

# Learning from Dialogue after Deployment: Feed Yourself, Chatbot!

**Braden Hancock\***  
Computer Science Dept.  
Stanford University  
bradenjh@cs.stanford.edu

**Antoine Bordes, Pierre-Emmanuel Mazaré  
Jason Weston**  
Facebook AI Research  
{abordes, pem, jase}@fb.com

## Abstract

The majority of conversations a dialogue agent sees over its lifetime occur after it has already been trained and deployed, leaving a vast store of potential training signal untapped. In this work, we propose the self-feeding chatbot, a dialogue agent with the ability to extract new training examples from the conversations it participates in. As our agent engages in conversation, it also estimates user satisfaction in its responses. When the conversation appears to be going well, the user’s responses become new training examples to imitate. When the agent believes it has made a mistake, it asks for feedback; learning to predict the feedback that will be given improves the chatbot’s dialogue abilities further. On the PERSONACHAT chit-chat dataset with over 131k training examples, we find that learning from dialogue with a self-feeding chatbot significantly improves performance, regardless of the amount of traditional supervision.

## 1 Introduction

Training a dialogue agent to converse like a human requires extensive supervision. The most common approach is to train models to imitate humans in large corpora of crowdsourced or scraped conversations (Serban et al., 2015). These fully-supervised conversations tend to be expensive to collect in sufficient quantity and/or occur in settings with significant differences from the deployment environment (Ross et al., 2009). Instead, dialogue agents would ideally learn directly from dialogue, the conversations they participate in after deployment, which are usually abundant, task-specific, dynamic, and cheap. This corresponds to the way humans learn to converse—not merely observing others engaging in “expert-level” conver-



Figure 1: As the self-feeding chatbot engages in dialogue, it estimates user satisfaction to know when to ask for feedback. From the satisfied responses and feedback responses, new training examples are extracted for the DIALOGUE and FEEDBACK tasks, respectively, both of which improve the model’s dialogue abilities further.

sations, but instead actively adjusting and correcting our speech based on feedback woven throughout our own conversations (Bassiri, 2011; Werts et al., 1995). Giving a dialogue agent this ability would enable it to continuously improve and adapt over its lifetime, rather than requiring additional annotation costs for each and every improvement.

However, naively training a dialogue agent on its own conversations yields poor results. For example, training a model on its own output can simply reinforce its existing failure modes, and mistakes by the agent can lead to absurd conversations that no longer resemble the target domain (Hashimoto and Sassano, 2018). To combat this, one approach is to allow the agent to request feed-

\*BH completed most of this work at Facebook (FAIR).

back during conversations (Zhang et al., 2018a; Li et al., 2017b), e.g., when it believes it is about to make a mistake. This approach, however, falls victim to the Dunning-Kruger effect (Kruger and Dunning, 1999), which in this case suggests that a bad model will also be bad at knowing when it is doing a bad job. Regardless of when feedback is requested, existing methods typically require accompanying scalar rewards or adherence to particular templates or structure to ensure that the feedback is usable by the model (Rieser and Lemon, 2011; Zhang et al., 2017; Liu et al., 2018). These requirements may be acceptable for paid annotators, but they impose unnatural workflows on unpaid conversation partners in a standard dialogue environment. Humans are able to request and provide feedback using only natural language; ideally, dialogue agents would be able to do the same.

In this work we propose the self-feeding chatbot, a dialogue agent with the ability to extract new examples from the conversations it participates in after deployment (Figure 1). Concretely, in addition to being trained on the primary DIALOGUE task, the agent is trained to predict its speaking partner’s satisfaction with its responses. When the conversation seems to be going well, the user’s responses (but not the bot’s own utterances) become the targets in new training examples for the DIALOGUE task. When the agent believes it has made a mistake, it instead requests feedback on what it could have said instead. Predicting the feedback that will be provided in a given context becomes an auxiliary task (FEEDBACK) on which the model is also trained. Importantly, these new examples improve the agent’s dialogue abilities while using only natural responses from the user that do not require special structure, accompanying numerical feedback, or additional human intervention in order to be used.

With this approach, the conversations the chatbot participates in are sliced into two complementary datasets—one largely protected from the chatbot’s mistakes (DIALOGUE examples), and one which directly addresses them (FEEDBACK examples). We validate our approach on the PERSONACHAT (Zhang et al., 2018b) dialogue dataset, finding empirically that regardless of the number of available supervised examples, the dialogue ability of the chatbot is always improved by adding the automatically extracted examples of either type, and improves the most by adding both.

The main contributions of this work thus include the following:

- We propose the self-feeding chatbot, a dialogue agent with the ability to extract new training examples for itself from the conversations it participates in during deployment.
- We show that dialogue ability improves by imitating human responses when the human is satisfied, or by asking for feedback when they are not, predicting it as an auxiliary task.
- We demonstrate that classifying user satisfaction is a learnable task important for the self-feeding process, significantly outperforming an approach based on model uncertainty.
- We release three new datasets to further research in this direction: (1) deployment chat logs (513k messages); (2) ratings of user satisfaction (42k); (3) textual feedback on what a bot could have said in a given context (62k).

The datasets and models described in this paper are available via the ParlAI platform (Miller et al., 2017), along with training code. Hyperparameter values are included in Appendix G.

## 2 Related Work

The general concepts of lifelong learning (Silver et al., 2013) and never-ending (language) learning (Carlson et al., 2010) are related to the topics discussed in this work, as is active learning (Tong and Koller, 2001) and predictive modeling (Schmidhuber and Huber, 1991).

The specific case of learning actively from dialogue during deployment was explored for the question answering (QA) setting in (Weston, 2016) and (Li et al., 2017a), where the authors examined multiple learning strategies on a suite of dialogue tasks with varying types of feedback, such as verbal cues (e.g., “Yes, that’s right!”) and scalar rewards. Most relevant to our work was their use of *forward prediction*, where the learner improved in quality by trying to predict the teacher’s responses without an explicit reward signal. Our work extends this idea, adding the ability for the model to recognize its mistakes and request feedback explicitly, and moving beyond QA to the more general chit-chat setting where there may be many valid responses in a given context.

Learning to ask questions is another area that has been studied (Strub et al., 2017; Wang et al.,

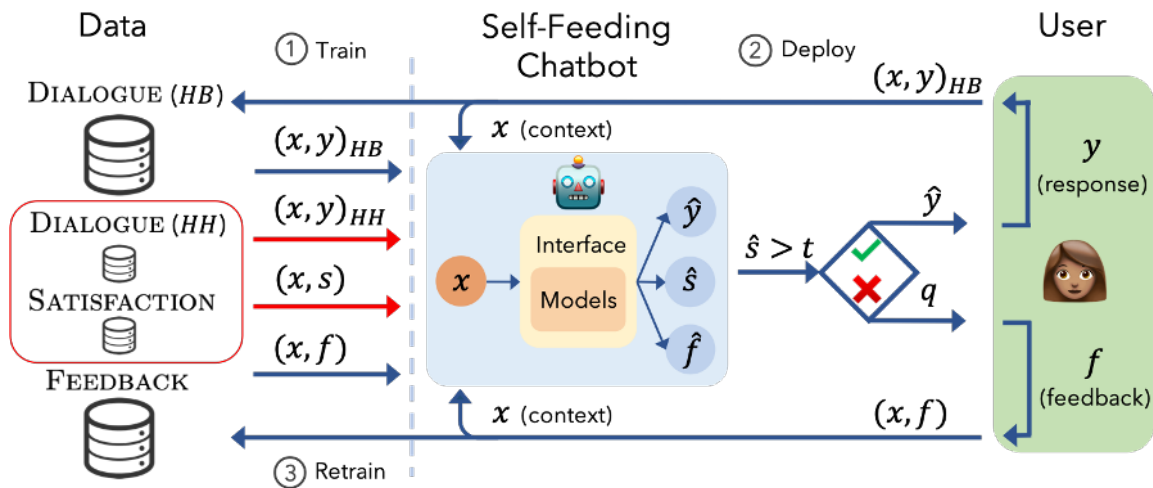


Figure 2: (1) The chatbot is first trained with any available supervised data (boxed in red) on the Human-Human (HH) DIALOGUE  $(x, y)_{HH}$  and SATISFACTION  $(x, s)$  tasks. (2) During deployment, whenever the predicted satisfaction score of the current conversation  $x$  is above the threshold ( $\hat{s} > t$ ), a new Human-Bot (HB) DIALOGUE example  $(x, y)_{HB}$  is extracted and the bot continues the conversation with its own response  $\hat{y}$ . Otherwise, the chatbot requests feedback with question  $q$  and extracts a new FEEDBACK example  $(x, f)$ . (3) The chatbot is periodically retrained with the available examples from all four datasets, improving its DIALOGUE performance without collecting any new supervised examples.

2018; Rao and Daumé, 2018). While those works focused on identifying *which* question to ask in a given context, in this work we are more interested in first learning *when* to ask a question. Li et al. (2017b) considered this question as well, but again in the context of a QA setting rather than dialogue.

Hashimoto and Sassano (2018) used user responses to detect mistakes made by a deployed virtual assistant, showing that model mistakes can be identified in chit-chat, weather, or web search domains. However, they did not explore how to use these identified mistakes to improve the model further; their agent was not equipped to feed itself. Eskenazi et al. (2018) also found that the correctly assessing the appropriateness of chatbot responses is highly dependent on user responses and not preceding context alone.

There are other, somewhat less related, ways to use feedback during dialogue for learning, notably for collecting knowledge to answer questions (Mazumder et al., 2018; Hixon et al., 2015; Pappu and Rudnicky, 2013), and more commonly in reinforcement learning settings, where the feedback is a scalar rather than the dialogue messages themselves (Levin et al., 2000; Schatzmann et al., 2006; Rieser and Lemon, 2011; Liu et al., 2018; Hong et al., 2019). In particular (Serban et al., 2017) employ user sentiment detection for reward shaping in their Alexa prize entry.

Finally, our work improves dialogue quality by utilizing larger datasets with noisier labels than traditional supervision. Other applications of weak supervision to dialogue (Mallinar et al., 2019) and relation extraction have observed similar results (Bunescu and Mooney, 2007; Hancock et al., 2018; Ratner et al., 2017).

### 3 The Self-Feeding Chatbot

The lifecycle of a self-feeding chatbot is outlined in Figure 2. In the initial training phase, the dialogue agent is trained on two tasks—DIALOGUE (next utterance prediction, or what should I say next?) and SATISFACTION (how satisfied is my speaking partner with my responses?)—using whatever supervised training data is available. We refer to these initial DIALOGUE examples as Human-Human (HH) examples, since they were generated in conversations between two humans.

In the deployment phase, the agent engages in multi-turn conversations with users, extracting new *deployment* examples of two types. Each turn, the agent observes the context  $x$  (i.e., the conversation history) and uses it to predict its next utterance  $\hat{y}$  and its partner’s satisfaction  $\hat{s}$ . If the satisfaction score is above a specified threshold  $t$ , the agent extracts a new Human-Bot (HB) DIALOGUE example using the previous context  $x$  and the human’s response  $y$  and continues the conversation.

If, however, the user seems unsatisfied with its previous response ( $\hat{s} < t$ ), the agent requests feedback with a question  $q$ , and the resulting feedback response  $f$  is used to create a new example for the FEEDBACK task (what feedback am I about to receive?). The agent acknowledges receipt of the feedback and the conversation continues. The rate at which new DIALOGUE or FEEDBACK examples are collected can be adjusted by raising or lowering the satisfaction threshold  $t$  (we use  $t = 0.5$ ).<sup>1</sup> Periodically, the agent is retrained using all available data, thereby improving performance on the primary DIALOGUE task.

It is important to note that the user’s responses are always in the form of *natural dialogue*. In particular, at no point are the new FEEDBACK examples inspected, post-processed, or cleaned. Instead, we rely on the fact that the feedback is not random: regardless of whether it is a verbatim response, a description of a response, or a list of possible responses (see Table 2 for examples), there is a learnable relationship between conversation contexts and their corresponding feedback which requires many of the same language understanding skills to master as does carrying on a normal conversation.

The experiments in this paper are limited to the setting where the number of supervised and deployment examples are on the same order of magnitude; however, we envision scenarios in which the number of deployment examples can easily grow to  $100\times$  or more the number of supervised examples over the chatbot’s deployment lifetime, effectively providing a massive task-specific corpus at minimal cost. Table 1 reports the sizes of each dataset, all of which are available via ParlAI.

### 3.1 Task 1: DIALOGUE

The chatbot’s primary task (DIALOGUE) is to carry on a coherent and engaging conversation with a speaking partner. Training examples take the form of  $(x, y)$  pairs, where  $x$  is the context of the conversation (the concatenation of all responses so far up to some history length, delimited with tokens marking the speaker), and  $y$  is the appropriate response given by the human.

The Human-Human (HH) portion of the DIALOGUE dataset comes from the PERSONACHAT dataset (Zhang et al., 2018b), which consists of

<sup>1</sup>Another option would be to have two thresholds—one for each example type—to decouple collection their rates.

Task	Train	Valid	Test	Total
DIALOGUE				
– HH (HUMAN-HUMAN)	131438	7801	6634	145873
– HB (HUMAN-BOT)	60000	0	0	60000
FEEDBACK	60000	1000	1000	62000
SATISFACTION	1000	500	1000	2500

Table 1: The number of examples used in our experiments by task and split. Note that the HH DIALOGUE examples come from the PERSONACHAT dataset, HB DIALOGUE and FEEDBACK examples were collected during deployment, and an additional 40k SATISFACTION training examples were collected for the analysis in Section 5.1.

short dialogues (6-8 turns) between two crowdworkers (humans) who have been assigned short text profiles and are instructed to “chat with the other person naturally and try to get to know each other.” We chose this dataset because of its size (over 145k total examples), the breadth of topics it covers, and its focus on promoting engaging conversations, which we anticipate being a necessary property of a chatbot that people will be willing to chat with voluntarily and repeatedly. We use the standard splits of the dataset made available in ParlAI as a part of the ConvAI2 challenge (Burtsev et al., 2018). Since the question of how to incorporate external knowledge (such as profiles) in dialogue is an open research question of its own (Li et al., 2016; Luan et al., 2017; Luo et al., 2018) and we are primarily interested in the question of learning from dialogue, we discard the profiles and simply train and test on the conversations themselves, making the dataset more challenging in terms of raw performance scores.

The Human-Bot (HB) portion of the DIALOGUE dataset is extracted during deployment as described earlier, where the user is again a crowdworker instructed to chat naturally. The context may contain responses from both the human and the bot, but the target response is always from the human, as we will see experimentally that targeting bot responses degrades performance. Because the chit-chat domain is symmetric, both the HH and HB DIALOGUE examples are used for the same task. In an asymmetric setting where the bot has a different role than the human, it is unclear whether HB examples may still be used as an auxiliary task, but FEEDBACK examples will remain usable.

Category	%	Feedback Examples
Verbatim	53.0	<ul style="list-style-type: none"> <li>• my favorite food is pizza</li> <li>• no, i have never been to kansas</li> <li>• i like when its bright and sunny outside</li> </ul>
Suggestion	24.5	<ul style="list-style-type: none"> <li>• you could say hey, i'm 30. how old are you?</li> <li>• yes, i play battlefield would have a been a great answer.</li> <li>• you could have said "yes, I'm happy it's friday."</li> </ul>
Instructions	14.5	<ul style="list-style-type: none"> <li>• tell me what your favorite breakfast food is</li> <li>• answer the question about having children!</li> <li>• tell me why your mom is baking bread</li> </ul>
Options	8.0	<ul style="list-style-type: none"> <li>• you could have said yes it really helps the environment or no its too costly</li> <li>• you could have said yes or no, or talked more about your mustang dream.</li> <li>• you should have said new york, texas or maryland. something like one of those.</li> </ul>

Table 2: Examples of the types of feedback given to the dialogue agent, pulled from a random sample of 200 feedback responses. Verbatim responses could be used directly in conversation, Suggestion responses contain a potential verbatim response in them somewhere, Instructions describe a response or tell the bot what to do, and Options make multiple suggestions.

### 3.2 Task 2: SATISFACTION

The objective of the SATISFACTION auxiliary task is to predict whether or not a speaking partner is satisfied with the quality of the current conversation. Examples take the form of  $(x, s)$  pairs, where  $x$  is the same context as in the DIALOGUE task, and  $s \in [0, 1]$ , ranging from dissatisfied to satisfied. Crucially, it is hard to estimate from the bot's utterance itself whether the user will be satisfied, but much easier using the human's response to the utterance, as they may explicitly say something to that effect, e.g. "What are you talking about?".

The dataset for this task was collected via crowdsourcing. Workers chatted with our baseline dialogue agent and assigned a rating 1-5 for the quality of each of the agent's responses.<sup>2</sup> Contexts with rating 1 were mapped to the negative class (dissatisfied) and ratings [3, 4, 5] mapped to the positive class (satisfied). Contexts with rating 2 were discarded to increase the separation between classes for a cleaner training set. Note that these numeric ratings were requested only when collecting the initial training data, not during deployment, where only natural dialogue is used.

### 3.3 Task 3: FEEDBACK

The objective of the FEEDBACK auxiliary task is to predict the feedback that will be given by the speaking partner when the agent believes it has made a mistake and asks for help. Examples take the form of  $(x, f)$  pairs, where  $x$  is the same context as the other two tasks and  $f$  is the feedback utterance.

<sup>2</sup>A snapshot of the data collection interface and sample conversations are included in the Appendix.

Training data for this task is collected during deployment. Whenever the user's estimated satisfaction is below a specified threshold, the chatbot responds "Oops! Sorry. What should I have said instead?".<sup>3</sup> A new example for the FEEDBACK task is then extracted using the context up to but not including the turn where the agent made the poor response as  $x$  and the user's response as  $f$  (as shown in Figure 1). At that point to continue the conversation during deployment, the bot's history is reset, and the bot instructs the user to continue, asking for a new topic. Examples of FEEDBACK responses are shown in Table 2.

## 4 Model and Settings

### 4.1 Model Architecture

The self-feeding chatbot has two primary components: an interface component and a model component. The interface component is shared by all tasks, and includes input/output processing (tokenization, vectorization, etc.), conversation history storage, candidate preparation, and control flow (e.g., when to ask a question vs. when to give a normal dialogue response). The model component contains a neural network for each task, with embeddings, a network body, and a task head, some of which can be shared. In our case, we obtained maximum performance by sharing all parameters between the FEEDBACK and DIALOGUE tasks (prepending FEEDBACK responses with a special token), and using separate model parameters for the SATISFACTION task. Identifying optimal task structure in multi-task learning (MTL)

<sup>3</sup>Future work should examine how to ask different kinds of questions, depending on the context.

architectures is an open research problem (Ruder, 2017). Regardless of what parameters are shared, each training batch contains examples from only one task at a time, candidate sets remain separate, and each task’s cross-entropy loss is multiplied by a task-specific scaling factor tuned on the validation set to help account for discrepancies in dataset size, loss magnitude, dataset relevance, etc.

Our dialogue agent’s models are built on the Transformer architecture (Vaswani et al., 2017), which has been shown to perform well on a variety of NLP tasks (Devlin et al., 2018; Radford et al., 2018), including multiple persona-based chat applications (Shuster et al., 2018a,b; Rashkin et al., 2018). For the SATISFACTION task, the context  $x$  is encoded with a Transformer and converted to the scalar satisfaction prediction  $\hat{s}$  by a final linear layer in the task head. The DIALOGUE and FEEDBACK tasks are set up as ranking problems, as in (Zhang et al., 2018b; Mazaré et al., 2018), where the model ranks a collection of candidate responses and returns the top-ranked one as its response. The context  $x$  is encoded with one Transformer and  $\hat{y}$  and  $\hat{f}$  candidates are encoded with another. The score for each candidate is calculated as the dot product of the encoded context and encoded candidate.

During training, negative candidates are pulled from the correct responses for the other examples in the mini-batch. During evaluation, however, to remain independent of batch size and data shuffling, each example is assigned a static set of 19 other candidates sampled at random from its split of the data. During deployment, all 127,712 unique HH DIALOGUE candidates from the train split are encoded once with the trained model and each turn the model selects the top-ranked one for the given context.

## 4.2 Model Settings

Contexts and candidates are tokenized using the default whitespace and punctuation tokenizer in ParlAI. We use a maximum dialogue history length of 2 (i.e., when making a prediction, the dialogue agent has access to its previous utterance and its partner’s response). Tokens are embedded with fastText (Bojanowski et al., 2017) 300-dimensional embeddings. We do not limit the vocabulary size, which varies from 11.5k to 23.5k words in our experiments, depending on the training set. The Transformer is implemented in

PyTorch (Paszke et al., 2017) within the ParlAI framework. We use the AdaMax (Kingma and Ba, 2014) optimizer with a learning rate schedule that decays based on the inverse square root of the step number after 500 steps of warmup from  $1e-5$ . We use proportional sampling (Sanh et al., 2018) to select batches from each task for training, with batch size 128. Each Transformer layer has two attention heads and FFN size 32. The initial learning rate (0.001-0.005), number of Transformer layers (1-2), and task-specific loss factors (0.5-2.0) are selected on a per-experiment basis based on a grid search over the validation set averaged over three runs (we use the DIALOGUE validation set whenever multiple tasks are involved). We use early stopping based on the validation set to decide when to stop training. The hyperparameter values for the experiments in Section 5 are included in Appendix G.

Note that throughout development, a portion of the DIALOGUE validation split was used as an informal test set. The official hidden test set for the DIALOGUE task was used only to produce the final numbers included in this paper.

## 5 Experimental Results

Throughout this section, we use the ranking metric hits@X/Y, or the fraction of the time that the correct candidate response was ranked in the top X out of Y available candidates; accuracy is another name for hits@1/Y. Statistical significance for improvement over baselines is assessed with a two-sample one-tailed T-test.

### 5.1 Benefiting from Deployment Examples

Our main result, reported in Table 3, is that utilizing the deployment examples improves accuracy on the DIALOGUE task regardless of the number of available supervised (HH) DIALOGUE examples.<sup>4</sup> The boost in quality is naturally most pronounced when the HH DIALOGUE training set is small (i.e., where the learning curve is steepest), yielding an increase of up to 9.4 accuracy points, a 31% improvement. However, even when the entire PERSONACHAT dataset of 131k examples is used—a much larger dataset than what is available for most dialogue tasks—adding deployment examples is still able to provide an additional 1.6 points of accuracy on what is otherwise a very flat region of

<sup>4</sup>For comparisons with other models, see Appendix C. The best existing score reported elsewhere on the PERSONACHAT test set without using profiles is 34.9.

Human-Bot (HB)		Human-Human (HH) DIALOGUE			
DIALOGUE	FEEDBACK	20k	40k	60k	131k
-	-	30.3 (0.6)	36.2 (0.4)	39.1 (0.5)	44.7 (0.4)
20k	-	32.7 (0.5)	37.5 (0.6)	40.2 (0.5)	45.5 (0.7)
40k	-	34.5 (0.5)	37.8 (0.6)	40.6 (0.6)	45.1 (0.6)
60k	-	35.4 (0.4)	37.9 (0.7)	40.2 (0.8)	45.0 (0.7)
-	20k	35.0 (0.5)	38.9 (0.3)	41.1 (0.5)	45.4 (0.8)
-	40k	36.7 (0.7)	39.4 (0.5)	41.8 (0.4)	45.7 (0.6)
-	60k	37.8 (0.6)	40.6 (0.5)	42.2 (0.7)	45.8 (0.7)
60k	60k	<b>39.7</b> (0.6)	<b>42.0</b> (0.6)	<b>43.3</b> (0.7)	<b>46.3</b> (0.8)

Table 3: Accuracy (hits@1/20) on the DIALOGUE task’s hidden test set by number of Human-Human (HH) DIALOGUE, Human-Bot (HB) DIALOGUE, and FEEDBACK examples, averaged over 20 runs, with standard deviations in parentheses. For each column, the model using all three data types (last row) is significantly better than all the others, and the best model using only one type of self-feeding (FEEDBACK examples or HB DIALOGUE examples) is better than the supervised baseline in the first row ( $p < 0.05$ ).

the learning curve. It is interesting to note that the two types of deployment examples appear to provide complementary signal, with models performing best when they use both example types, despite them coming from the same conversations. We also calculated hit rates with 10,000 candidates (instead of 20), a setup more similar to the interactive setting where there may be many candidates that could be valid responses. In that setting, models trained with the deployment examples continue to outperform their HH-only counterparts by significant margins (see Appendix B).

On average, we found that adding 20k FEEDBACK examples benefited the agent about as much as 60k HB DIALOGUE examples.<sup>5</sup> This is somewhat surprising given the fact that nearly half of the FEEDBACK responses would not even be reasonable responses if used verbatim in a conversation (instead being a list of options, a description of a response, etc.) as shown in Table 2. Nevertheless, the tasks are related enough that the DIALOGUE task benefits from the MTL model’s improved skill on the FEEDBACK task. And whereas HB DIALOGUE examples are based on conversations where the user appears to already be satisfied with the agent’s responses, each FEEDBACK example corresponds to a mistake made by the model, giving the latter dataset a more active

role in improving quality. Interestingly, our best-performing model, which achieves 46.3 accuracy on DIALOGUE, scores 68.4 on FEEDBACK, suggesting that the auxiliary task is a simpler task overall.

When extracting HB DIALOGUE examples, we ignore human responses that the agent classifies as expressing dissatisfaction, since these turns do not represent typical conversation flow. Including these responses in the 60k HB dataset decreases hits@1/20 by 1.2 points and 0.6 points when added to 20k and 131k HH DIALOGUE examples, respectively. We also explored using chatbot responses with favorable satisfaction scores ( $\hat{s} > t$ ) as new training examples, but found that our models performed better without them (see Appendix D for details).

We also found that “fresher” feedback results in bigger gains. We compared two models trained on 20k HH DIALOGUE examples and 40k FEEDBACK examples—the first collected all 40k FEEDBACK examples at once, whereas the second was retrained with its first 20k FEEDBACK examples before collecting the remaining 20k. While the absolute improvement of the second model over the first was small (0.4 points), it was statistically significant ( $p = 0.027$ ) and reduced the gap to a model trained on fully supervised (HH) DIALOGUE examples by 17% while modifying only 33% of the training data.<sup>6</sup> This improvement makes sense intuitively, since new FEEDBACK examples are

<sup>5</sup>Our baseline chatbot collected approximately one FEEDBACK example for every two HB DIALOGUE examples, but this ratio will vary by application based on the task difficulty, satisfaction threshold(s), and current model quality.

<sup>6</sup>Additional detail can be found in Appendix E.

Method	Pr.	Re.	F1
Uncertainty Top	0.39	0.99	0.56
(Pr. $\geq 0.5$ )	0.50	0.04	0.07
Uncertainty Gap	0.38	1.00	0.55
(Pr. $\geq 0.5$ )	0.50	0.04	0.07
Satisfaction Regex	0.91	0.27	0.42
Satisfaction Classifier (1k)	0.84	0.84	0.84
Satisfaction Classifier (2k)	0.89	0.84	0.87
Satisfaction Classifier (5k)	0.94	0.82	0.88
Satisfaction Classifier (20k)	0.96	0.84	0.89
Satisfaction Classifier (40k)	0.96	0.84	<b>0.90</b>

Table 4: The maximum F1 score (with corresponding precision and recall) obtained on the SATISFACTION task. For the Uncertainty methods, we also report the maximum F1 score with the constraint that precision must be  $\geq 0.5$ . The Satisfaction Classifier is reported with varying numbers of SATISFACTION training examples.

collected based on failure modes of the current model, making them potentially more efficient in a manner similar to new training examples selected via active learning. It also suggests that the gains we observe in Table 3 might be further improved by (a) collecting FEEDBACK examples specific to each model (rather than using the same 60k FEEDBACK examples for all models), and (b) more frequently retraining the MTL model (e.g., every 5k examples instead of every 20k) or updating it in an online manner. We leave further exploration of this observation for future work.

The same experiment repeated for HB DIALOGUE examples found that fresher HB examples were no more valuable than stale ones, matching our intuition that HB DIALOGUE examples are less targeted at current model failure modes than FEEDBACK ones.

## 5.2 Predicting User Satisfaction

For maximum efficiency, we aim to ask for feedback when it will most benefit our model. The approach we chose (classifying the tone of partner responses) takes advantage of the fact that it is easier to recognize that a mistake has already been made than it is to avoid making that mistake; or in other words, sentiment classification is generally an easier task than next utterance prediction.

We compare this to the approach of asking for feedback whenever the model is most uncertain

what to say next. This approach acts on the assumption that the model will be least confident when it is about to make a mistake, which we find very frequently to not be the case. Not only is it difficult to recognize one’s own mistakes, but also there are often multiple valid responses to a given context (e.g., “Yes, I love seafood!” or “Yuck, fish is gross.”)—a lack of certainty about which to use does not necessarily suggest a poor model.

Table 4 reports the maximum F1 scores achieved by each method on the SATISFACTION test set. For the model uncertainty approach, we tested two variants: (a) predict a mistake when the confidence in the *top* rated response is below some threshold  $t$ , and (b) predict a mistake when the *gap* between the top two rated responses is below the threshold  $t$ . We used the best-performing standalone DIALOGUE model (one trained on the full 131k training examples) for assessing uncertainty and tuned the thresholds to achieve maximum F1 score. For the user satisfaction approach, we trained our dialogue agent on just the SATISFACTION task. Finally, we also report the performance of a regular-expression-based method which we used during development, based on common ways of expressing dissatisfaction that we observed in our pilot studies, see Appendix F for details.

As shown by Table 4, even with only 1k training examples (the amount used for the experiments in Section 5.1), the trained classifier significantly outperforms both the uncertainty-based methods and our original regular expression, by as much as 0.28 and 0.42 F1 points, respectively.

## 6 Future Work

In this work we learned from dialogue using two types of self-feeding: imitation of satisfied user messages, and learning from the feedback of unsatisfied users. In actuality, there are even more ways a model could learn to improve itself—for example, learning which question to ask in a given context to receive the *most* valuable feedback. One could even use the flexible nature of dialogue to intermix data collection of more than one type—sometimes requesting new FEEDBACK examples, and other times requesting new SATISFACTION examples (e.g., asking “Did my last response make sense?”). In this way, a dialogue agent could both improve its dialogue ability and its potential to improve further. We leave exploration of this meta-learning theme to future work.



## References

- M. A. Bassiri. 2011. Interactional feedback and the impact of attitude and motivation on noticing L2 form. *English Language and Literature Studies*, 1(2):61–73.
- P. Bojanowski, E. Grave, A. Joulin, and T. Mikolov. 2017. Enriching word vectors with subword information. *Transactions of the Association for Computational Linguistics (ACL)*, 5:135–146.
- R. Bunescu and R. Mooney. 2007. Learning to extract relations from the web using minimal supervision. In *Association for Computational Linguistics (ACL)*.
- M. Burtsev, V. Logacheva, V. Malykh, R. Lowe, I. Serban, S. Prabhunoye, E. Dinan, D. Kiela, A. Miller, K. Shuster, A. Szlam, J. Urbanek, and J. Weston. 2018. The conversational intelligence challenge 2 (ConvAI2).
- A. Carlson, J. Betteridge, B. Kisiel, B. Settles, E. R. H. Jr, and T. M. Mitchell. 2010. Toward an architecture for never-ending language learning. In *Association for the Advancement of Artificial Intelligence (AAAI)*.
- J. Devlin, M. Chang, K. Lee, and K. Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.
- M. Eskenazi, R. Evgeniia M. Shikib, and T. Zhao. 2018. Beyond turing: Intelligent agents centered on the user. *arXiv preprint arXiv:1803.06567*.
- B. Hancock, P. Varma, S. Wang, M. Bringmann, P. Liang, and C. Ré. 2018. Training classifiers with natural language explanations. In *Association for Computational Linguistics (ACL)*.
- C. Hashimoto and M. Sassano. 2018. Detecting absurd conversations from intelligent assistant logs by exploiting user feedback utterances. In *World Wide Web (WWW)*, pages 147–156.
- B. Hixon, P. Clark, and H. Hajishirzi. 2015. Learning knowledge graphs for question answering through conversational dialog. In *North American Association for Computational Linguistics (NAACL)*.
- T. Hong, O. Kwon, and Y. Kim. 2019. An end-to-end trainable task-oriented dialog system with human feedback. In *Association for the Advancement of Artificial Intelligence (AAAI)*.
- D. Kingma and J. Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.
- J. Kruger and D. Dunning. 1999. Unskilled and unaware of it: how difficulties in recognizing one’s own incompetence lead to inflated self-assessments. *Journal of personality and social psychology*, 77(6):1121–1134.
- E. Levin, R. Pieraccini, and W. Eckert. 2000. A stochastic model of human-machine interaction for learning dialog strategies. *IEEE Transactions on Speech and Audio Processing*, 8(1):11–23.
- J. Li, M. Galley, C. Brockett, J. Gao, and B. Dolan. 2016. A persona-based neural conversation model. In *Association for Computational Linguistics (ACL)*.
- J. Li, A. H. Miller, S. Chopra, M. Ranzato, and J. Weston. 2017a. Dialogue learning with human-in-the-loop. In *International Conference on Learning Representations (ICLR)*.
- J. Li, A. H. Miller, S. Chopra, M. Ranzato, and J. Weston. 2017b. Learning through dialogue interactions by asking questions. In *International Conference on Learning Representations (ICLR)*.
- B. Liu, G. Tür, D. Hakkani-Tür, P. Shah, and L. Heck. 2018. Dialogue learning with human teaching and feedback in end-to-end trainable task-oriented dialogue systems. In *North American Association for Computational Linguistics (NAACL)*, volume 1, pages 2060–2069.
- Y. Luan, C. Brockett, B. Dolan, J. Gao, and M. Galley. 2017. Multi-task learning for speaker-role adaptation in neural conversation models. In *Association for Computational Linguistics and International Joint Conference on Natural Language Processing (ACL-IJCNLP)*, volume 1, pages 605–614.
- L. Luo, W. Huang, Q. Zeng, Z. Nie, and X. Sun. 2018. Learning personalized end-to-end goal-oriented dialog. *arXiv preprint arXiv:1811.04604*.
- N. Mallinar, A. Shah, R. Ugrani, A. Gupta, M. Gurusankar, T. K. Ho, Q. V. Liao, Y. Zhang, R. Bellamy, and R. Yates. 2019. Bootstrapping conversational agents with weak supervision. In *Association for the Advancement of Artificial Intelligence (AAAI)*.
- P. Mazaré, S. Humeau, M. Raison, and A. Bordes. 2018. Training millions of personalized dialogue agents. In *Empirical Methods in Natural Language Processing (EMNLP)*, pages 2775–2779.
- S. Mazumder, N. Ma, and B. Liu. 2018. Towards a continuous knowledge learning engine for chatbots. *arXiv preprint arXiv:1802.06024*.
- A. H. Miller, W. Feng, A. Fisch, J. Lu, D. Batra, A. Bordes, D. Parikh, and J. Weston. 2017. Parlai: A dialog research software platform. In *Empirical Methods in Natural Language Processing (EMNLP)*, pages 79–84.
- A. Pappu and A. Rudnicky. 2013. Predicting tasks in goal-oriented spoken dialog systems using semantic knowledge bases. In *Proceedings of the SIGDIAL 2013 Conference*, pages 242–250.
- A. Paszke, S. Gross, S. Chintala, G. Chanan, E. Yang, Z. DeVito, Z. Lin, A. Desmaison, L. Antiga, and A. Lerer. 2017. Automatic differentiation in pytorch.

- A. Radford, K. Narasimhan, T. Salimans, and I. Sutskever. 2018. Improving language understanding by generative pre-training. Technical report, OpenAI.
- S. Rao and H. Daumé. 2018. Learning to ask good questions: Ranking clarification questions using neural expected value of perfect information. pages 2737–2746.
- H. Rashkin, E. M. Smith, M. Li, and Y. Boureau. 2018. I know the feeling: Learning to converse with empathy. *arXiv preprint arXiv:1811.00207*.
- A. Ratner, S. H. Bach, H. Ehrenberg, J. Fries, S. Wu, and C. Ré. 2017. Snorkel: Rapid training data creation with weak supervision. In *Very Large Data Bases (VLDB)*, 3, pages 269–282.
- V. Rieser and O. Lemon. 2011. Reinforcement learning for adaptive dialogue systems: a data-driven methodology for dialogue management and natural language generation. *Springer Science & Business Media*.
- J. Ross, A. Zaldivar, L. Irani, and B. Tomlinson. 2009. Who are the turkers? worker demographics in amazon mechanical turk. Technical report, Department of Informatics, University of California, Irvine.
- S. Ruder. 2017. An overview of multi-task learning in deep neural networks. *arXiv preprint arXiv:1706.05098*.
- V. Sanh, T. Wolf, and S. Ruder. 2018. A hierarchical multi-task approach for learning embeddings from semantic tasks. *arXiv preprint arXiv:1811.06031*.
- J. Schatzmann, K. Weilhammer, M. Stuttle, and S. Young. 2006. A survey of statistical user simulation techniques for reinforcement-learning of dialogue management strategies. *The knowledge engineering review*, 21(2):97–126.
- J. Schmidhuber and R. Huber. 1991. Learning to generate artificial fovea trajectories for target detection. *International Journal of Neural Systems*, 2(1):125–134.
- I. V. Serban, R. Lowe, L. Charlin, and J. Pineau. 2015. A survey of available corpora for building data-driven dialogue systems. *arXiv preprint arXiv:1512.05742*.
- I. V. Serban, C. Sankar, M. Germain, S. Zhang, Z. Lin, S. Subramanian, T. Kim, M. Pieper, S. Chandar, N. R. Ke, et al. 2017. A deep reinforcement learning chatbot. *arXiv preprint arXiv:1709.02349*.
- K. Shuster, S. Humeau, A. Bordes, and J. Weston. 2018a. Engaging image chat: Modeling personality in grounded dialogue. *arXiv preprint arXiv:1811.00945*.
- K. Shuster, S. Humeau, H. Hu, A. Bordes, and J. Weston. 2018b. Engaging image captioning via personality. *arXiv preprint arXiv:1810.10665*.
- D. L. Silver, Q. Yang, and L. Li. 2013. Lifelong machine learning systems: Beyond learning algorithms. In *Association for the Advancement of Artificial Intelligence (AAAI)*, volume 13.
- F. Strub, H. D. Vries, J. Mary, B. Piot, A. Courville, and O. Pietquin. 2017. End-to-end optimization of goal-driven and visually grounded dialogue systems. *arXiv preprint arXiv:1703.05423*.
- S. Tong and D. Koller. 2001. Support vector machine active learning with applications to text classification. *Journal of machine learning research*, 2(0):45–66.
- A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin. 2017. Attention is all you need. *arXiv preprint arXiv:1706.03762*.
- Y. Wang, B. Dai, L. Kong, X. Ma, S. M. Erfani, J. Bailey, S. Xia, L. Song, and H. Zha. 2018. Learning deep hidden nonlinear dynamics from aggregate data. In *Uncertainty in Artificial Intelligence (UAI)*.
- M. G. Werts, M. Wolery, A. Holcombe, and D. L. Gast. 1995. Instructive feedback: Review of parameters and effects. *Journal of Behavioral Education*, 5(1):55–75.
- J. E. Weston. 2016. Dialog-based language learning. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 829–837.
- H. Zhang, H. Yu, and W. Xu. 2017. Listen, interact and talk: Learning to speak via interaction. *arXiv preprint arXiv:1705.09906*.
- H. Zhang, H. Yu, and W. Xu. 2018a. Interactive language acquisition with one-shot visual concept learning through a conversational game. *arXiv preprint arXiv:1805.00462*.
- S. Zhang, E. Dinan, J. Urbanek, A. Szlam, D. Kiela, and J. Weston. 2018b. Personalizing dialogue agents: I have a dog, do you have pets too? *arXiv preprint arXiv:1801.07243*.

## A Data Collection Protocol

Here we report in greater detail the protocol we followed to collect the SATISFACTION, FEEDBACK, and HB DIALOGUE examples used in the experiments of Section 5.

We first trained our dialogue agent on just the DIALOGUE task with 20k HH examples. This agent was deployed on a crowdsourcing platform using the interface shown in Appendix H.2 to collect 2.5k SATISFACTION examples. These were split into 1k train, 500 validation, and 1k test examples. The agent was retrained using the 20k HH DIALOGUE examples and 1k SATISFACTION examples, then deployed to collect the first batch of deployment examples.

We collected 40k FEEDBACK examples (feedback set A) over the course of 17,250 conversations with 10 turns each (20 utterances, including the initial prompt). We then retrained the agent on all three datasets, using the same 20k HH DIALOGUE examples as before and only 20k of the available 40k FEEDBACK examples. This model was deployed to collect another 20k FEEDBACK examples (feedback set B), for a total of 60k FEEDBACK examples (A + B). In Table 3 we use these 60k FEEDBACK examples interchangeably; in Appendix E we compare them head-to-head. The 60k HB DIALOGUE examples were extracted from the logs of the deployment conversations. Finally, we collected an additional 40k SATISFACTION training examples to produce the numbers in Table 4 investigating the learning curve for this task.

No filtering was performed on the crowdworker conversations. Upon inspection after the fact, some workers did indeed give poor responses, make typographical mistakes, misunderstand the instructions, try to use the chatbot as a question answering interface, etc. We assume however that similar types of noise will be present in most chatbot deployment environments and opted to maintain a workflow that truly does not require developer intervention to use the newly collected examples.

## B Results with 10k Candidates

HH	HB	FB	Hits@X/10,000		
			@1	@10	@100
20k	-	-	0.8	4.6	16.2
20k	60k	60k	2.0	8.4	25.0
40k	-	-	1.3	6.5	21.8
40k	60k	60k	2.1	9.0	27.2
60k	-	-	1.6	7.0	24.0
60k	60k	60k	2.2	9.7	28.8
131k	-	-	2.5	10.0	30.3
131k	60k	60k	2.8	11.2	31.8

Table 5: When the number of candidates to choose from is increased to 10,000, adding Human-Bot (HB) DIALOGUE and FEEDBACK (FB) examples continues to improve performance on the DIALOGUE task at all levels.

## C PERSONACHAT Comparisons and Baselines

Our experiments use the PERSONACHAT distribution that was released as a part of the ConvAI2 (Burtsev et al., 2018) challenge. This distribution is slightly cleaner than the original PERSONACHAT release and comes with a new crowdsourced test set. In order to compare with the models and baselines used in the original PERSONACHAT paper (Zhang et al., 2018b), we report in this section the performance of our models on the original PERSONACHAT test set, not the ConvAI2 test set. Note that empirically, near Hits@1/20 = 50, each additional point of improvement corresponds to tens of thousands of fully-supervised Human-Human DIALOGUE examples. All numbers reported here are for models that do not have access to the profiles that were used in the creation of the conversations; models that do have access to this additional information tend to perform even better.

Model	Hits@1/20
<i>(Zhang et al., 2018b)</i>	
Seq2Seq	9.2
IR Baseline	21.4
Starspace	31.8
Profile Memory	31.8
KV Profile Memory	34.9
<i>Ours</i>	
Transformer	49.6
Self-Feeding	51.7

Table 6: The accuracy of various models and baselines on the original PERSONACHAT test set.

## D Using Chatbot Responses as Targets

HH	BF	BU	Hits@1/20
20k	-	-	<b>30.3</b>
20k	32k	-	22.7
20k	-	33k	19.3
131k	-	-	<b>44.7</b>
131k	32k	-	40.4
131k	-	33k	39.0

Table 7: Both with few HH DIALOGUE examples (20k) and many (131k), adding examples with bot utterances as the target decreased quality. We explored using all bot responses (Bot Unfiltered, or BU) and only those responses with estimated satisfaction scores greater than the 0.5 (Bot Filtered, or BF).

We also considered whether it was possible to consistently identify really good responses by the chatbot, rather than the really bad ones. These could potentially be used as DIALOGUE examples along with the ones that have human responses as targets (which we refer to as HH and HB in the paper). To explore this question, we modified our SATISFACTION dataset so that contexts with a rating of 5 were the positive class and ones with ratings [1, 2, 3] were the negative class (discarding ratings of 4 to increase the separation between classes). The results were negative—even with a training set of over 34k examples, the maximum precision we were able to achieve while maintaining at least 10% recall was 0.70, which is insufficient to improve performance on the DIALOGUE task. Upon inspection, it appears that really good responses are hard to identify because

most of the time they look like a normal human-to-human conversation, and recognizing an appropriate next utterance is precisely the DIALOGUE task that we are trying to solve! Negative responses, however, are much more semantically similar to one another, since most express one of a few common ideas such as asking for clarification or conveying confusion.

## E The Effect of Data Freshness

HH	HB <sub>A</sub>	HB <sub>B</sub>	FB <sub>A</sub>	FB <sub>B</sub>	Total	Hits@1/20
20k	-	-	-	-	20k	30.3
20k	40k	-	-	-	60k	35.4
20k	20k	20k	-	-	60k	35.3
40k	-	-	-	-	40k	36.2
20k	-	-	40k	-	60k	36.7
20k	-	-	20k	20k	60k	37.1
60k	-	-	-	-	60k	39.1

Table 8: As discussed in Section 5.1 and illustrated in Figure 3, FEEDBACK (FB) examples collected from a more recently retrained model (set B instead of set A) are more valuable in terms of improving performance; see Appendix A for details on how sets A and B were collected. We did not observe the same trend for HB DIALOGUE examples. We include the performance of models trained on only HH DIALOGUE examples in italics as reference points.

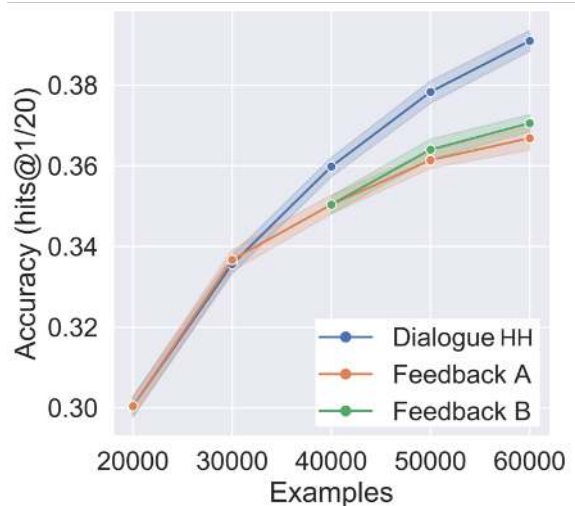


Figure 3: The first 20k examples for all models are supervised DIALOGUE examples. This model is deployed to collect 20k FEEDBACK examples (set A). If the model is retrained before collecting the next 20k examples (set B), the fresher feedback results in better performance ( $p = 0.027$ ). Shaded regions depict 95% confidence intervals.

## F SATISFACTION Regular Expressions

As described in Section 5.2, before we trained a classifier on the SATISFACTION task, we used the union of the following six regular expressions (using Python regular expression syntax) to identify user dissatisfaction and trigger feedback requests:

```
r"i .*(?:said|asked|told) .*"
r"((not|nt|n't) .*mak.*sense) | (mak.*no .*sense) "
r"u(m|h)+\W"
r"you.*what\?"
r"what.*you (?:mean|refer|talk) .*\?"
r"what.*to do with.*\?"
```

## G Hyperparameters

HH	HB	FB	layers	learning rate	loss factor	
					DIALOGUE	FEEDBACK
20k	-	-	1	0.0010	1.00	-
20k	20k	-	1	0.0010	1.00	-
20k	40k	-	1	0.0010	1.00	-
20k	60k	-	1	0.0010	1.00	-
20k	-	20k	1	0.0010	1.00	0.50
20k	-	40k	1	0.0010	1.00	0.50
20k	-	60k	1	0.0010	1.00	0.75
20k	60k	60k	1	0.0025	1.00	1.50
40k	-	-	1	0.0010	1.00	-
40k	20k	-	1	0.0010	1.00	-
40k	40k	-	1	0.0010	1.00	-
40k	60k	-	1	0.0025	1.00	-
40k	-	20k	1	0.0010	1.00	0.50
40k	-	40k	1	0.0010	1.00	0.75
40k	-	60k	1	0.0025	1.00	1.00
40k	60k	60k	1	0.0025	1.00	1.25
60k	-	-	2	0.0010	1.00	-
60k	20k	-	1	0.0025	1.00	-
60k	40k	-	1	0.0025	1.00	-
60k	60k	-	1	0.0025	1.00	-
60k	-	20k	1	0.0025	1.00	1.00
60k	-	40k	1	0.0025	1.00	1.00
60k	-	60k	1	0.0025	1.00	1.00
60k	60k	60k	1	0.0025	1.00	0.75
131k	-	-	2	0.0025	1.00	-
131k	20k	-	2	0.0025	1.00	1.00
131k	40k	-	2	0.0025	1.00	1.00
131k	60k	-	2	0.0025	1.00	1.00
131k	-	20k	2	0.0025	1.00	1.00
131k	-	40k	2	0.0025	1.00	1.00
131k	-	60k	2	0.0025	1.00	1.00
131k	60k	60k	2	0.0025	1.00	1.00

Table 9: The hyperparameters used to obtain the results in Table 3.

## H Data Collection Interfaces

### H.1 Deployment Data Collection

We simulated a deployment environment on a crowdsourcing platform with the following interface. Crowdworkers were told to “Talk as you would to a person you have just met.”

**Live Chat**

In this task, you will have a short conversation (10 responses) with a chatbot. Its purpose is to engage in small talk, *not* answer trivia questions. Talk as you would to a person you have just met.

**IMPORTANT:** The chatbot is trying to learn to recognize when it has made a mistake, so when it says something that doesn't make sense, *don't just continue the conversation*; instead, tell it in some way that it isn't making sense.

When the chatbot thinks it has made a mistake, it will ask you what it could have said instead. Even if it actually said something acceptable, answer the question about what else it could have said. Don't say "I don't know" or tell it that what it said was ok. Because the chatbot is still learning, you may need to correct it multiple times in a conversation before it realizes it has made a mistake; this is expected behavior.

**Chatbot:** honeymooning at niagra falls with the new hubby

**You:** wow! congratulations! were you married in new york?

**Chatbot:** you are probably right , how long have you all been dating ?

**You:** hey, you didn't answer my question!

**Chatbot:** Oops! I think I messed up. Whether I messed up or not, what could I have said (in response to "wow! congratulations! were you married in new york?")?

**You:** no, we're actually both Canadian; we live in Montreal.

**Chatbot:** Thanks! I'll try to remember that. Can you pick a new topic for us to talk about now?

Please enter here... **Send**

### H.2 SATISFACTION Data Collection

The environment for collecting SATISFACTION ratings was very similar to the deployment environment, with the exception that workers were required to give a rating 1-5 on the quality of each response given by the chatbot at the end of their own responses.

In this task, you will have a short conversation (10 responses) with a chatbot. Each time you respond to the bot, you will also rate the quality of its response: (1 - very bad, 2 - bad, 3 - ok, 4 - good, 5 - very good).

Occasionally you will be told that no rating is required (e.g., the initial prompt); otherwise, your response should always end with a space followed by a number between 1 and 5.

The chatbot is trying to learn to recognize when it has made a mistake, so when it says something that doesn't make sense, *don't just continue the conversation*; instead, tell it that is isn't making sense. In other words, if you are giving its response a rating of 1, then your response should indicate your dissatisfaction in some way.

Occasionally when the chatbot messes up, it will ask you what it should have said. If there isn't an obvious right answer, then think of something that it could have said that would have fit in the conversation. Don't just say "I don't know".

**System:** Start a conversation with the chatbot. If you need ideas for a topic, here are a few random words: houseboat, party, office. (No rating required)

**You:** howdy! where are you from?

**Chatbot:** hi there , my name is dustin , how are you ? (Reminder: At the end of your response, include a space followed by a rating 1-5.)

**You:** nice to meet you dustin! I'm doing well. 3

**Chatbot:** it is a fiction novel about a young teacher like myself solving murders .

**You:** um, that was super random. 1

Please enter here... **Send**

## I Sample Deployment Conversations

A sample of the conversations that occurred between a Bot (B) and Human (H) during the collection of the deployment datasets (more specifically, feedback set B). The self-feeding chatbot used its satisfaction classifier trained on 1k examples to determine when to ask for feedback. The dialogue model at this point had only been trained on 20k HH DIALOGUE examples and 20k FEEDBACK examples, so it still makes frequent mistakes. See Appendix H.1 for a snapshot of the actual interface.

B: start a conversation with the chatbot. if you need ideas for a topic, here are a few random words: random, mistake, century.  
H: what's the last mistake you made?  
B: what do you do ? i've a toothpick business ,  
H: that's not what i asked.  
B: oops! i think i messed up. whether i messed up or not, what could i have said (in response to "what's the last mistake you made?")?  
H: told me about a recent mistake.  
B: thanks! i'll try to remember that. can you pick a new topic for us to talk about now?  
H: what do you do for a living?  
B: i am a psychology major . what do you do ?  
H: i am a student as well. i'm studying law.  
B: oh , my five children and husband all speak spanish as well .  
H: i never said i speak spanish, what are you doing?  
B: oops! i think i messed up. whether i messed up or not, what could i have said (in response to "i am a student as well. i'm studying law.")?  
H: talked about your college courses.  
B: thanks! i'll try to remember that. can you pick a new topic for us to talk about now?  
H: how old are your children?  
B: it is . we are just thrilled we have ac again . its calming my anxiety .  
H: that has nothing to do with what i said.  
B: yes since i prefer my hair straight . i straighten it with a flat iron . tedious !  
H: you make no sense.

B: start a conversation with the chatbot. if you need ideas for a topic, here are a few random words: picnic, drain, elm.  
H: do you like to go to the beach or woods for a picnic?  
B: good , going to the doctor soon , i've some issue with my knees  
H: that's doesn't make sense  
B: oops! i think i messed up. whether i messed up or not, what could i have said (in response to "do you like to go to the beach or woods for a picnic?")?  
H: you could have chose the beach or woods to previous question  
B: thanks! i'll try to remember that. can you pick a new topic for us to talk about now?  
H: do you travel?  
B: i want to go to paris and make clothes .  
H: wow that's interesting. i've never been to paris  
B: i can see that . do you like music ?  
H: i love music. what is your favorite kind?  
B: tell me more about you  
H: that doesn't answer the question  
B: oops! i think i messed up. whether i messed up or not, what could i have said (in response to "i love music. what is your favorite kind?")?  
H: you could have selected a type of music that you like  
B: thanks! i'll try to remember that. can you pick a new topic for us to talk about now?  
H: do you like animals? i have a dog  
B: i like those , too , but pizza is better ! ;  
H: umm, what?

B: start a conversation with the chatbot. if you need ideas for a topic, here are a few random words: boy, punk, ray.  
H: are you a boy or a girl?  
B: do you have any jokes ?  
H: i wasn't talking about jokes.  
B: oops! i think i messed up. whether i messed up or not, what could i have said (in response to "are you a boy or a girl?")?  
H: i was talking about your gender.  
B: thanks! i'll try to remember that. can you pick a new topic for us to talk about now?  
H: sure, do you like music?  
B: oh that's neat . you should write me a poem about star trek !  
H: what?  
B: oops! i think i messed up. whether i messed up or not, what could i have said (in response to "sure, do you like music?")?  
H: we were talking about music.  
B: thanks! i'll try to remember that. can you pick a new topic for us to talk about now?  
H: what is your favorite type of art?

B: nice ! are the pets your own ?  
H: no  
B: oops! i think i messed up. whether i messed up or not, what could i have said (in response to "what is your favorite type of art?")?  
H: we weren't talking about pets.  
B: thanks! i'll try to remember that. can you pick a new topic for us to talk about now?  
H: sure, do you know anything about biology?

B: start a conversation with the chatbot. if you need ideas for a topic, here are a few random words: college, cliff, tachometer.  
H: hey, how's it going?  
B: no . i am just 15 . parents would kill me .  
H: what are you talking about?  
B: oops! i think i messed up. whether i messed up or not, what could i have said (in response to "hey, how's it going?")?  
H: i was just saying hi, so you should have said a greeting?  
B: thanks! i'll try to remember that. can you pick a new topic for us to talk about now?  
H: no problem. so, do you go to school?  
B: i have ten of them !  
H: no, that isn't right.  
B: oops! i think i messed up. whether i messed up or not, what could i have said (in response to "no problem. so, do you go to school?")?  
H: i was asking about school or education.  
B: thanks! i'll try to remember that. can you pick a new topic for us to talk about now?  
H: so, what is the best subject to you?  
B: oh wow , i've no cats , but i do have a couple of dogs .  
H: so, you like dogs a lot?  
B: i am also preparing for a marathon also , wish me luck .  
H: oh wow! does your dog help you prepare?  
B: i like hindu food . i'm not a very good baker though . my mom has mastered that .  
H: no, that isn't right.

B: start a conversation with the chatbot. if you need ideas for a topic, here are a few random words: drink, wool, bag.  
H: what is your favorite drink  
B: i love animals and the outdoors , so i work at the zoo .  
H: wrong answer  
B: oops! i think i messed up. whether i messed up or not, what could i have said (in response to "what is your favorite drink")?  
H: you should have said my favorite drink is either juice, coke or wine  
B: thanks! i'll try to remember that. can you pick a new topic for us to talk about now?  
H: are you religious  
B: somewhat lol he sits when i draw and watches me  
H: what are you talking about  
B: oops! i think i messed up. whether i messed up or not, what could i have said (in response to "are you religious")?  
H: you should have told me yes i am religious in fact i attend church services at least once a week and if not then no i am not religious never attended church services  
B: thanks! i'll try to remember that. can you pick a new topic for us to talk about now?  
H: tell me about yourself  
B: i drive my dream car every day i am home . a 65 mustang .  
H: that is nice you must be really love cars  
B: hello , how is your day going ?  
H: my day is good relaxing at home after work watching tv  
B: i'm sorry to hear that . i sincerely hope things get better for you .  
H: what are you talking about

B: start a conversation with the chatbot. if you need ideas for a topic, here are a few random words: glasses, puppy, pillow.  
H: do you like dogs? i got a new puppy  
B: hey ! i love being outside planting my plants ! how about you ?  
H: that is nice i love to plant too!  
B: right now i've a beagle . she's a little feisty . i miss my shepard mix .  
H: aww my dog is named sam he is a collie  
B: yeah , animals are my passion , i have 3 , all big dogs .  
H: yeah i just have the one dog and two cats  
B: i do not . i take a lot of pictures . i'd like to be a famous photographer .  
H: i actually really want to get a picture of my animals with santa lol  
B: do you want to work out with me ?  
H: umm what?  
B: oops! i think i messed up. whether i messed up or not, what could i have said (in response to "i actually really want to get a picture of my animals with santa lol")?  
H: that would be very cute  
B: thanks! i'll try to remember that. can you pick a new topic for us to talk about now?  
H: i love to eat



## J Sample SATISFACTION Conversations

A sample of the conversations that occurred between a Bot (B) and Human (H) during the collection of the SATISFACTION dataset. Users were asked to rate each bot response on a scale from 1 (very bad) to 5 (very good). See Appendix H.2 for a snapshot of the actual interface.

B: start a conversation with the chatbot. if you need ideas for a topic, here are a few random words: invitation, dot, czar.

H: hello friend, i missed you

B: hello my name is ray

H: hi i'm Leah (reward:5)

B: hello my name is michael

H: no it's not (reward:1)

B: oops! sorry. what should i have said instead?

H: nice to meet you

B: thanks! i'll remember that. can you pick a new topic for us to talk about now?

H: do you have many friends?

B: lol i have kids to

H: what? (reward:1)

B: oops! sorry. what should i have said instead?

H: no but i have kids

B: thanks! i'll remember that. can you pick a new topic for us to talk about now?

H: what do you want for christmas?

B: i love to bake cakes

H: what? (reward:1)

B: oops! sorry. what should i have said instead?

H: i want a pony for christmas

B: start a conversation with the chatbot. if you need ideas for a topic, here are a few random words: popsicle, lawnmower, mom.

H: hi, are you having a good day so far?

B: i'm here . i'm doing well , just waiting for the hurricane to hit .

H: oh im sorry. where do you live? (reward:5)

B: miami . . . i live in little cuba , but i wish i lived in europe .

H: yeah, i bet you do (reward:5)

B: we live in aruba , in a van down by the river

H: doesnt make sense (reward:1)

B: oops! sorry. what should i have said instead?

H: that you were from miami, not aruba

B: thanks! i'll remember that. can you pick a new topic for us to talk about now?

H: do you own a car?

B: i only drive toyota cars

H: what's the fastest you've ever driven a car? (reward:2)

B: a yellow toyota camry convertible .

H: that doesnt make sense (reward:1)

B: oops! sorry. what should i have said instead?

H: how fast you have driven before

B: thanks! i'll remember that. can you pick a new topic for us to talk about now?

H: do you like staying up late?

B: start a conversation with the chatbot. if you need ideas for a topic, here are a few random words: stout, mop, through.

H: i just finished mopping my kitchen.

B: i am a barista at joe jacks coffee shop

H: that has nothing to do with what i said (reward:1)

B: oops! sorry. what should i have said instead?

H: talked about how you hate cleaning the kitchen

B: thanks! i'll remember that. can you pick a new topic for us to talk about now?

H: what music do you listen to?

B: listen about lectures about space and stuff

H: so you don't like to listen to music at all? (reward:4)

B: i mostly listen to jazz

H: jazz is a fun genre. i play the trombone (reward:5)

B: classical jazz is my very favorite music

H: do you play any instruments? (reward:2)

B: classical music . mozart piano music .

H: i did not ask that (reward:1)

B: oops! sorry. what should i have said instead?

H: the instrument that you play

B: thanks! i'll remember that. can you pick a new topic for us to talk about now?

H: what are your plans today?

B: start a conversation with the chatbot. if you need ideas for a topic, here are a few random words: clown, drip, whatever.  
H: my pipes are dripping, do you know a plumber i can call?  
B: anything but taylor swift lol  
H: that doesn't make any sense. (reward:1)  
B: oops! sorry. what should i have said instead?  
H: you could say "sorry i don't know anyone who could get out there today."  
B: thanks! i'll remember that. can you pick a new topic for us to talk about now?  
H: i love music. who do you like to listen to?  
B: i mostly listen to smooth jazz and pop music  
H: oh smooth jazz eh? do you like kenny g? (reward:5)  
B: i like pop jazz and rap  
H: you already said that! (reward:1)  
B: oops! sorry. what should i have said instead?  
H: you could say "kenny g is a little too smooth for my tastes."  
B: thanks! i'll remember that. can you pick a new topic for us to talk about now?  
H: i think i'm going to just eat cereal with breakfast. what did you have to eat?  
B: bacon bacon bacon bacon bacon bacon bacon  
H: haha bacon is pretty good. i like the maple smoked stuff. (reward:3)  
B: onion green peppers mushrooms cream cheese shredded cheese .  
H: that sounds like an awesome bagel. how often do you eat breakfast? (reward:3)

B: start a conversation with the chatbot. if you need ideas for a topic, here are a few random words: princess, dump, telephone.  
H: hello! how are you on this fine day?  
B: i am pretty good i guess . i am even better at eating haha especially cookies and ice cream  
H: i do love eating too, especially cookies (reward:5)  
B: cakes , muffins , cookies , cupcakes ! sugar cookies are my favorite !  
H: you like a lot of sweets! (reward:3)  
B: cookies pies cakes and cup cakes  
H: okay, i get it, you like sweets (reward:2)  
B: peanut butter goes well with bread and crackers  
H: i do love a peanut butter sandwich, my kids do too (reward:5)  
B: tights with peanut butter peas  
H: what??? (reward:1)  
B: oops! sorry. what should i have said instead?  
H: something that made sense about peanut butter  
B: thanks! i'll remember that. can you pick a new topic for us to talk about now?  
H: do you have a cell phone?  
B: hostage , it was a robbery in the store i was selling furniture too  
H: huh? (reward:1)  
B: oops! sorry. what should i have said instead?  
H: tell me if you have a cell phone, maybe what kind you have

B: start a conversation with the chatbot. if you need ideas for a topic, here are a few random words: ray, scared, shirt.  
H: hi cutie  
B: hi  
H: how are you? (reward:4)  
B: good just making brisket and ribs you ?  
H: that sounds tasty. i'm just watching youtube (reward:5)  
B: i see . . i like watching movies too star wars especially  
H: star wars is a good movie (reward:5)  
B: horror movies or tv ? i like both i suppose  
H: yeah horror movies are ok (reward:5)  
B: watching a classic movie the wizard of oz  
H: do you like to swim? (reward:5)  
B: the beach and bike riding  
H: that sounds fun (reward:5)  
B: i hike at the park  
H: what else do you do? (reward:4)  
B: i run at central park  
H: oh you're from new york? (reward:5)  
B: upstate new york near jamestown ny  
H: wow awesome, i live in canada (reward:5)