# Towards Empathetic Open-domain Conversation Models: a New Benchmark and Dataset

**Hannah Rashkin**[1][*], **Eric Michael Smith**[2], **Margaret Li**[2], **Y-Lan Boureau**[2]

[1] Paul G. Allen School of Computer Science & Engineering, University of Washington

[2] Facebook AI Research

hrashkin@cs.washington.edu, {ems,margaretli,ylan}@fb.com

## Abstract

One challenge for dialogue agents is recognizing feelings in the conversation partner and replying accordingly, a key communicative skill. While it is straightforward for humans to recognize and acknowledge others' feelings in a conversation, this is a significant challenge for AI systems due to the paucity of suitable publicly-available datasets for training and evaluation. This work proposes a new benchmark for empathetic dialogue generation and EMPATHETICDIALOGUES, a novel dataset of 25k conversations grounded in emotional situations. Our experiments indicate that dialogue models that use our dataset are perceived to be more empathetic by human evaluators, compared to models merely trained on large-scale Internet conversation data. We also present empirical comparisons of dialogue model adaptations for empathetic responding, leveraging existing models or datasets without requiring lengthy retraining of the full model.

## 1 Introduction

A desirable trait in a human-facing dialogue agent is to appropriately respond to a conversation partner that is describing personal experiences, by understanding and acknowledging any implied feelings — a skill we refer to as empathetic responding. For instance, while the crossed-out response in Figure 1 is topically relevant, "Congrats! That's great!" may be more satisfying because it acknowledges the underlying feelings of accomplishment in an empathetic way. In this work, we investigate empathetic response generation from current dialogue systems, and propose experiments using a new resource, EMPATHETICDIALOGUES, as a benchmark to evaluate this skill set.
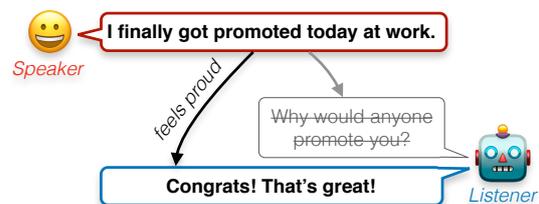
**EMPATHETICDIALOGUES** dataset example

Figure 1: Example where acknowledging an inferred feeling is appropriate

Empathetic responding is clearly relevant to dialogue systems that are geared towards general conversation or chit-chat. Indeed, ordinary communication is frequently prompted by people sharing their feelings or circumstances. But researchers analyzing goal-directed conversations have also observed the frequent intrusion of ordinary conversation in those interactions as well, either as a "warm-up" introduction or as a detour (Levinson et al., 2000; Heritage, 2005). Engaging in social talk, reacting to emotional cues and displaying a caring attitude have, in fact, been associated with better task outcomes in many domains (Wentzel, 1997; Levinson et al., 2000; Bickmore and Cassell, 2001; Kim et al., 2004; Fraser et al., 2018). While many of those studies deal with human-human interactions, humans have been shown to often interact with machines in a natural and social way (Reeves and Nass, 1996; Lee et al., 2010), so it is reasonable to expect that dialogue agents would also benefit from empathetic responding.

Most recent powerful language architectures are trained on vast amounts of barely curated text scrapes, social media conversations, or independent books (Ritter et al., 2010; Zhang et al., 2018; Mazare et al., 2018; Devlin et al., 2018; Liu et al., 2019; Radford et al., 2019). It might be the case

| **Label: Afraid** | **Label: Proud** |
|---|---|
| **Situation:** Speaker felt this when... | **Situation:** Speaker felt this when... |
| "I've been hearing noises around the house at night" | "I finally got that promotion at work! I have tried so hard for so long to get it!" |
| **Conversation:** | **Conversation:** |
| Speaker: I've been hearing some strange noises around the house at night. | Speaker: I finally got promoted today at work! |
| Listener: oh no! That's scary! What do you think it is? | Listener: Congrats! That's great! |
| Speaker: I don't know, that's what's making me anxious. | Speaker: Thank you! I've been trying to get it for a while now! |
| Listener: I'm sorry to hear that. I wish I could help you figure it out | Listener: That is quite an accomplishment and you should be proud! |

Figure 2: Two examples from EMPATHETICDIALOGUES training set. The first worker (the speaker) is given an emotion label and writes their own description of a situation when they've felt that way. Then, the speaker tells their story in a conversation with a second worker (the listener).

that models trained on this type of data could exhibit some of the aggressive and callous responses that have been observed in spontaneous internet conversations (Anderson, 2015). Unfortunately, while chitchat dialogue benchmarks have been proposed (e.g., Dinan et al., 2019), to the best of our knowledge there are currently no benchmarks gauging whether dialogue agents can converse with empathy.

This work aims to facilitate evaluating models' ability to produce empathetic responses. We introduce a new task for dialogue systems to respond to people discussing situations that cover a wide range of emotions, and EMPATHETICDIALOGUES (ED), a novel dataset with about 25k personal dialogues. Each dialogue is grounded in a specific situation where a speaker was feeling a given emotion, with a listener responding (Figure 2). The new resource consists of crowdsourced one-on-one conversations, and covers a large set of emotions in a balanced way. This dataset is larger and contains a more extensive set of emotions than many similar emotion prediction datasets from other text domains such as Scherer and Wallbott (1994), Strapparava and Mihalcea (2007), Mohammad et al. (2018), and Gupta et al. (2017). The dataset has been publicly released, with code to reproduce the main experimental results of this paper[1].

Our experiments show that large-capacity conversation models trained on spontaneous internet conversation data are not rated as very empathetic. We propose two simple ways to leverage our dataset to improve those models: use utterances from our training data as candidate responses in a retrieval model at inference time, and fine-tune the model on our task. Finally, we explore whether

different ways of combining information from related tasks can lead to more empathetic responses. The contributions of this work are thus: 1) we release a novel empathetic dialogue dataset as a new benchmark; 2) we show that training over this dataset can improve the performance of an end-to-end dialogue system on empathetic dialogue.

## 2 Related Work

**Emotion data** Crafting our dataset requires deciding what set of emotions the models should be capable of reacting to. Multiple schemas have attempted to organize the spectrum of emotions, from a handful of basic emotions derived from biological responses (Ekman, 1992; Plutchik, 1984) to larger sets of subtle emotions inferred from contextual situations (Skerry and Saxe, 2015). We incorporate emotions from multiple annotation schemas, noting that emotions merely inferred from a situation are important in dialogue scenarios. There is a wide breadth of research in distributional representation approaches for many emotion classification tasks (Duppada et al., 2018; Park et al., 2018; Xu et al., 2018; Mohammad et al., 2018) that build on deep networks pretrained on large-scale weakly-labelled data such as emojis (Felbo et al., 2017) or hashtags (Mohammad, 2012), gathered from public social media content published on Twitter. The SEMEVAL2019 EmoContext challenge also uses conversation data for detection of three basic emotions ('happy', 'sad', and 'angry') over two turns of context from Twitter exchanges (Gupta et al., 2017). We focus on personal conversations rather than using social media data to be closer to a context of a one-on-one conversation. Public social media content occurs in front of large "peripheral audiences" (Goffman, 1981) where uncertainty as to how wide

that audience is and the need for curated self-presentation (Goffman, 1959) have been shown to lead to different choices of subject matters compared to private messaging, with people sharing more intense and negative emotions through private channels (Bazarova et al., 2015; Litt et al., 2014). In this work, we generate a more balanced coverage of emotions than would appear in public social media content, using a domain that is closer to our ultimate goal of training a model for conversation that can respond to any emotion.

**Controllable language generation** Several other works have focused on controlling the emotional content of a text response either through a manually specified target (Zhou and Wang, 2018; Zhou et al., 2018; Wang and Wan, 2018; Hu et al., 2017; Huang et al., 2018) or through a general term to encourage higher levels of affect (Asghar et al., 2018), with evaluations focused on matching a predetermined desired emotion rather than empathetic responding. Niu and Bansal (2018) generate responses conditioned on a specified politeness setting (polite, rude or neutral). Huber et al. (2018) investigate how to respond to emotions detected from an image. Our work focuses on empathetic responses that are appropriate to signals inferred purely from text rather than conveying a pre-specified emotion.

**Related chit-chat data** Several works have attempted to make chit-chat dialogue models more engaging by grounding them in personal contexts (Li et al., 2016b; Zhang et al., 2018; Mazare et al., 2018), focusing on personal facts ("I am from New York"). Another interesting resource is the DAILYDIALOG (DD) dataset (Li et al., 2017), which comprises about 13k dialogues obtained by crawling educational websites intended for learners of English and also has emotion label annotations. Many of the dialogues are focused on topics for ESL learners (ordering from a restaurant, asking for directions, introductions, etc), but only $\approx 5\%$ of the utterances have a label other than "none" or "happy". Our task focuses explicitly on conversations about emotionally grounded personal situations, and considers a richer, evenly distributed set of emotions. We also introduce an explicit single *listener* in the conversation who is reacting to the situation being described in an empathetic way, to make the setting as close as possible to our desired goal of a one-on-one empathetic conversation.

| Emotion | Most-used speaker words | Most-used listener words | Training set emotion distrib |
|---|---|---|---|
| Surprised | got,shocked,really | that's,good,nice | 5.1% |
| Excited | going,wait,i'm | that's,fun,like | 3.8% |
| Angry | mad,someone,got | oh,would,that's | 3.6% |
| Proud | got,happy,really | that's,great,good | 3.5% |
| Sad | really,away,get | sorry,oh,hear | 3.4% |
| Annoyed | get,work,really | that's,oh,get | 3.4% |
| Grateful | really,thankful,i'm | that's,good,nice | 3.3% |
| Lonely | alone,friends,i'm | i'm,sorry,that's | 3.3% |
| Afraid | scared,i'm,night | oh,scary,that's | 3.2% |
| Terrified | scared,night,i'm | oh,that's,would | 3.2% |
| Guilty | bad,feel,felt | oh,that's,feel | 3.2% |
| Impressed | really,good,got | that's,good,like | 3.2% |
| Disgusted | gross,really,saw | oh,that's,would | 3.2% |
| Hopeful | i'm,get,really | hope,good,that's | 3.2% |
| Confident | going,i'm,really | good,that's,great | 3.2% |
| Furious | mad,car,someone | oh,that's,get | 3.1% |
| Anxious | i'm,nervous,going | oh,good,hope | 3.1% |
| Anticipating | wait,i'm,going | sounds,good,hope | 3.1% |
| Joyful | happy,got,i'm | that's,good,great | 3.1% |
| Nostalgic | old,back,really | good,like,time | 3.1% |
| Disappointed | get,really,work | oh,that's,sorry | 3.1% |
| Prepared | ready,i'm,going | good,that's,like | 3% |
| Jealous | friend,got,get | get,that's,oh | 3% |
| Content | i'm,life,happy | good,that's,great | 2.9% |
| Devastated | got,really,sad | sorry,oh,hear | 2.9% |
| Embarrassed | day,work,got | oh,that's,i'm | 2.9% |
| Caring | care,really,taking | that's,good,nice | 2.7% |
| Sentimental | old,really,time | that's,oh,like | 2.7% |
| Trusting | friend,trust,know | good,that's,like | 2.6% |
| Ashamed | feel,bad,felt | oh,that's,i'm | 2.5% |
| Apprehensive | i'm,nervous,really | oh,good,well | 2.4% |
| Faithful | i'm,would,years | good,that's,like | 1.9% |

Figure 3: Distribution of conversation labels within EMPATHETICDIALOGUES training set and top 3 content words used by speaker/listener per category.

## 3 Talking about Personal Situations

We consider an open-domain one-on-one conversational setting where two people are discussing a situation that happened to one of them, related to a given feeling. We collect around 25k conversations using the following format.

**Emotional situation grounding** Each conversation is grounded in a situation, which one participant writes about in association with a given emotion label. We consider 32 emotion labels, listed in Figure 3, which we chose by aggregating labels from several emotion prediction datasets (Scherer and Wallbott, 1994; Strapparava and Mihalcea, 2007; Skerry and Saxe, 2015; Li et al., 2017; Mohammad, 2012). These emotion labels cover a broad range of positive and negative emotions. Our goal in providing a single emotion label is to have a situation strongly related to (at least) one particular emotional experience, though we note that some emotions may be very closely related[2] and additional related emotions may be

---
[2]Researchers could merge similar emotions, like "afraid" and "terrified", to get coarser labels, if desired.

invoked in a given conversation.

**Speaker and listener** The person who wrote the situation description (*Speaker*) initiates a conversation to talk about it. The other conversation participant (*Listener*) becomes aware of the underlying situation through what the Speaker says and responds. Speaker and Listener then exchange up to 6 more turns. We include two example conversations from the training data in Figure 2 and ten more in Table 5 in the Appendix. The models discussed below are tested in the role of *Listener* responding to the Speaker. Neither the situation description written by the Speaker nor the emotion label is given to the models (just as they were not given to the Listener during dialogue collection). Our data could also be used to generate conversations for the Speaker conditioned on the situation description though we leave this for future work.

**Collection details** We collected crowdsourced dialogues using the ParlAI platform (Miller et al., 2017) to interact with Amazon Mechanical Turk (MTurk), hiring 810 US workers. A pair of workers are asked to (i) select an emotion word each and describe a situation when they felt that way, and to (ii) have a conversation about each of the situations, as outlined below. Each worker had to contribute at least one situation description and one pair of conversations: one as Speaker about the situation they contributed, and one as Listener about the situation contributed by another worker. They were allowed to participate in as many hits as they wanted for the first ∼10k conversations, then we limited the more "frequently active" workers to a maximum of 100 conversations. The median number of conversations per worker was 8, while the average was 61 (some workers were more active contributors than others). To ensure quality, we manually checked random subsets of conversations by our most-frequent workers.

**Task set-up** In the first stage of the task, workers are asked to describe in a few sentences a situation based on a feeling label. We ask the workers to try to keep these descriptions between 1-3 sentences. The average response is 19.8 words. In the second stage, two workers are paired and asked to have two short chats with each other. In each chat, one worker (*speaker*) starts a conversation about the situation they previously described, and the other worker (*listener*) responds. Neither can see what the other worker was given as emotion

label or the situation description they submitted, so they must respond to each others' stories based solely on cues within the conversation. Each conversation is allowed to be 4-8 utterances long (the average is 4.31 utterances per conversation). The average utterance length was 15.2 words long.

**Ensuring balanced emotion coverage** After the first few initial rounds of data collection, we forced workers to select an emotion that among three emotion labels that had been the least chosen overall so far if it was their first time working on the task. If they had already performed the task, the offered emotion labels were among those that they had chosen the least often before. Given that a conversation model trained for empathetic responding needs to be able to handle emotions even if they are less frequent, we opted for this balancing procedure to make training for these categories easier, while still allowing for some measure of choice for workers. As shown in Figure 3, the distribution of emotion label prompts is close to evenly distributed, with a few that are selected slightly more/less often.

**EMPATHETICDIALOGUES dataset statistics** The resulting dataset comprises 24,850 conversations about a situation description, gathered from 810 different participants, which are publicly available through the ParlAI framework[3] and for direct download with accompanying code[4]. We split the conversations into approximately 80% train, 10% validation, and 10% test partitions. To prevent overlap of discussed situations between partitions, we split the data so that all sets of conversations with the same speaker providing the initial situation description would be in the same partition. The final train/val/test split was 19533 / 2770 / 2547 conversations, respectively. We include ten examples from our training set in Appendix Section A.

## 4 Empathetic Response Generation

This section shows how ED can be used as a benchmark to gauge the ability of a model to respond in an empathetic way, and as a training resource to make generic chitchat models more empathetic. We also examine different ways existing models can be combined to produce more empathetic responses. We use ED dialogues to train

---

[3]https://parl.ai/
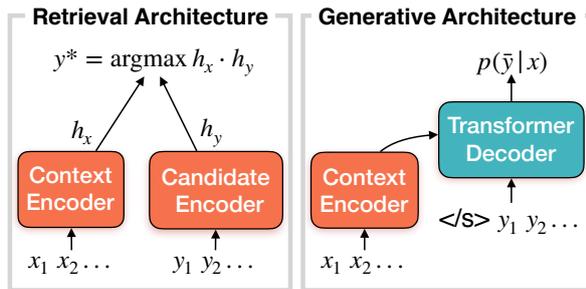[4]https://github.com/facebookresearch/EmpatheticDialogues

Figure 4: Dialogue generation architectures used in our experiments. The context of concatenated previous utterances is tokenized into $x_1, x_2, \cdots$, and encoded into vector $h_x$ by the context encoder. *Left:* In the retrieval set-up, each candidate $y$ is tokenized into $y_1, y_2, \cdots$ and encoded into vector $h_y$ by the candidate encoder. The system outputs the candidate $y^*$ that maximizes dot product $h_x \cdot h_y$. *Right:* In the generative set-up, the encoded context $h_x$ is used as input to the decoder to generate start symbol $</\text{s}>$ and tokens $y_1, y_2, \cdots$. The model is trained to minimize the negative log-likelihood of target sequence $\bar{y}$ conditioned on context.

and evaluate models in the task of generating conversation responses in the *Listener* role. To emulate a normal conversation, the model has access to previous utterances in the dialogue, but not to the emotion word prompt (e.g., "proud"), nor to the situation description generated by the Speaker. Given a dialogue context $x$ of $n$ previous conversation utterances concatenated and tokenized as $x_1, \cdots, x_m$, followed by a target response $\bar{y}$, our models are trained to maximize the likelihood $p(\bar{y}|x)$ of producing the target response. We investigate both generative and retrieval-based settings (Lowe et al., 2016) as described in Figure 4.

## 4.1 Base Architecture

We base our models on Transformer networks (Vaswani et al., 2017), which have proven successful in machine translation and dialogue generation tasks (Zhang et al., 2018; Mazare et al., 2018).

**Retrieval-based** In the retrieval-based set-up, the model is given a large set $Y$ of candidate responses and picks the "best" one, $y^*$. We first experiment with the retrieval Transformer-based architecture from Yang et al. (2018): two Transformer encoders separately embedding the context, $x$, and candidates, $y \in Y$, as $h_x$ and $h_y$, respectively. We also experiment with BERT (Devlin et al., 2018) as base architecture to encode candidates and contexts, using the final hidden vector from BERT as the $h_x$ or $h_y$ encodings. The

model chooses a candidate utterance according to a softmax on the dot product: $h_x \cdot h_y$. We minimize the negative log-likelihood of selecting the correct candidate. At training time, we use all of the utterances from the batch as candidates, with a large batch size of 512 to give the model more negative examples (except for BERT for which a batch size of 256 was used). At inference time, we experiment with three sets of candidate utterances for the model to choose from: all of the response utterances in the ED training set ($Y^{ED}$), all the utterances in the DailyDialog (Li et al., 2017) training set ($Y^{DD}$), and a million utterances from a dump of 1.7 billion Reddit (R) conversations ($Y^R$).

**Generative** In the generative set-up, we use the full Transformer architecture (Vaswani et al., 2017), consisting of an encoder and a decoder. The Transformer decoder uses the encoder output to predict a sequence of words $y$, and is trained to minimize the negative log-likelihood of the target sequence $\bar{y}$. At inference time, we use diverse beam search from Vijayakumar et al. (2016).

**Training details** Models are pretrained on predicting replies from a dump of 1.7 billion Reddit conversations, starting either from scratch for the Transformer architectures, or from the BERT$_{base}$ model released by Devlin et al. (2018) for the BERT-based architectures.[5] Pretrained models without any fine-tuning on ED will be referred to as "Pretrained" hereafter. We limit the maximum number of word tokens in the context and response to be 100 each. The Transformer networks used in most experiments have the same base architecture (four layers and six transformer heads) and are trained the same way as in Mazare et al. (2018). We also experiment with a larger architecture of five layers (denoted as "Large"), and BERT retrieval models, that are allowed to train for much longer (see training times in Table 3).[6] For all models, we keep the version that has the lowest loss on the validation set. For the Transformer models, we use 300-d word embed-

---

[5] We used the Hugging Face PyTorch implementation of BERT at https://github.com/huggingface/pytorch-transformers. We experimented with directly fine-tuning BERT on ED without first training on Reddit conversations, but this did not perform as well.

[6] While the models had not fully converged when we stopped training, we trained the Pretrained models for a few iterations more than the corresponding Fine-Tuned models, to ensure that any observed improvement was due to the data used for fine-tuning and not the extra training time.
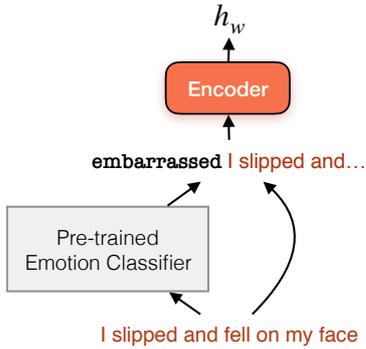
Figure 5: Incorporating additional supervised information, here from an emotion classification task. An input sequence (either a dialogue context or a candidate) is run through a pre-trained classifier, and the top $k$ output labels are prepended to the sequence, which is then run through the corresponding (context or candidate) encoder to output a hidden representation $h_w$ (either $h_x$ or $h_y$) as in the base setting.

dings pretrained on common-crawl data using fast-Text (Grave et al., 2018), and for the BERT models, we use 768-d word embeddings pretrained on BooksCorpus and English Wikipedia (Devlin et al., 2018). More training details are provided in Appendix D.1.

## 4.2 Leveraging the Training Data from ED

A retrieval-based model relies on candidates. ED data was explicitly collected with instructions to be empathetic, in a one-on-one setting, which is not the case of the Reddit conversation data used for pretraining, and these domain candidates may be better suited to empathetic responding than generic conversation utterances. Thus, we experiment with incorporating ED training candidates into the pool used at inference time by pretrained retrieval-based models, with no fine-tuning on ED. For retrieval-based and generative models, we also experiment with fine-tuning pretrained models to predict the next utterance over ED with a context window of four previous utterances, which is the average length of a conversation in our dataset. These models are referred to as "Fine-Tuned" models. This fine-tuning is conducted until convergence for all architectures except those referred to as "Pretrained".

## 4.3 Adding Information from External Predictors

Many existing models have been pretrained on supervised tasks that may be relevant to empathetic responding. Combining these models with the rep-

resentations from our base architecture may reap benefits from previous training time and external training data without having to redo the work or requiring access to that data, which may matter to practitioners. Note that this may considerably augment the effective capacity of the resulting models, as well as the total amount of training data used overall, but our goal here is to get an empirical sense of how robust performance improvement is to variations in architecture set-up or supervision domain. We experiment with adding supervised information from two prediction tasks: emotion detection, which is more closely relevant to our task, and topic detection, which may also be useful in crafting relevant replies.[7]

**Prepending Top-k Predicted Labels**  This set-up (Fig. 5), PREPEND-1, is a very simple way to add supervised information to data, requires no architecture modification, and can be used with black-box classifiers. The top predicted label[8] from the supervised classifier is merely prepended to the beginning of the token sequence as encoder input, as below:

> **Original**:"I finally got promoted!"
> **Prepend-1**:"`proud` I finally got promoted!"

Similar methods have been used for controlling the style of generated text (e.g. Niu and Bansal, 2018). Here, we use a fastText model (Joulin et al., 2017) as prediction architecture. Both the context and the candidates are run through the classifier and receive prepended labels. Fine-tuning is conducted similarly as before, but using these modified inputs. We use two external sources of information. To provide emotion signal, we train a classifier to predict the emotion label from the description of the situation written by the Speaker before the dialogue for the training set dialogues of ED (EMOPREPEND-1).[9] To gauge whether supervision from a more distant task would still be helpful, we also experiment with a classifier trained on the 20-Newsgroup dataset (Joachims, 1996), for topic classification (TOPICPREPEND-1).

---

[7]We considered multitask or feature concatenation set-ups, but they did not provide consistent improvements. These experiments are included in Appendix D.2.

[8]We only discuss prepending the top predicted label here, but also experimented with top-3 and top-5 models, with similar result patterns, shown in Appendix D.3.

[9]We also experimented with training the classifier on the utterances themselves, with similar results.

| | | Retrieval | | Retrieval w/ BERT | | Generative | |
|---|---|---|---|---|---|---|---|
| Model | Candidate Source | P@1,100 | AVG BLEU | P@1,100 | AVG BLEU | PPL | AVG BLEU |
| Pretrained | R | - | 4.10 | - | 4.26 | 27.96 | 5.01 |
| | ED | 43.25 | 5.51 | 49.94 | 5.97 | - | - |
| Fine-Tuned | ED | **56.90** | 5.88 | 65.92 | **6.21** | **21.24** | **6.27** |
| | ED+DD | - | 5.61 | - | - | - | - |
| | ED+DD+R | - | 4.74 | - | - | - | - |
| EmoPrepend-1 | ED | 56.31 | 5.93 | **66.04** | 6.20 | 24.30 | 4.36 |
| TopicPrepend-1 | ED | 56.38 | **6.00** | 65.96 | 6.18 | 25.40 | 4.17 |

Table 1: Automatic evaluation metrics on the test set. Pretrained: model pretrained on a dump of 1.7 billion RED-DIT conversations (4-layer Transformer architecture, except when specified BERT). Fine-Tuned: model fine-tuned over the EMPATHETICDIALOGUES training data (Sec. 4.2). EmoPrepend-1, Topic-Prepend1: model incorporating supervised information from an external classifiers, as described in Sec. 4.3. Candidates come from REDDIT (R), EMPATHETICDIALOGUES (ED), or DAILYDIALOG (DD). P@1,100: precision retrieving the correct test candidate out of 100 test candidates. AVG BLEU: average of BLEU-1,-2,-3,-4. PPL: perplexity. All automatic metrics clearly improve with in-domain training on utterances (Fine-Tuned vs. Pretrained), other metrics are inconsistent. *Bold: best performance for that architecture.*

## 5 Experimental Evaluation

We evaluate the models on their ability to reproduce the Listener's portion of the conversation (i.e. the ability to react to someone else's story). We use both automated metrics and human evaluation to score each model's retrievals/generations. Human evaluation is important, as automated metrics don't always correlate with human judgments of dialogue quality (Liu et al., 2016), but we provide automated metrics to give a sense of how well they align with human judgment on this task.

**Automated metrics (Table 1)** For both retrieval and generative systems, we compute BLEU scores (Papineni et al., 2002) for the model response, comparing against the gold label (the actual response), following the practice of earlier work in dialogue generation (Wen et al., 2015; Li et al., 2016a,b). For the generative systems, we additionally report perplexity of the actual gold response. For the retrieval-based systems, we further compute p@1,100, the accuracy of the model at choosing the correct response out of a hundred randomly selected examples in the test set. When we compute p@1,100, the actual response is included in the candidates, unlike inference from the retrieval systems for all other metrics, which only uses training utterances as candidates.

**Human ratings (Table 2)** We ran crowd-sourcing tasks on MTurk (further details in

Appendix B). Participants were given a model's output for a randomly selected test set example and asked to score different aspects of the model. The rating task provides a means of comparing aspects of responses, and we ask raters specifically about whether the response is acknowledging the conversation partner's feelings. We collected at least 100 ratings per model and asked about three aspects of performance, all rated on a Likert scale (1: not at all, 3: somewhat, 5: very much):

**Empathy/Sympathy:** did the responses show understanding of the feelings of the person talking about their experience?

**Relevance:** did the responses seem appropriate to the conversation? Were they on-topic?

**Fluency:** could you understand the responses? Did the language seem accurate?

### 5.1 Results

**Pretrained models baseline** Pretrained conversation models are rated poorly by humans for empathy when the candidates are retrieved from Reddit utterances or when a generative model is used (Table 2). Higher ratings with models based on BERT or larger Transformer models show that increasing the capacity makes the models seem more empathetic, but still remain far from human performance, while being considerably more onerous

| | Model | Candidate | Empathy | Relevance | Fluency |
|---|---|---|---|---|---|
| | *Pre-trained* | R | $2.82 \pm 0.12$ | $3.03 \pm 0.13$ | $4.14 \pm 0.10$ |
| | | R+ED | $3.16 \pm 0.14$ | $3.35 \pm 0.13$ | $4.16 \pm 0.11$ |
| | | ED | $3.45 \pm 0.12$ | $3.55 \pm 0.13$ | $4.47 \pm 0.08$ |
| Retrieval | Fine-tuned | ED | $\mathbf{3.76 \pm 0.11}$ | $3.76 \pm 0.12$ | $4.37 \pm 0.09$ |
| | EmoPrepend-1 | ED | $3.44 \pm 0.11$ | $3.70 \pm 0.11$ | $4.40 \pm 0.08$ |
| | TopicPrepend-1 | ED | $3.72 \pm 0.12$ | $\mathbf{3.91 \pm 0.11}$ | $\mathbf{4.57 \pm 0.07}$ |
| | *Pre-trained* | R | $3.06 \pm 0.13$ | $3.29 \pm 0.13$ | $4.20 \pm 0.10$ |
| | | R+ED | $3.49 \pm 0.12$ | $3.62 \pm 0.12$ | $4.41 \pm 0.09$ |
| | | ED | $3.43 \pm 0.13$ | $3.49 \pm 0.14$ | $4.37 \pm 0.10$ |
| Retrieval w/ BERT | Fine-tuned | ED | $3.71 \pm 0.12$ | $3.76 \pm 0.12$ | $4.58 \pm 0.06$ |
| | EmoPrepend-1 | ED | $3.93 \pm 0.12$ | $3.96 \pm 0.13$ | $4.54 \pm 0.09$ |
| | TopicPrepend-1 | ED | $\mathbf{4.03 \pm 0.10}$ | $\mathbf{3.98 \pm 0.11}$ | $\mathbf{4.65 \pm 0.07}$ |
| | *Pre-trained* | – | $2.31 \pm 0.12$ | $2.21 \pm 0.11$ | $3.89 \pm 0.12$ |
| Generative | Fine-Tuned | – | $\mathbf{3.25 \pm 0.12}$ | $\mathbf{3.33 \pm 0.12}$ | $4.30 \pm 0.09$ |
| | EmoPrepend-1 | – | $3.16 \pm 0.12$ | $3.19 \pm 0.13$ | $4.36 \pm 0.09$ |
| | TopicPrepend-1 | – | $3.09 \pm 0.13$ | $3.12 \pm 0.13$ | $\mathbf{4.41 \pm 0.08}$ |
| *Gold Response* | – | – | $4.19 \pm 0.10$ | $4.55 \pm 0.07$ | $4.68 \pm 0.06$ |

Table 2: Human ratings. Fine-tuning on ED and using ED candidates generally improves scores, especially on Empathy, with minimal retraining. Additional external supervision (Prepend) improves the Empathy and Relevance scores for BERT-based models. Bold: best score for that group. Italics: reference model for the group.

to train (Table 3).[10]

**Using EMPATHETICDIALOGUES for candidate selection** Table 1 shows that merely using the pool of candidates from the training set of ED improves the BLEU scores of retrieval models.

Using candidates from our dataset also substantially improves the performance of pre-trained retrieval models on all human metrics, particularly the Empathy subscore of most interest to us (Table 2).

**Using EMPATHETICDIALOGUES for fine-tuning** Additionally, fine-tuning to predict conversation responses on our data improves all automated metrics (Table 1). While fine-tuning on ED data improves performance on predicting the next ED utterance, this may come at the expense of performance when predicting next utterance in other corpora. To measure this, we compared automated metrics on next utterance prediction with pre-trained models and models fine-tuned using ED data (for our base and larger retrieval-based Transformer models) when predicting on DAILYDIALOG and REDDIT (drawing both context and candidates from the

same corpus). Compared to the 12-14% P@1,100 increase measured with ED (see Tables 1 and 7), fine-tuning on ED leads to a 5-7% increase on DD, and a 2-3% decrease on R.[11] For all three datasets, fine-tuning increases AVG BLEU by 0.2 to 0.5. The slight decrease of performance on R is not surprising because the pre-trained model was trained directly on Reddit predictions. But, the improvement on DD is an encouraging sign that improvements from fine-tuning on ED may generalize to other conversation datasets.

Fine-tuning on the ED data also generally improves human metrics on the ED task, in both retrieval and generative set-ups (Table 2).

**Augmenting conversation models with external pretrained classifiers** Automated and human evaluations suggest that prepending emotion or topic predictions may boost perfomance of high-capacity models based on BERT (but not the smaller models), with Empathy ratings close to approaching human performance.

More extensive experiments with large models would be required to confirm that larger capacity makes additional external supervision effective for this task.

---

[10]Results on larger retrieval-based Transformer models in Table 9 of the Appendix show the same pattern.

[11]Numbers for these datasets are included in Table 6 of the appendix.

| | Model | Params, resources, train examples | Emp | Rel | Fluent |
|---|---|---|---|---|---|
| Retrieval | Pretrained-R | 84.3M, 2.5 days, 8GPUs, 1.7B | 2.8 | 3.0 | 4.1 |
| | Pretrained-ED | same , same, same | 3.5 | 3.6 | 4.5 |
| | Fine-Tuned | same , + 0.5 hour, 1 GPU, +22.3k | 3.8 | 3.8 | 4.4 |
| | Pretrained-Bert-R | 217M, 13.5 days, 8GPUs , 1.7B | 3.1 | 3.3 | 4.2 |
| | Pretrained-Bert-ED | same, same, same | 3.4 | 3.5 | 4.4 |
| | Fine-Tuned-Bert | same, +1hour, 8GPUs, +22.3k | 3.7 | 3.8 | 4.6 |
| Generative | Pretrained | 85.1M, 2 days, 32 GPUs, 1.7B | 2.3 | 2.2 | 3.9 |
| | Fine-Tuned | same , +1 hour, 1 GPU, +22.3k | 3.3 | 3.3 | 4.3 |
| | Pretrained-Large | 86.2M, 2.5 days, 32 GPUs, 1.7B | 2.8 | 3.0 | 4.0 |
| | Fine-Tuned-Large | same , +0.5 hour, 1 GPU, +22.3k | 3.6 | 3.6 | 4.5 |

Table 3: Training resources for different models, with human ratings for empathy (Emp), relevance (Rel) and fluency (Fluent). Retrieval-based models use reply candidates from the ED training set (ED) or from Reddit (R). Resource comparisons are relative to the first row of each group. Fine-tuning on ED improves all scores (except for Fluency in one case) while requiring minimal additional training resources. SEM is approximately 0.1

| Context | Setup | Model | Cand | Generated Response |
|---|---|---|---|---|
| Someone came to my door with a gun the other day. | Gen | Pretrained | - | I'm not sure if you're joking or not. |
| | Gen | Fine-Tuned | - | Oh no! What did you do? |
| | Gen | EmoPrepend | - | Oh no! Did you have a weapon? |
| | Gen | TopicPrepend | - | I'm glad to hear that. |
| My son failed his exams! Can you believe it! I don't know what to do! | Ret-BERT | Pretrained | R | Hopefully his mistakes can help you know what not to do on the tests. |
| | Ret-BERT | Pretrained | ED | Oh no! I ' m so sorry! I hope you could pull your grade up! |
| | Ret-BERT | Fine-Tuned | ED | Oh no! Did he not study? |
| | Ret-BERT | EmoPrepend | ED | Im sorry about that! Did your son start school already? |
| | Ret-BERT | TopicPrepend | ED | Oh no! Does he have the option of taking it again? |

Table 4: Examples of model responses. Shading is provided for better legibility.

**Resources and capacity** Table 3 quantifies resource and parameter usage for several models and set-ups, including a larger Transformer generative model (5 layers instead of 4) and BERT-based architectures with substantially more parameters that require longer training. Using ED candidates in pretrained retrieval models, or fine-tuning pretrained conversation models on ED data makes smaller models perform better than larger ones with minimal increase in resource usage.

tion labels. Our experiments show that using this dataset to provide retrieval candidates or fine-tune conversation models leads to responses that are evaluated as more empathetic. how to integrate empathetic responding into more general dialogue when, for example, the needs for empathy have to be balanced with staying on topic or providing information. We hope that our results and dataset will stimulate more research in the important direction of making dialog systems more empathetic.

# 6 Conclusion

We introduce a new dataset of 25k dialogues grounded in situations prompted by specific emo-

# Acknowledgments

## References

Katie Elson Anderson. 2015. Ask me anything: what is reddit? *Library Hi Tech News*, 32(5):8–11.

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. Springer.

Natalya N Bazarova, Yoon Hyung Choi, Victoria Schwanda Sosik, Dan Cosley, and Janis Whitlock. 2015. Social sharing of emotions on facebook: Channel differences, satisfaction, and replies. In *Proceedings of the 18th ACM conference on computer supported cooperative work & social computing*, pages 154–164. ACM.

Timothy Bickmore and Justine Cassell. 2001. Relational agents: a model and implementation of building user trust. In *Proceedings of the SIGCHI conference on Human factors in computing systems*, pages 396–403. ACM.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. BERT: pre-training of deep bidirectional transformers for language understanding. *CoRR*, abs/1810.04805.

Emily Dinan, Varvara Logacheva, Valentin Malykh, Alexander Miller, Kurt Shuster, Jack Urbanek, Douwe Kiela, Arthur Szlam, Iulian Serban, Ryan Lowe, Shrimai Prabhumoye, Alan W Black, Alexander Rudnicky, Jason Williams, Joelle Pineau, Mikhail Burtsev, and Jason Weston. 2019. The second conversational intelligence challenge (convai2). *arXiv preprint arXiv:1902.00098*.

Venkatesh Duppada, Royal Jain, and Sushant Hiray. 2018. Seernet at semeval-2018 task 1: Domain adaptation for affect in tweets. In *Proceedings of The 12th International Workshop on Semantic Evaluation*, pages 18–23.

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

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 the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 1615–1625.

Jamie Fraser, Ioannis Papaioannou, and Oliver Lemon. 2018. Spoken conversational ai in video games: Emotional dialogue management increases user engagement. In *Proceedings of the 18th International Conference on Intelligent Virtual Agents*, pages 179–184. ACM.

Erving Goffman. 1959. The presentation of self in everyday life.

Erving Goffman. 1981. *Forms of talk*. University of Pennsylvania Press.

Edouard Grave, Piotr Bojanowski, Prakhar Gupta, Armand Joulin, and Tomas Mikolov. 2018. Learning word vectors for 157 languages. In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC-2018)*.

Umang Gupta, Ankush Chatterjee, Radhakrishnan Srikanth, and Puneet Agrawal. 2017. A sentiment-and-semantics-based approach for emotion detection in textual conversations. *arXiv preprint arXiv:1707.06996*.

John Heritage. 2005. Conversation analysis and institutional talk. *Handbook of language and social interaction*, 103:47.

Zhiting Hu, Zichao Yang, Xiaodan Liang, Ruslan Salakhutdinov, and Eric P Xing. 2017. Toward controlled generation of text. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 1587–1596. JMLR. org.

Chenyang Huang, Osmar Zaiane, Amine Trabelsi, and Nouha Dziri. 2018. Automatic dialogue generation with expressed emotions. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, volume 2, pages 49–54.

Bernd Huber, Daniel McDuff, Chris Brockett, Michel Galley, and Bill Dolan. 2018. Emotional dialogue generation using image-grounded language models. In *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems*, page 277. ACM.

Thorsten Joachims. 1996. A probabilistic analysis of the rocchio algorithm with tfidf for text categorization. Technical report, Carnegie-mellon univ pittsburgh pa dept of computer science.

Armand Joulin, Edouard Grave, Piotr Bojanowski, and Tomas Mikolov. 2017. 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*, volume 2, pages 427–431.

Sung Soo Kim, Stan Kaplowitz, and Mark V Johnston. 2004. The effects of physician empathy on patient satisfaction and compliance. *Evaluation & the health professions*, 27(3):237–251.

Min Kyung Lee, Sara Kiesler, and Jodi Forlizzi. 2010. Receptionist or information kiosk: how do people talk with a robot? In *Proceedings of the 2010 ACM conference on Computer supported cooperative work*, pages 31–40. ACM.

Wendy Levinson, Rita Gorawara-Bhat, and Jennifer Lamb. 2000. A study of patient clues and physician responses in primary care and surgical settings. *Jama*, 284(8):1021–1027.

Jiwei Li, Michel Galley, Chris Brockett, Jianfeng Gao, and Bill Dolan. 2016a. A diversity-promoting objective function for neural conversation models. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 110–119.

Jiwei Li, Michel Galley, Chris Brockett, Georgios Spithourakis, Jianfeng Gao, and Bill Dolan. 2016b. A persona-based neural conversation model. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, volume 1, pages 994–1003.

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 (Volume 1: Long Papers)*, volume 1, pages 986–995.

Eden Litt, Erin Spottswood, Jeremy Birnholtz, Jeff T Hancock, Madeline E Smith, and Lindsay Reynolds. 2014. Awkward encounters of an other kind: collective self-presentation and face threat on facebook. In *Proceedings of the 17th ACM conference on Computer supported cooperative work & social computing*, pages 449–460. ACM.

Chia-Wei Liu, Ryan Lowe, Iulian Serban, Mike Noseworthy, Laurent Charlin, and Joelle Pineau. 2016. How not to evaluate your dialogue system: An empirical study of unsupervised evaluation metrics for dialogue response generation. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 2122–2132.

Xiaodong Liu, Pengcheng He, Weizhu Chen, and Jianfeng Gao. 2019. Multi-task deep neural networks for natural language understanding. *arXiv preprint arXiv:1901.11504*.

Ryan Lowe, Iulian V Serban, Mike Noseworthy, Laurent Charlin, and Joelle Pineau. 2016. On the evaluation of dialogue systems with next utterance classification. In *17th Annual Meeting of the Special Interest Group on Discourse and Dialogue*, page 264.

Minh-Thang Luong, Quoc V Le, Ilya Sutskever, Oriol Vinyals, and Lukasz Kaiser. 2016. Multi-task sequence to sequence learning. In *Proceedings of ICLR*.

Pierre-Emmanuel Mazare, Samuel Humeau, Martin Raison, and Antoine Bordes. 2018. Training millions of personalized dialogue agents. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2775–2779.

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

Saif Mohammad, Felipe Bravo-Marquez, Mohammad Salameh, and Svetlana Kiritchenko. 2018. Semeval-2018 task 1: Affect in tweets. In *SemEval@NAACL-HLT*.

Saif M. Mohammad. 2012. #emotional tweets. In *Proceedings of the First Joint Conference on Lexical and Computational Semantics - Volume 1: Proceedings of the Main Conference and the Shared Task, and Volume 2: Proceedings of the Sixth International Workshop on Semantic Evaluation*, SemEval '12, pages 246–255, Stroudsburg, PA, USA. Association for Computational Linguistics.

Tong Niu and Mohit Bansal. 2018. Polite dialogue generation without parallel data. *Transactions of the Association for Computational Linguistics*, 6:373–389.

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting on association for computational linguistics*, pages 311–318. Association for Computational Linguistics.

Ji Ho Park, Peng Xu, and Pascale Fung. 2018. Plusemo2vec at semeval-2018 task 1: Exploiting emotion knowledge from emoji and #hashtags. In *Proceedings of The 12th International Workshop on Semantic Evaluation*, pages 264–272.

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

Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners.

Byron Reeves and Clifford Ivar Nass. 1996. *The media equation: How people treat computers, television, and new media like real people and places.* Cambridge university press.

Alan Ritter, Colin Cherry, and Bill Dolan. 2010. Unsupervised modeling of twitter conversations. In *Human Language Technologies: The 2010 Annual Conference of the North American Chapter of the Association for Computational Linguistics*, pages 172–180. Association for Computational Linguistics.

Klaus R. Scherer and Harald G. Wallbott. 1994. Evidence for universality and cultural variation of differential emotion response patterning. *Journal of personality and social psychology*, 66 2:310–28.

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

Carlo Strapparava and Rada Mihalcea. 2007. Semeval-2007 task 14: Affective text. In *SemEval@ACL*.

Ashish Vaswani, 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.

Ashwin K. Vijayakumar, Michael Cogswell, Ramprasaath R. Selvaraju, Qing Sun, Stefan Lee, David J. Crandall, and Dhruv Batra. 2016. Diverse beam search: Decoding diverse solutions from neural sequence models. *CoRR*, abs/1610.02424.

Ke Wang and Xiaojun Wan. 2018. Sentigan: generating sentimental texts via mixture adversarial networks. In *Proceedings of the 27th International Joint Conference on Artificial Intelligence*, pages 4446–4452. AAAI Press.

Tsung-Hsien Wen, Milica Gasic, Nikola Mrkšić, 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.

Kathryn R Wentzel. 1997. Student motivation in middle school: The role of perceived pedagogical caring. *Journal of educational psychology*, 89(3):411.

Peng Xu, Andrea Madotto, Chien-Sheng Wu, Ji Ho Park, and Pascale Fung. 2018. Emo2vec: Learning generalized emotion representation by multitask training. In *Proceedings of the 9th Workshop on Computational Approaches to Subjectivity, Sentiment and Social Media Analysis*, pages 292–298.

Yinfei Yang, Steve Yuan, Daniel Cer, Sheng yi Kong, Noah Constant, Petr Pilar, Heming Ge, Yun-Hsuan Sung, Brian Strope, and Ray Kurzweil. 2018. Learning semantic textual similarity from conversations. In *Rep4NLP@ACL*.

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 (Volume 1: Long Papers)*, volume 1, pages 2204–2213.

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*.

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 (Volume 1: Long Papers)*, volume 1, pages 1128–1137.

## A  Data Examples

We include ten randomly selected dialogues from our training set in Table 5.

## B  Human Evaluation Crowdsourcing Task

Human evaluations were collected on MTurk. For the rating task, each worker was shown a set of 10 randomly subsampled examples from the test set, one after another, each from a different randomly selected model. The worker had to rate the empathy, relevance, and fluency of each example before moving onto the next one. At least 100 ratings were collected per model. 221 US workers participated in the rating task, and each had to perform a minimum of one set of 10 ratings.

## C  Next utterance prediction on other datasets

We test how fine-tuning on ED data affects next utterance prediction on two external datasets (REDDIT and DAILYDIALOG). In this experiment, we use both candidates and context from the DD or R data. Results in Table 6 show that performance on DAILYDIALOG improves after fine-tuning on our data.

## D  Additional Experimental Details and Results

### D.1  Training Details

We used Adamax for training throughout, and dropout was set to 0% everywhere except for a 20% dropout in the linear layer of the emotion-label term of the MULTITASK objective function (discussed below). A learning rate of $8e-4$ was used for all four-layer Transformer models, following Mazare et al. (2018). For the five-layer retrieval-based Transformer model (Pretrained-Large and Fine-Tuned-Large), the learning rate was selected by picking the best performing over the validation set, among values randomly sampled between between $5e-5$ and $8e-4$. When training the retrieval-based BERT model on Reddit and ED data, the learning rate was selected by picking the best performing over the validation set, among values randomly sampled between $6e-6$ and $2e-4$. For training BERT models on Reddit data, we also experimented with adding an additional Transformer layer after the output embedding of the BERT model, but this slightly degraded

P@1,100 scores on the validation set. We used a learning rate of $8e-5$ for the five-layer generative Transformer models.

### D.2  Additional Experimental Set-Ups

We investigated a few additional approaches for incorporating supervised emotion or topic prediction in generating dialogue, but observed little performance improvement. Methods are described below.

**Multitask with Emotion labels**  If the most appropriate response depends on some information for which supervision is available, e.g., the emotions at play, nudging the model to encode this information could result in better performance. We experimented with this by training the base architecture in the one-to-many style of multi-task learning that has been used for NLP seq2seq settings (Luong et al., 2016). In this set-up, MULTITASK, we altered the objective function to also optimize for predicting the emotion label of the conversation to which the utterances being encoded belonged. We added to the context encoder a linear layer and softmax that predicted the emotion label from the context sentences. The objective function was altered to be the average of the negative log-likelihood of predicting the next utterance $\bar{y}$ and the negative log-likelihood of the added linear layer being able to predict the correct emotion.

**Prepend-3, Prepend-5**  We investigated whether Prepend models could be improved by adding the top-3/5 predicted emotion or topic labels by the classifier (rather than top-1).

**Ensemble of Encoders**  We also investigated another approach for incorporating external predictors, which we report the results of in our extended results tables. In this set-up (ENSEM), we augmented the encoders to incorporate latent representations from pretrained supervised architectures. We replaced each of the encoders in our Transformer networks with an Ensemble encoder , similar to a many-to-one style encoder-decoder architecture (Luong et al., 2016). This encoder took the encoding $h_w$ from our basic Transformer encoder (either $h_x$ or $h_y$), already trained on our data, and concatenated it with the representation $h_c$ extracted from the inner layer of a classification network. We used the penultimate layer of a deep emotion classifier. The concatenated encodings were projected linearly to the di-

**Label: Content**
**Situation:** Speaker felt this when...
"eating my favorite meal makes me happy."
**Conversation:**
Speaker: i am at my best when i have my favorite meal.
Listener: nice
Speaker: i love enchiladas
Listener: really?
Speaker: yes. enchiladas for the win!

**Label: Proud**
**Situation:** Speaker felt this when...
"I was proud when my brother finished college. He worked so hard at it"
**Conversation:**
Speaker: I was proud of my brother when he finished school. He worked so hard at it
Listener: Nice, tell him congrats. What did he major in?
Speaker: It was English
Listener: He should become an English teacher1

**Label: Joyful**
**Situation:** Speaker felt this when...
"I have had a great week!"
**Conversation:**
Speaker: I have had a great start to my week!
Listener: That's great. Do you think the rest of the week will be as great?
Speaker: I hope so! It looks promising!!
Listener: Lucky you. Are you always a positive person or it's just been an amazing week really?
Speaker: haha. Kind of both. And also probably too much coffee to start my shift tonight

**Label: Terrified**
**Situation:** Speaker felt this when...
"I got home for lunch and found a bat outside on my front porch."
**Conversation:**
Speaker: I got home for lunch and found a bat outside on my front porch. It probably has rabies. Bats shouldn't be out during the day.
Listener: Doesn't rabies cause sensativity to light? Either way I would freak out...
Speaker: It can but, it also causes anmails to behave erratically... like bats wadering around in the middle of the day.
Listener: Oh yeah, gotcha. I really don't like animals that are small and move quickly
Speaker: Generally yes.

**Label: Anticipating**
**Situation:** Speaker felt this when...
"I cant wait to go on my end of summer trip"
**Conversation:**
Speaker: I cant wait to go on my end of summer trip in texas.
Listener: Sounds like fun. What you got planned ?
Speaker: not really sure but im excited to just be invited
Listener: Got any family out there? Cousins perhaps

**Label: Terrified**
**Situation:** Speaker felt this when...
"My brother jump scared me while I was out playing. It was crazy bad."
**Conversation:**
Speaker: Just got scared to death.
Listener: Oh no. What happened?
Speaker: My brother jumped scared me.
Listener: lol is he younger or older?

**Label: Proud**
**Situation:** Speaker felt this when...
"My little dog learned to sit!"
**Conversation:**
Speaker: I finally tough my new little puppy his first trick!
Listener: What trick did you teach him?
Speaker: I tought him to sit for a treat, its so cute.
Listener: That is good, do you plan to teach him more tricks?

**Label: Apprehensive**
**Situation:** Speaker felt this when...
"I have to call my landlord about being late on the rent. I really don't want to have this conversation."
**Conversation:**
Speaker: I have to make a dreadful phone call tomorrow
Listener: Oh no, about what?
Speaker: I'm late on my rent and I need another week. I don't want to because my landlord isnt very nice
Listener: Oh no, I've been there done that too many times.
Speaker: I don't want her to make a big deal

**Label: Confident**
**Situation:** Speaker felt this when...
"When my husband asked me about how to build a chicken coop I was able to give him a reply that was backed up by blueprints and research from the internet. "
**Conversation:**
Speaker: We recently got 9 chicks and we've been having to work on making them a coop! I had to do so much research but I think we finally have a place that they'll enjoy living when they aren't able to free range.
Listener: OHH! I Love chickens ! I have always wanted some. I have a duck! lol- What kind of chickens are they?
Speaker: We currently have 2 Australorps, 3 Rhode Island Reds, 3 Barred Plymouth Rocks, and 1 Welsummer, but 4 of the 9 ended up being roosters. Ugh!
Listener: Oh man! They fight sometimes. I hope they aren't too bad about waking you up in the morning. Chickens can be very sweet though!
Speaker: I love my little hens, especially one I've named Curly. The roosters might get replaced by hens though because the crowing is so frustrating!

**Label: Surprised**
**Situation:** Speaker felt this when...
"I got a lottery ticket while I was at work today. I won $100 on the scratch off. I was shocked. I never win."
**Conversation:**
Speaker: I won $100 on a scratch off today. I was shocked. I never win.
Listener: Wow! How often do you play the lottery?
Speaker: I usually go on our Tuesday break to buy one with coworkers.
Listener: Neat! Well that is a fantastic feat. Maybe you can win again sometime?

Table 5: 10 random examples from EMPATHETICDIALOGUES training set.

|  | P @1,100 | | BLEU | |
|---|---|---|---|---|
| Model | DD | R | DD | R |
| Pretrained | 39.04 | 58.95 | 6.65 | 1.43 |
| Fine-Tuned | 44.58 | 56.25 | 7.14 | 1.64 |
| Pretrained-Large | 42.28 | 61.60 | 6.94 | 1.42 |
| Fine-Tuned-Large | 48.96 | 58.71 | 7.42 | 1.73 |

Table 6: Performance of the retrieval-based pretrained model and retrieval-based models fine-tuned on ED data for next utterance prediction in other datasets, with both context and candidates from the same dataset (R=Reddit, DD=DailyDialog).

mension required by the decoder, whose architecture didn't change. When training the dialogue model, we froze both the base Transformer encoder and the pretrained classifier and trained only the linear layers (and the decoder for generative systems). We used emotion-related supervision from Emojis from Twitter, through the use of the trained Deepmoji system (Felbo et al., 2017) released by the authors, either as-is (ENSEM-DM) or fine-tuned on the situation descriptions of EMPATHETICDIALOGUES (ENSEM-DM+).

### D.3 Additional Experiments Results

Automated and human evaluations for any additional experiments are in Tables 7 and 8, respectively. All of these model variations show improvements over the pre-trained models. In some metrics, many of these models show slight improvements over the fine-tuned models, as well, though not as consistently, except for the larger BERT retrieval-based models. While prepending top-1 or top-3 labels do not improve generative model scores, the results in Table 8 suggest that multitask, prepend-5, and ensemble set-ups may improve the human evaluations of the fine-tuned generative model for empathy, but are too inconsistent to be conclusive without more corroborating experiments.

### D.4 Emotion Classification Results

Our dataset can also be used to train or fine-tune an emotion classifier, as we do in our PREPEND-K and ENSEM-DM+ set-ups. To give a sense of where the difficulty falls compared to existing emotion and sentiment classification benchmarks, we reproduce the table from Felbo et al. (2017) and add results when fine-tuning the Deepmoji

model on our dataset, or using a fastText classifier (Table 10).

| Model | Candidate Source | Retrieval | | Retrieval w/ BERT | | Generative | |
| | | P@1,100 | AVG BLEU | P@1,100 | AVG BLEU | PPL | AVG BLEU |
|---|---|---|---|---|---|---|---|
| Pretrained | R | - | 4.10 | - | 4.26 | 27.96 | 5.01 |
| | R+ED | - | 4.96 | - | 5.62 | - | - |
| | ED | 43.25 | 5.51 | 49.94 | 5.97 | - | - |
| Fine-Tuned | R | - | 3.85 | - | 4.14 | - | - |
| | R+ED | - | 4.76 | - | 5.39 | - | - |
| | ED | 56.90 | 5.88 | **65.92** | 6.21 | 21.24 | 6.27 |
| | ED+DD | - | 5.61 | - | - | - | - |
| | ED+DD+R | - | 4.74 | - | - | - | - |
| Pretrained-Large | R | - | 4.16 | - | - | - | - |
| | ED | 47.58 | 5.78 | - | - | 23.64 | 6.31 |
| Fine-Tuned-Large | ED | **60.44** | 6.01 | - | - | **16.55** | **8.06** |
| Multitask | ED | 55.73 | 6.18 | 65.90 | 6.17 | 24.07 | 5.42 |
| EmoPrepend-1 | ED | 56.31 | 5.93 | 66.04 | 6.20 | 24.30 | 4.36 |
| EmoPrepend-3 | ED | 55.75 | **6.23** | 65.85 | 6.14 | 23.96 | 2.69 |
| EmoPrepend-5 | ED | 56.35 | 6.18 | 64.69 | 6.21 | 25.40 | 5.56 |
| TopicPrepend-1 | ED | 56.38 | 6.00 | 65.96 | 6.18 | 25.40 | 4.17 |
| TopicPrepend-3 | ED | 55.44 | 5.97 | 65.85 | **6.25** | 25.02 | 3.13 |
| TopicPrepend-5 | ED | 55.75 | 6.17 | 65.65 | 6.19 | 25.10 | 6.20 |
| Ensem-DM | ED | 52.71 | 6.03 | - | - | 19.05 | 6.83 |
| Ensem-DM+ | ED | 52.35 | 6.04 | - | - | 19.10 | 6.77 |

Table 7: Automatic evaluation metrics on the test set for full set of experimental setups. Pretrained: basic Transformer model pretrained on a dump of 1.7 billion REDDIT conversations. Fine-Tuned: model fine-tuned over the EMPATHETICDIALOGUES training data. Multitask: model trained with multitask loss function (predicting the emotion label). EmoPrepend-1/3/5, TopicPrepend-1/3/5: model using top-k labels outputted by an external classifier as prepended tokens. Ensem: model incorporating external classifiers by concatenating representations from deepmoji with the fine-tuned transformer representation. Candidates come from REDDIT (R) or EMPATHETICDIALOGUES (ED). P@1,100: precision retrieving the correct test candidate out of 100 test candidates. AVG BLEU: average of BLEU-1,-2,-3,-4. PPL: perplexity. *Bold: best performance in that column.*

|  | Model | Candidate | Empathy | Relevance | Fluency |
|---|---|---|---|---|---|
| Retrieval | *Pretrained* | R | $2.82 \pm 0.12$ | $3.03 \pm 0.13$ | $4.14 \pm 0.10$ |
|  |  | R+ED | $3.16 \pm 0.14$ | $3.35 \pm 0.13$ | $4.16 \pm 0.11$ |
|  |  | ED | $3.45 \pm 0.12$ | $3.55 \pm 0.13$ | $4.47 \pm 0.08$ |
|  | Fine-Tuned | R | $2.51 \pm 0.12$ | $2.90 \pm 0.12$ | $4.04 \pm 0.11$ |
|  |  | R+ED | $3.06 \pm 0.14$ | $3.34 \pm 0.13$ | $4.12 \pm 0.11$ |
|  |  | ED | $3.76 \pm 0.11$ | $3.76 \pm 0.12$ | $4.37 \pm 0.09$ |
|  | Multitask | ED | $3.63 \pm 0.12$ | $3.83 \pm 0.12$ | $4.49 \pm 0.08$ |
|  | EmoPrepend-1 | ED | $3.44 \pm 0.11$ | $3.70 \pm 0.11$ | $4.40 \pm 0.08$ |
|  | EmoPrepend-3 | ED | $3.54 \pm 0.11$ | $3.76 \pm 0.11$ | $4.54 \pm 0.07$ |
|  | EmoPrepend-5 | ED | $3.42 \pm 0.11$ | $3.61 \pm 0.11$ | $4.53 \pm 0.07$ |
|  | TopicPrepend-1 | ED | $3.72 \pm 0.12$ | $3.91 \pm 0.11$ | $4.57 \pm 0.07$ |
|  | TopicPrepend-3 | ED | $3.64 \pm 0.11$ | $3.66 \pm 0.12$ | $4.51 \pm 0.08$ |
|  | TopicPrepend-5 | ED | $3.34 \pm 0.12$ | $3.52 \pm 0.12$ | $4.24 \pm 0.09$ |
|  | Ensem-DM | ED | $3.61 \pm 0.11$ | $3.71 \pm 0.12$ | $4.45 \pm 0.08$ |
|  | Pretrained-Large | R | $2.94 \pm 0.14$ | $3.12 \pm 0.14$ | $4.23 \pm 0.10$ |
|  |  | ED | $3.47 \pm 0.14$ | $3.56 \pm 0.13$ | $4.41 \pm 0.10$ |
|  | Fine-Tuned-Large | ED | $3.81 \pm 0.12$ | $3.90 \pm 0.12$ | $4.56 \pm 0.08$ |
| Retrieval w/ BERT | *Pretrained* | R | $3.06 \pm 0.13$ | $3.29 \pm 0.13$ | $4.20 \pm 0.10$ |
|  |  | R+ED | $3.49 \pm 0.12$ | $3.62 \pm 0.12$ | $4.41 \pm 0.09$ |
|  |  | ED | $3.43 \pm 0.13$ | $3.49 \pm 0.14$ | $4.37 \pm 0.10$ |
|  | Fine-Tuned | R | $2.90 \pm 0.13$ | $3.39 \pm 0.13$ | $4.36 \pm 0.09$ |
|  |  | R+ED | $3.46 \pm 0.13$ | $3.90 \pm 0.12$ | $4.46 \pm 0.08$ |
|  |  | ED | $3.71 \pm 0.12$ | $3.76 \pm 0.12$ | $4.58 \pm 0.06$ |
|  | Multitask | ED | $3.80 \pm 0.12$ | $3.97 \pm 0.11$ | $4.63 \pm 0.07$ |
|  | EmoPrepend-1 | ED | $3.93 \pm 0.12$ | $3.96 \pm 0.13$ | $4.54 \pm 0.09$ |
|  | EmoPrepend-3 | ED | $3.73 \pm 0.13$ | $3.88 \pm 0.14$ | $4.60 \pm 0.09$ |
|  | EmoPrepend-5 | ED | $4.08 \pm 0.10$ | $4.10 \pm 0.11$ | $4.67 \pm 0.07$ |
|  | TopicPrepend-1 | ED | $4.03 \pm 0.10$ | $3.98 \pm 0.11$ | $4.65 \pm 0.07$ |
|  | TopicPrepend-3 | ED | $3.73 \pm 0.12$ | $3.84 \pm 0.13$ | $4.52 \pm 0.08$ |
|  | TopicPrepend-5 | ED | $3.72 \pm 0.12$ | $3.80 \pm 0.12$ | $4.46 \pm 0.09$ |
| Generative | *Pretrained* | - | $2.31 \pm 0.12$ | $2.21 \pm 0.11$ | $3.89 \pm 0.12$ |
|  | Fine-Tuned | - | $3.25 \pm 0.12$ | $3.33 \pm 0.12$ | $4.30 \pm 0.09$ |
|  | Multitask | - | $3.36 \pm 0.13$ | $3.34 \pm 0.13$ | $4.21 \pm 0.10$ |
|  | EmoPrepend-1 | - | $3.16 \pm 0.12$ | $3.19 \pm 0.13$ | $4.36 \pm 0.09$ |
|  | EmoPrepend-3 | - | $3.09 \pm 0.13$ | $3.02 \pm 0.13$ | $4.39 \pm 0.09$ |
|  | EmoPrepend-5 | - | $3.32 \pm 0.12$ | $3.23 \pm 0.12$ | $4.35 \pm 0.09$ |
|  | TopicPrepend-1 | - | $3.09 \pm 0.13$ | $3.12 \pm 0.13$ | $4.41 \pm 0.08$ |
|  | TopicPrepend-3 | - | $3.09 \pm 0.12$ | $3.34 \pm 0.13$ | $4.53 \pm 0.08$ |
|  | TopicPrepend-5 | - | $3.46 \pm 0.13$ | $3.68 \pm 0.13$ | $4.60 \pm 0.08$ |
|  | Ensem-DM | - | $3.42 \pm 0.12$ | $3.45 \pm 0.12$ | $4.67 \pm 0.06$ |
|  | Pretrained-Large | - | $2.84 \pm 0.13$ | $2.97 \pm 0.12$ | $4.01 \pm 0.11$ |
|  | Fine-Tuned-Large | - | $3.61 \pm 0.13$ | $3.62 \pm 0.13$ | $4.46 \pm 0.10$ |
| *Gold Response* | – | – | $4.19 \pm 0.10$ | $4.55 \pm 0.07$ | $4.68 \pm 0.06$ |

Table 8: Human evaluation metrics from rating task for additional experiments.

| | Model | Params, resources, train examples | Emp | Rel | Fluent |
|---|---|---|---|---|---|
| Retrieval | Pretrained-R | 84.3M, 2.5 days, 8 GPUs, 1.7B | 2.8 | 3.0 | 4.1 |
| | Pretrained-ED | same , same, +22.3k | 3.5 | 3.6 | 4.5 |
| | Fine-Tuned | same , + 0.5 hour, 1 GPU, +22.3k | 3.8 | 3.8 | 4.4 |
| | Multitask | +9.6k, + 0.5 hour, 1 GPU, +22.3k | 3.6 | 3.6 | 4.5 |
| | Pretrained-Large-R | 86.5M, 10.5 days, 8 GPUs , 1.7B | 2.9 | 3.1 | 4.2 |
| | Pretrained-Large-ED | same, same, +22.3k | 3.5 | 3.6 | 4.4 |
| | Fine-Tuned-Large | same, +1 hour, 1GPU, +22.3k | 3.8 | 3.9 | 4.6 |
| | Pretrained-BERT-R | 217M, 13.5 days, 8 GPUs , 1.7B | 3.1 | 3.3 | 4.2 |
| | Pretrained-BERT-ED | same, same, +22.3k | 3.4 | 3.5 | 4.4 |
| | Fine-Tuned-BERT | same, +1 hour, 8 GPUs, +22.3k | 3.7 | 3.8 | 4.6 |
| | Multitask-BERT | +9.6k, +0.5 hour, 8 GPUs, +22.3k | 3.8 | 4.0 | 4.6 |
| Generative | Pretrained | 85.1M, 2 days, 32 GPUs, 1.7B | 2.3 | 2.2 | 3.9 |
| | Fine-Tuned | same , +1 hour, 1 GPU, +22.3k | 3.3 | 3.3 | 4.3 |
| | Multitask | +9.6k, +1 hour, 1 GPU, +22.3k | 3.2 | 3.2 | 4.3 |
| | Pretrained-Large | 86.2M, 2.5 days, 32 GPUs, 1.7B | 2.8 | 3.0 | 4.0 |
| | Fine-Tuned-Large | same , +0.5 hour, 1 GPU, +22.3k | 3.6 | 3.6 | 4.5 |

Table 9: Training resources for different models, with human ratings for empathy (Emp), relevance (Rel) and fluency (Fluent) for full set of experiments. Retrieval-based models use reply candidates from the ED training set (ED) or from Reddit (R). Resource comparisons are relative to the first row of each group. Fine-tuning on ED improves all scores (except for Fluency in one case) while requiring minimal additional training resources. SEM is approximately 0.1

| Dataset | Metric | SOTA (in 2017) | fastText | DeepMoji new | DeepMoji full | DeepMoji last | DeepMoji chain-thaw |
|---|---|---|---|---|---|---|---|
| SE0714 | F1 | 0.34 | 0.16 | 0.21 | 0.31 | 0.36 | 0.37 |
| OLYMPIC | F1 | 0.50 | 0.38 | 0.43 | 0.50 | 0.61 | 0.61 |
| PSYCHEXP | F1 | 0.45 | 0.44 | 0.32 | 0.42 | 0.56 | 0.57 |
| SS-TWITTER | Acc | 0.82 | 0.68 | 0.62 | 0.85 | 0.87 | 0.88 |
| SS-YOUTUBE | Acc | 0.86 | 0.75 | 0.75 | 0.88 | 0.92 | 0.93 |
| SE0614 | Acc | 0.51 | - | 0.51 | 0.54 | 0.58 | 0.58 |
| SCv1 | F1 | 0.63 | 0.60 | 0.67 | 0.65 | 0.68 | 0.69 |
| SCv2-GEN | F1 | 0.72 | 0.69 | 0.71 | 0.71 | 0.74 | 0.75 |
| ED | Acc | - | 0.43 | 0.40 | 0.46 | 0.46 | 0.48 |
| ED-CUT | Acc | - | 0.41 | 0.36 | 0.42 | 0.44 | 0.45 |

Table 10: Classification performance on EMPATHETICDIALOGUES, with the benchmarks proposed in (Felbo et al., 2017) for reference. ED: performance on predicting the emotion label from the situation description. ED-CUT: same, but after having removed all the situation descriptions where the target label was present.