg Tiedemann, Milen Kouylekov, et al.		651
	2019. Open subtitles 2018: Statistical rescoring of		652
	sentence alignments in large, noisy parallel corpora.		653
	In <i>LREC 2018, Eleventh International Conference</i>		654
	<i>on Language Resources and Evaluation</i> . European		655
	Language Resources Association (ELRA).		656
	Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Man-		657
	dar Joshi, Danqi Chen, Omer Levy, Mike Lewis,		658
	Luke Zettlemoyer, and Veselin Stoyanov. 2019.		659
	Roberta: A robustly optimized bert pretraining ap-		660
	proach. <i>arXiv preprint arXiv:1907.11692</i> .		661
	Robert Plutchik. 1984. Emotions: A general psy-		662
	choevolutionary theory. <i>Approaches to emotion</i> ,		663
	1984:197–219.		664
	Hannah Rashkin, Eric Michael Smith, Margaret Li, and		665
	Y-Lan Boureau. 2018. Towards empathetic open-		666
	domain conversation models: A new benchmark and		667
	dataset. <i>arXiv preprint arXiv:1811.00207</i> .		668
	Chinnadhurai Sankar and Sujith Ravi. 2019. <a href="#">Deep re-</a>		669
	<a href="#">inforcement learning for modeling chit-chat dialog</a>		670
	<a href="#">with discrete attributes</a> . In <i>Proceedings of the 20th</i>		671
	<i>Annual SIGdial Meeting on Discourse and Dialogue</i> ,		672
	pages 1–10, Stockholm, Sweden. Association for		673
	Computational Linguistics.		674
	Sashank Santhanam, Zhuo Cheng, Brodie Mather, Bon-		675
	nie Dorr, Archana Bhatia, Bryanna Hebenstreit, Alan		676
	Zemel, Adam Dalton, Tomek Strzalkowski, and		677
	Samira Shaikh. 2020. <a href="#">Learning to plan and real-</a>		678
	<a href="#">ize separately for open-ended dialogue systems</a> . In		679
	<i>Findings of the Association for Computational Lin-</i>		680
	<i>guistics: EMNLP 2020</i> , pages 2736–2750, Online.		681
	Association for Computational Linguistics.		682
	Lifeng Shang, Zhengdong Lu, and Hang Li. 2015. <a href="#">Neu-</a>		683
	<a href="#">ral responding machine for short-text conversation</a> .		684
	In <i>Proceedings of the 53rd Annual Meeting of the</i>		685
	<i>Association for Computational Linguistics and the</i>		686
	<i>7th International Joint Conference on Natural Lan-</i>		687
	<i>guage Processing (Volume 1: Long Papers)</i> , pages		688
	1577–1586, Beijing, China. Association for Compu-		689
	tational Linguistics.		690
	Amy E Skerry and Rebecca Saxe. 2015. Neural repre-		691
	sentations of emotion are organized around abstract		692
	event features. <i>Current biology</i> , 25(15):1945–1954.		693
	Zhenqiao Song, Xiaoqing Zheng, Lu Liu, Mu Xu, and		694
	Xuanjing Huang. 2019. <a href="#">Generating responses with</a>		695
	<a href="#">a specific emotion in dialog</a> . In <i>Proceedings of the</i>		696
	<i>57th Annual Meeting of the Association for Com-</i>		697
	<i>putational Linguistics</i> , pages 3685–3695, Florence,		698
	Italy. Association for Computational Linguistics.		699

700	Alessandro Sordani, Michel Galley, Michael Auli,	Xianda Zhou and William Yang Wang. 2018. <a href="#">Mo-</a>	756
701	Chris Brockett, Yangfeng Ji, Margaret Mitchell,	<a href="#">jiTalk: Generating emotional responses at scale</a> . In	757
702	Jian-Yun Nie, Jianfeng Gao, and Bill Dolan. 2015.	<i>Proceedings of the 56th Annual Meeting of the As-</i>	758
703	<a href="#">A neural network approach to context-sensitive gen-</a>	<i>sociation for Computational Linguistics (Volume 1:</i>	759
704	<a href="#">eration of conversational responses</a> . In <i>Proceedings</i>	<i>Long Papers)</i> , pages 1128–1137, Melbourne, Aus-	760
705	<i>of the 2015 Conference of the North American Chap-</i>	tralia. Association for Computational Linguistics.	761
706	<i>ter of the Association for Computational Linguis-</i>		
707	<i>tics: Human Language Technologies</i> , pages 196–		
708	205, Denver, Colorado. Association for Computa-		
709	tional Linguistics.		
710	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob		
711	Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz		
712	Kaiser, and Illia Polosukhin. 2017. Attention is all		
713	you need. In <i>Advances in Neural Information Pro-</i>		
714	<i>cessing Systems</i> , volume 30. Curran Associates, Inc.		
715	Oriol Vinyals and Quoc Le. 2015. A neural conversa-		
716	tional model. In <i>ICML Deep Learning Workshop</i> .		
717	Anuradha Welivita and Pearl Pu. 2020. A taxon-		
718	omy of empathetic response intents in human so-		
719	cial conversations. In <i>Proceedings of the 28th Inter-</i>		
720	<i>national Conference on Computational Linguistics</i> ,		
721	pages 4886–4899.		
722	Anuradha Welivita, Yubo Xie, and Pearl Pu. 2020.		
723	<a href="#">Fine-grained emotion and intent learning in movie</a>		
724	<a href="#">dialogues</a> .		
725	Anuradha Welivita, Yubo Xie, and Pearl Pu. 2021. A		
726	large-scale dataset for empathetic response genera-		
727	tion. In <i>Findings of the Association for Computa-</i>		
728	<i>tional Linguistics: EMNLP 2021</i> . Association for		
729	Computational Linguistics.		
730	Wei Wu, Can Xu, Yu Wu, and Zhoujun Li. 2018.		
731	<a href="#">Towards interpretable chit-chat: Open domain dia-</a>		
732	<a href="#">logue generation with dialogue acts</a> .		
733	Wei Wu and Rui Yan. 2018. Deep chit-chat: Deep		
734	learning for ChatBots. In <i>Proceedings of the 2018</i>		
735	<i>Conference on Empirical Methods in Natural Lan-</i>		
736	<i>guage Processing: Tutorial Abstracts</i> , Melbourne,		
737	Australia. Association for Computational Linguis-		
738	tics.		
739	Zeqiu Wu, Michel Galley, Chris Brockett, Yizhe Zhang,		
740	Xiang Gao, Chris Quirk, Rik Koncel-Kedziorski,		
741	Jianfeng Gao, Hannaneh Hajishirzi, Mari Osten-		
742	dorf, et al. 2020. A controllable model of		
743	grounded response generation. <i>arXiv preprint</i>		
744	<i>arXiv:2005.00613</i> .		
745	Yubo Xie and Pearl Pu. 2021. Generating empathetic		
746	responses with a large scale dialog dataset. In <i>The</i>		
747	<i>SIGNLL Conference on Computational Natural Lan-</i>		
748	<i>guage Learning (CoNLL)</i> . Association for Computa-		
749	tional Linguistics.		
750	Hao Zhou, Minlie Huang, Tianyang Zhang, Xiaoyan		
751	Zhu, and Bing Liu. 2018. Emotional chatting ma-		
752	chine: Emotional conversation generation with in-		
753	ternal and external memory. In <i>Proceedings of</i>		
754	<i>the AAAI Conference on Artificial Intelligence</i> , vol-		
755	ume 32.		