

# Representation Learning for Conversational Data using Discourse Mutual Information Maximization

Anonymous ACL submission

## Abstract

Although many pretrained models exist for text or images, there have been relatively fewer attempts to train representations specifically for dialog understanding. Prior works usually relied on finetuned representations based on generic text representation models like BERT or GPT-2. But such language modeling pre-training objectives do not take the structural information of conversational text into consideration. Although generative dialog models can learn structural features too, we argue that the structure-unaware word-by-word generation is not suitable for effective conversation modeling. We empirically demonstrate that such representations do not perform consistently across various dialog understanding tasks. Hence, we propose a structure-aware Mutual Information based loss-function DMI (Discourse Mutual Information) for training dialog-representation models, that additionally captures the inherent uncertainty in response prediction. Extensive evaluation on nine diverse dialog modeling tasks shows that our proposed DMI-based models outperform strong baselines by significant margins.

## 1 Introduction

Representation learning has transformed how we can apply machine learning to solve real-world problems. However, despite a vast body of research on pretrained language representations, there have been relatively fewer attempts to train representations specifically for dialog understanding. Prior works mostly relied on finetuned representations based on generic models like BERT (Devlin et al., 2019) or GPT-2 (Radford et al., 2019). In our experiments, we demonstrate that such representations do not perform uniformly across various dialog understanding tasks such as dialog-act classification, intent detection or dialog evaluation.

On the other hand, prior works on pretraining large-scale dialog models focused mainly on open-domain generation. These works evaluated their

models only on dialog generation (Zhang et al., 2020; Roller et al., 2021; Adiwardana et al., 2020) or tasks related directly to the pretraining objective (Henderson et al., 2020; Gao et al., 2020). Their effectiveness on other dialog understanding tasks like act classification or intent detection remains unexplored. So we ask the following research question: *Can we learn enriched representations directly at the pretraining phase that are specifically helpful for dialog understanding?*

Existing language modeling (causal or masked) pretraining objectives unfortunately are not the best to model dialogs for these reasons: (1) The model is not directly trained to learn the content discourse structure (e.g., context-response in dialogs). (2) Such models are trained to generate the response word-by-word rather than predicting a larger unit. (3) The inherent one-to-many nature of dialog generation implies that the encoding model should be able to capture uncertainty in the response prediction task, that such models ignore.

Hence, in this paper, we propose pretraining objectives for improved dialog modeling that turn the discourse-level organizational structure of texts from natural sources (e.g., documents, dialogs, or monologues) into a learnable objective. We call this objective the Discourse Mutual Information (DMI). The key insight towards the design of our pretraining objective is to capture representations that can account for a meaningful conversation out of a specific ordered sequences of utterances. We hope that a discourse-level pretraining objective with conversational data would guide the model to learn complex context-level features. For example, in Fig. 1, we illustrate the differences between standard language modeling (causal or masked) based pretraining objectives and a discourse-level reasoning task.

The second research question that we ask is *whether discourse-level features learned using self-supervised pretraining outperform word-level pre-*

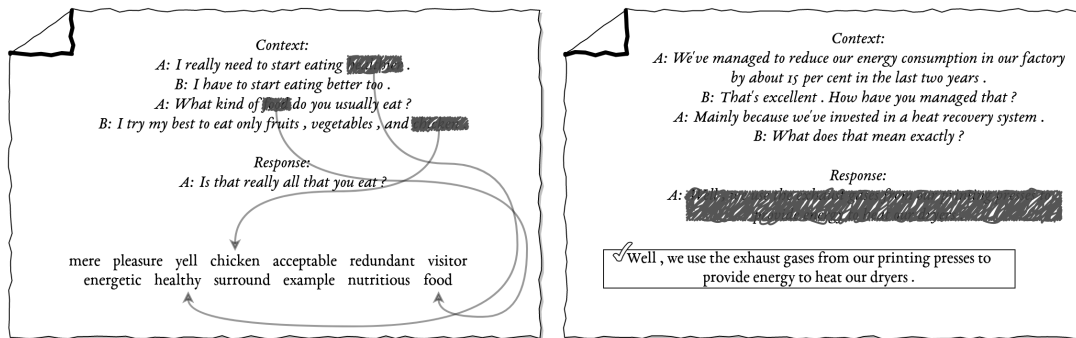


Figure 1: Possible reasoning involved in two types of pretraining: Word-level (left), Discourse-level (right). In a discourse-level reasoning task, the immediately preceding utterance may not be enough for understanding the full context. To predict the correct response, the model will need to capture both the larger context, in this case the topic of discussion, and the intent (e.g., asking for details) of the preceding utterance. In comparison, word-level reasoning is often easier and can be solved using local reasoning. Each of the three masked-words, in the left image, could have been predicted with reasonable confidence without any more information than the utterance itself.

084 *training objectives for downstream dialog under-*  
 085 *standing tasks.* Experimentally, we show that rep-  
 086 *resentations learned using the proposed objective*  
 087 *function are highly effective compared to both ex-*  
 088 *isting discriminative as well as generative dialog*  
 089 *models.* In terms of various dialog understanding  
 090 tasks, our models achieve state-of-the-art perfor-  
 091 mances in several tasks (absolute improvements up  
 092 to 8.5% and 3.5% in task accuracies in probing and  
 093 finetuning setups, resp.) and perform consistently  
 094 well across a variety of dialog understanding tasks,  
 095 whereas baseline models usually have a rather im-  
 096 balanced performance across tasks.

097 Overall, our main contributions are: (1) We pro-  
 098 pose DMI, a novel information-theoretic objective  
 099 function for pretraining dialog representation. (2)  
 100 We train models in two sizes (small and base) based  
 101 on our proposed self-supervised training objectives.  
 102 We will release our pretrained dialog representation  
 103 models, on acceptance of this paper. We make the  
 104 code publicly available<sup>1</sup>. (3) We extensively eval-  
 105 uate our DMI based representations on multiple  
 106 open-domain downstream tasks like intent detec-  
 107 tion, dialog-act classification, response retrieval,  
 108 dialog reasoning, and response-generation evalua-  
 109 tion, and beat state-of-the-art across 9 tasks in both  
 110 probe as well as finetune setups.

## 111 2 Literature Review

### 112 2.1 Dialog System Pretraining

113 There have been quite a few efforts towards uti-  
 114 lizing existing representations or developing new

<sup>1</sup><https://anonymous.4open.science/r/2022-DMI-anonymous>

115 pretrained models for dialog systems. While BERT  
 116 (Devlin et al., 2019), ELMo (Peters et al., 2018),  
 117 GPT-2 (Radford et al., 2019) and other general pur-  
 118 pose large-scale pretrained networks are not spe-  
 119 cific to dialogs, transfer learning from such models  
 120 could be reasonable. Basic language understanding  
 121 capability available through these representations  
 122 helps to get decent performance on many dialog-  
 123 understanding tasks (Hosseini-Asl et al., 2020).

124 On the other hand, there have been various works  
 125 on pretraining dialog specific representations or  
 126 large-scale generation models. We summarize the  
 127 properties of various previously proposed dialog-  
 128 representation learning models in Table 1. Di-  
 129 aloGPT (Zhang et al., 2020), Meena (Adiwardana  
 130 et al., 2020) and Blenderbot (Roller et al., 2021) are  
 131 large-scale Transformer-based language models,  
 132 which are trained to generate the gold-response (as  
 133 per the dataset) given a dialog context. ContextPre-  
 134 train (Mehri et al., 2019), ConveRT (Henderson  
 135 et al., 2020) and ConvFiT (Vulić et al., 2021) are  
 136 trained on the response retrieval task using Multi-  
 137 Woz or Reddit conversations. DEB or Dialog Eval-  
 138 uation using BERT (Sai et al., 2020) is a model  
 139 based on extended pretraining of the BERT archi-  
 140 tecture using Reddit data. DialogRPT (Gao et al.,  
 141 2020), on the other hand, is pretrained to predict  
 142 human-feedback (e.g., upvotes and downvotes) on  
 143 comments to Reddit threads. This model is initial-  
 144 ized using the weights of DialoGPT model. Wu  
 145 et al. (2020) thoroughly investigate these existing  
 146 pretrained representations, both generic and dialog  
 147 specific, for understanding their effectiveness on  
 148 various goal-oriented dialog-understanding tasks.

Model	Training Data Size	Pretraining Obj.	Architecture	Param	Downstream Task
<i>DialoGPT-small</i>	147M Dialogs	CE	GPT-2	125M	Generation /w MMI
<i>DialoGPT</i>	133M CR pairs	Response Ranking	DialoGPT	345M	Human Feedback Prediction
<i>Blenderbot-small</i>	1.5B comments	CE	Tr. S2S	90M	Generation
<i>Meena</i> ‡	40B words, 341 GB text	CE	Evolved Tr. S2S	2.4B	Generation
<i>ContextPretrain</i> ‡	10k Dialogs, MultiWoz	NUR, NUG, MUR, I2	HRED	-	Multiwoz (DST, Act, NUG, NUR)
<i>DEB</i>	727M Dialogs	MLM, NSP	BERT	110M	Adv/Random Dialog Evaluation
<i>ConveRT</i> †	727M Dialogs	Response Selection	Tr. Encoder	29M	Response Selection
<i>ConvFIT</i> ‡	8% of 727M Dialogs + Intent data	Response Selection	BERT	110M	Intent Detection
<i>DMI_Base</i>	7.5-10% of 727M Dialogs	InfoNCE-S	Tr. Encoder	124M	9 Dialog-NLU tasks

Table 1: Survey of Pretrained Dialog Models. NUR: next utterance retrieval, NUG: next utterance generation, MUR: masked utterance retrieval, I2: inconsistency identification, CR: Context-response, S2S: Seq2Seq, Tr.: Transformer, CE: Cross-entropy, HRED: Hierarchical RNN Encoder-Decoder. † Pretrained checkpoints available but only for inference. ‡ Both source-code and checkpoints are not available.

## 2.2 Self-supervised Representation Learning with InfoMax

Mutual Information maximization (InfoMax) is one of the popular approaches for self-supervised learning, first used by Oord et al. (2018) and Belghazi et al. (2018). Oord et al. (2018) proposed InfoNCE loss which is an estimator for lower bound to mutual information (MI) between two continuous-valued random variables. InfoNCE has also been used for other NLP applications like training sentence embeddings (SIMCSE (Gao et al., 2021)), question answering (QA-InfoMax (Yeh and Chen, 2019)), etc. Other estimators for mutual information have also been proposed like MINE (Mutual Information Neural Estimator) (Belghazi et al., 2018) and SMILE (Song and Ermon, 2020). In general, these estimators are also broadly studied in contrastive Learning (CL) literature for training both self-supervised (Mikolov et al., 2013; Devlin et al., 2019; Liu et al., 2019; Gao et al., 2021; Henderson et al., 2020; Vulić et al., 2021) and supervised models (Schroff et al., 2015; Gunel et al., 2020). In the next section, we derive our pretraining loss function DMI for conversational texts from an information-theoretic perspective.

## 3 Discourse Mutual Information

We define Discourse Mutual Information (DMI) as the mutual information<sup>2</sup> between two random variables representing two different segments within the same discourse. This is a general concept that can be applied to any form of discourse, no matter the domain or type of signal. In this paper, we focus on dialog type discourses and representation learn-

<sup>2</sup>Mutual Information between two random variables is defined as the reduction in uncertainty/entropy of one of the random variables by having knowledge about the value of the other random variable. Mathematically, this is written as  $I(X; Y) = H(X) - H(X|Y) = H(Y) - H(Y|X)$ .

ing for conversational texts. We define two random variables for the contexts ( $C$ ) and responses ( $R$ ) that jointly construct a valid conversation. Conversations between humans represent samples from the joint distribution  $P_{CR}$  of  $C$  and  $R$ . We pose the following learning problem, “*learn continuous representations for the textual random variables  $C$  and  $R$  such that the true mutual information between  $C$  and  $R$  can be closely estimated.*”

In the remainder of this section we show that, if the lower bound on MI estimated by some representations of context and response is close to the true value, the representation of the context would be as predictive of the response as the natural language form itself. Existing generative training objectives as used in DialoGPT or Blenderbot are extremely focused on predicting target response only. Per-word cross-entropy loss, used for training these models, fails to take into account the inherent uncertainty in the context-to-response generation function. Adapting context representations so as to predict the target responses optimally, helps our proposed DMI-based models learn better dialog representations applicable to a versatile set of dialog understanding tasks.

**Objective Function Formulation** Let  $E_c$  and  $E_r$  be the representations<sup>3</sup> for  $C$  and  $R$  based on some encoder. Using the data processing inequality from Information theory (Cover, 1999), we have

$$I(C; R) \geq I(E_c; E_r) \quad (1)$$

This tells us that MI between any encoded version of  $C$  and  $R$  will always be less or equal than the true mutual information. The equality will hold if  $E_c$  and  $E_r$  are both fully-invertible encoding processes (as opposed to representations which are

<sup>3</sup> $C, R, E_c, E_r$  in caps denote the random variables, whereas the lowercased versions  $c, r, e_c, e_r$  denote samples.

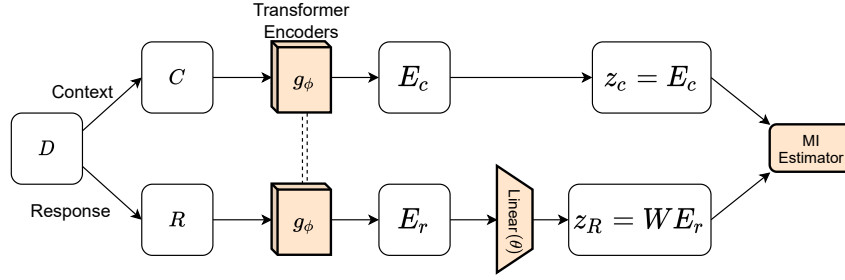


Figure 2: Base Pretraining Architecture for DMI. In our implementations of the model,  $f_\phi$  denotes the transformer (Vaswani et al., 2017) based encoders. Context and response encoders share all parameters for efficient learning.  $d$  denotes sample dialogs from the training dataset.

lossy or compressive, and inversion is thus not possible). However, neural networks generally embed the data points in a low dimensional manifold by learning robust features that can represent the data points efficiently. Because of this, neural representations are usually not invertible.<sup>4</sup> Now, the exact computation of MI is not possible for continuous-valued random variables. In recent years, various variational lower bounds have been proposed for estimating MI between continuous-valued random variables. Including the MI estimator ( $\hat{I}_\theta$ ), the overall relation becomes

$$I(C; R) \geq I(E_c; E_r) \geq \hat{I}_\theta(E_c; E_r) \quad (2)$$

This leads us to the proposed learning objective DMI:

$$\max_{\theta, \phi} \hat{I}_\theta(E_c^{(\phi)}; E_r^{(\phi)}) \quad (3)$$

where  $\hat{I}_\theta(E_c; E_r)$  is a variational lower bound estimate of  $I(E_c; E_r)$  (Equation 1) parametrized by  $\theta$  and  $\phi$  denotes the parameters of the encoder used for encoding  $C$  (or  $R$ ) to  $E_c$  (or  $E_r$ ).

**Loss function** For training our models, we minimize a loss function depending on the estimator being used.

$$\min_{\theta, \phi} \left[ L_{\theta, \phi}(C, R) = -\hat{I}_{\theta, estimator}(E_c^{(\phi)}, E_r^{(\phi)}) \right]$$

We experimented with various MI estimators from literature, namely, MINE (Belghazi et al., 2018), InfoNCE (Oord et al., 2018), JSD (Hjelm et al., 2019) and SMILE (Song and Ermon, 2020). These MI estimators generally compute samples of  $E_c$

<sup>4</sup>One general exception to this is a neural model/representation overfitted on some training data. In such cases, the model may exactly memorize the input/output pairs.

and  $E_r$  using  $C$  and  $R$  drawn from the joint distribution  $P_{CR}$ . Based on our preliminary experiments, we found that InfoNCE estimator produces better representations. The InfoNCE MI-estimate is computed as,

$$I(C; R) \geq \log N - L_N \quad (4)$$

$$L_N = -\frac{1}{2} \mathbb{E}_{P_{CR}} \left[ \log \frac{e^{f(c,r)}}{\sum_{r' \in R} e^{f(c,r')}} \right]$$

where  $N$  denotes the batch size, and  $f(c, r)$  is a scoring function for the  $\langle c, r \rangle$  pair.

**InfoNCE-S:** InfoNCE considers negative samples only for the response random variable. But, since the context and response variables have different sample spaces, we also introduce a symmetric version that considers negative samples ( $c'$  and  $r'$ ) for both context and response. This considerably improves the speed of training and convergence, and also gives a boost to downstream task performance.

$$L_N = -\frac{1}{2} \mathbb{E}_{P_{CR}} \left[ \log \frac{e^{f(c,r)}}{\sum_{r' \in R} e^{f(c,r')}} \right] - \frac{1}{2} \mathbb{E}_{P_{CR}} \left[ \log \frac{e^{f(c,r)}}{\sum_{c' \in C} e^{f(c',r)}} \right] \quad (5)$$

For other loss functions, more detailed discussion can be found in the Appendix.

**Comparison with ConveRT (Henderson et al., 2020):** There are a couple of differences between ConveRT's contrastive loss and our DMI objective. ConveRT models the problem as a response selection task and focuses on modeling cosine similarity between the context and the response. On the other hand, we propose a generic similarity computation function  $f(c, r)$  in Eqn. 4 and 5. Another difference is in encoding the input. ConveRT splits the context into previous turns and current query, and



277 encodes them independently. Our model encodes  
278 the entire context jointly and hence is capable of  
279 better learning the correlations between previous  
280 turns and current query.

## 281 4 DMI vs. Language Modeling Objectives

282 In this work, we focus on utilizing DMI for pre-  
283 training dialog representations incorporating strong  
284 discourse-level features. *But why should the DMI*  
285 *objective learn better discourse-level features than*  
286 *models trained on conversational data using MLM*  
287 *or LM objectives?* We can find the answer by look-  
288 ing at various LM-based objectives through the lens  
289 of InfoMax, as shown by Kong et al. (2020). They  
290 connected various pretraining objectives for natural  
291 language representations, including the ones used  
292 for training Skipgram, BERT and XLNet, to the  
293 InfoMax learning principle.

294 If we consider an input text  $T$  and a masking  
295 function  $M$  that returns a masked text  $\tilde{T}$  and the  
296 masked word  $w$ , the MLM objective is equivalent  
297 to  $\mathcal{L}_{MLM} = -\hat{I}_\theta(E^{(\phi)}(\tilde{T}), e_w)$  where,  $E^{(\phi)}$  is the  
298 language encoder (e.g., a Transformer encoder)  
299 and  $e_w$  is the embedding of the token  $w$ . Simi-  
300 larly, in the case of auto-regressive LMs like GPT-  
301 2, the InfoMax objective equivalent to the loss is  
302  $\mathcal{L}_{autoLM} = -\hat{I}_\theta(E^{(\phi)}(T_{1:t-1}), e_{T_t})$ , where  $T_{1:t-1}$   
303 is the input sequence till  $t - 1^{th}$  token and  $T_t$  is the  
304  $t^{th}$  token.

305 Compared to these LM objectives, DMI focuses  
306 on optimizing  $I(E_c, E_r)$ , where  $c$  and  $r$  are two  
307 structural components from the discourse with des-  
308 ignated roles. This enables DMI to discover more  
309 important features at the discourse level.

## 310 5 Experiments

### 311 5.1 Architecture

312 The exact encoder architecture and the pretraining  
313 pipeline has been shown in Figure 2. We use a  
314 dual encoder architecture for encoding the contexts  
315 and responses separately. We observe that sharing  
316 parameters between the two encoders leads to a  
317 more efficient learning process and faster conver-  
318 gence. We use vanilla transformer-based encoders<sup>5</sup>  
319 (Vaswani et al., 2017) for encoding the natural lan-  
320 guage inputs. The first tokens for both context and  
321 response sequences are the special [CLS] tokens  
322 whose contextual embeddings from the encoder  
323 are used as the context or response representations.

<sup>5</sup>We implemented all models and experiments using the PyTorch and Huggingface libraries.

The utterances in the context are delimited by an-  
other special token [EOU] (for end-of-utterance).  
We construct the context using as many utterances  
from the dialog history as possible up to a maxi-  
mum of 300 subword tokens. We use the Word-  
Piece tokenizer from BERT for tokenizing the input  
texts, with a vocabulary size of 30,522.

The scoring function  $f(c, r)$  in Eqs. 4 and 5  
is implemented using a Bilinear dot product be-  
tween the context and response representations:  
 $f(c, r) = e_c^T W e_r$  where,  $W$  is a square weight  
matrix trained along with other parameters in the  
model. This function can take any real value, pos-  
itive or negative, thus allowing the  $\hat{I}_\theta(E_C^{(\phi)}; E_R^{(\phi)})$   
function to take any positive real value. While  
any complicated function with that range could be  
chosen, we chose this as a simple formulation sat-  
isfying the range constraint and left most of the  
learning to the transformer and the projection ma-  
trix  $W$ .

### 344 5.2 Model Variants

345 We train two different scales of the DMI model:  
346 **DMI\_Small** with 8 layers, 8 attention heads, 768  
347 emb. size and **DMI\_Base** with 12 layers, 12 at-  
348 tention heads, 768 emb. size. **DMI\_Small** is  
349 trained from scratch on the pretraining dataset (see  
350 §5.4). **DMI\_Base** is initialized by weights from  
351 RoBERTa-base pretrained checkpoint, and further  
352 pretrained on the pretraining dataset. Both these  
353 models are trained using the InfoNCE-S estimator,  
354 unless specified otherwise.

### 355 5.3 Hyper-parameter Settings

356 We use Adam optimizer with a linear learning rate  
357 schedule for training both the models. Learning  
358 rate is first linearly increased to a max value of  $5 \times$   
359  $10^{-5}$  during the warm-up phase (first 1000 steps).  
360 Following this, in the remaining training period,  
361 learning rate is linearly decayed down back to zero.  
362 Before training **DMI\_Base**, we reset the parameters  
363 of the 12th self-attention layer, and it is trained  
364 again from scratch along with the weight matrix  
365  $W$  using our DMI objective. The embedding layer  
366 and initial 11 self-attention layers of the RoBERTa-  
367 base encoder are finetuned at a slower learning rate  
368 ( $5 \times 10^{-6}$ ) during our pretraining phase.

369 As the mutual information value obtained by  
370 the InfoNCE loss is upper bounded by  $\log(N)$ ,  $N$   
371 being the batch size, we try to keep the value of  $N$   
372 as large as possible. Both 8 and 12-layer models are  
373 trained on 4-GPU (4x32 GB V100s) systems with

Task	Description	Train	Valid	Test	Metric
Banking77	Intent 77-class Classification	8,002	2,001	3,080	Accuracy
SWDA	Dialog Act 41-class Classification	213,543	56,729	4,514	Accuracy
MuTual	Reasoning as Response Selection	25,516	2,836	3,544	R@1, R@2, MRR
MuTual Plus	MuTual + Safe response candidate	25,516	2,836	3,544	R@1, R@2, MRR
DD++	Dialog Evaluation	92,590	10,280	11,420	Accuracy
DD++/Adv	Train: Adv. neg., Test: Adv neg. samples	92,590	10,280	11,420	Accuracy
DD++/Cross	Train: Random neg., Test: Adv neg. samples	92,590	10,280	11,420	Accuracy
DD++/Full	Train: All samples, Test: Adv. neg. samples	138,885	10,280	11,420	Accuracy
Empathetic Intent	Emotion and Intent 44-class Classification	25,023	3,544	3,225	Accuracy

Table 2: Downstream task details. Adv.: Adversarial, Neg.: Negative

overall batch size of 480 and 384, respectively<sup>6</sup>. All the trained models will be publicly shared upon publication.

## 5.4 Pretraining Dataset

We pretrained all our models using the Reddit corpus (**Reddit-727M conversational-data**) released by Henderson et al., 2019. We ran the scripts released by the authors to recreate the dataset of 727M English conversations. Out of these 727M conversations, we utilize around 7.5% to 10% of the dataset to train our models, after which the validation loss generally saturates. In the rest of this paper, we will refer to this dataset as **rMax**, in short.

**Dialog Unrolling for Pretraining** For training our models, we need samples of context-response (CR) pairs. Each dialog is unrolled to create context-response pairs with each utterance in the dialog as a response, except the first one. Hence, for each dialog  $D = \{U_1, U_2, \dots, U_T\}$ , we generate the following set of samples  $S = \{(C_t : U_1, \dots, U_{t-1}; R_t : U_t) : t \in [2, T]\}$ . If we process the full rMax dataset, this leads to, approximately, 2.7B CR pairs.

## 5.5 MI Estimation

During pretraining, we compare the checkpoints from different epochs and across hyperparameter settings in terms of the bits of mutual information extracted by the trained representation on an unseen set of dialogs. This is calculated as  $MI_{valid} = \log(N) - L_N$  (see §3 for more details). As per the Information Bottleneck theory (Tishby et al., 2000), the mutual information learned between the two observed random variables can be factorized into two components, namely, predictive

<sup>6</sup>Training time: A maximum of 2 weeks of training time was allowed for 8-layer and 12-layer models. Though, the training process saturates long before the maximum allowed time, and we evaluate our models based on checkpoints when the best validation scores are first obtained.

and redundant information. Predictive information generally identifies whether the features learned by the representation are useful for a downstream task. The redundant information is caused by features that do not help in any downstream tasks. Such features can exist due to noise or spurious correlation in the dataset, or even overfitting. Hence, we train our final models on a fraction of the rMax dataset but only for one epoch (i.e., we never repeat the samples) which removes any possibility of overfitting.

Predictive features identified based on a fixed set of downstream tasks (Tishby et al., 2000; Alemi et al., 2017) may not be a sufficient to assess other features learned in the training process. Since, ideally, we want to maximize the amount of predictive information in the representation, we compare the bits of MI on the training set against the bits of MI on an unseen validation set, as captured by the learned representation. To make sure that we do not assume anything about the domain or the conversation topics, we use the validation set of dialogs from the open-domain Daily Dialog dataset (Li et al., 2017).

## 5.6 Downstream Tasks

Instead of focusing on a single downstream task like many previous works on dialog representation learning, we consider a more versatile range of tasks to evaluate the learned representations from DMI or the baseline models. To find out whether a certain representation is effective for some downstream task, we evaluate in two setups: probe and finetune. In both cases, the pretrained model is used along with an MLP classifier of fixed complexity (Pimentel et al., 2020). In probing setup, we only train the parameters of the MLP classifier. In finetuning setup, we also train the pretrained model parameters along with the MLP classifier parameters. We use the context and response representations from our models as the input to the MLP classifier.

	B77	SWDA	E-Intent	MuTual			MuTual Plus			DD++	DD++/adv	DD++/cross	DD++/full	
	Acc.	Acc.	Acc.	R@1	R@2	MRR	R@1	R@2	MRR	Acc.	Acc.	Acc.	Acc.	
Probing	RoBERTa_Base	72.84	67.18	50.45	<b>49.70</b>	<b>75.20</b>	<b>70.00</b>	43.60	66.60	65.10	55.75	84.20	65.11	68.76
	BERT_Base	72.74	67.99	46.84	45.40	72.80	67.30	42.60	67.70	64.90	60.39	86.56	65.25	72.50
	T5	60.82	68.79	44.50	43.20	69.40	65.60	38.30	65.70	62.20	57.46	84.14	61.23	63.35
	GPT-2	76.64	69.17	49.94	44.92	70.54	66.60	40.75	66.70	63.46	67.37	82.06	67.53	73.93
	DialoGPT_Small	53.00	65.10	43.42	29.80	53.50	55.15	25.51	57.56	54.05	63.63	78.02	<b>70.61</b>	70.77
	Blender_Small	70.39	70.11	48.52	41.42	68.06	64.29	42.89	68.85	65.18	60.07	65.14	57.76	68.20
	ConveRT	<b>89.88</b>	<b>71.36</b>	<b>55.47</b>	45.30	72.00	67.00	40.90	<b>69.00</b>	64.30	<b>79.14</b>	<b>88.67</b>	69.59	<b>80.86</b>
	DialogRPT	81.54	67.92	50.74	39.50	66.80	63.00	34.20	61.50	59.20	74.11	81.29	68.49	67.20
	DEB	79.18	68.50	45.31	45.10	74.00	67.50	<b>45.00</b>	67.70	<b>66.00</b>	70.66	86.07	67.25	67.77
	DMI_Small	88.80	71.71	53.89	51.35	76.86	71.07	44.81	71.67	66.75	85.79	90.88	76.81	87.31
	DMI_Base	91.43	72.73	60.00	52.48	76.41	71.65	48.98	71.33	68.73	86.91	91.98	79.15	88.32
	$\Delta$	1.55	1.37	4.53	2.78	1.21	1.65	3.98	2.33	2.73	7.77	3.31	8.54	7.46
	Finetuning	RoBERTa_Base	<b>92.75</b>	73.61	<b>62.81</b>	48.42	<b>77.20</b>	69.70	<b>49.55</b>	<b>73.70</b>	<b>69.50</b>	90.00	95.70	<b>73.76</b>
BERT_Base		92.27	72.29	60.12	47.86	73.93	68.80	49.10	72.35	69.00	87.05	94.33	67.70	88.82
T5		89.11	<b>73.77</b>	60.66	49.77	73.93	69.80	43.00	66.93	64.90	82.03	90.89	65.85	85.63
GPT-2		92.49	72.62	58.44	48.42	72.69	68.90	45.71	70.99	67.10	85.69	93.60	68.43	87.83
DialoGPT_Small		92.59	73.48	59.33	49.32	75.17	69.80	47.86	73.02	68.44	83.68	91.99	64.06	85.54
Blender_Small		91.59	71.10	58.31	<b>52.93</b>	75.85	<b>71.80</b>	47.97	70.99	68.30	86.83	92.29	66.39	87.82
DialogRPT		92.70	72.02	62.13	52.14	76.19	71.40	46.95	70.54	67.66	<b>90.26</b>	<b>95.81</b>	73.34	<b>91.25</b>
DEB		92.53	72.14	59.69	48.19	74.49	69.00	46.95	70.65	67.80	85.74	94.05	64.42	89.02
DMI_Small		91.95	72.02	59.41	53.84	75.28	72.15	46.95	72.91	68.21	86.97	91.76	74.41	86.87
DMI_Base		93.93	74.50	64.62	56.43	79.91	74.27	52.14	75.06	71.09	91.03	96.39	76.01	92.61
$\Delta$		1.18	0.73	1.81	3.50	2.71	2.47	2.59	1.36	1.59	0.77	0.59	2.25	1.35

Table 3: Results from probing (top) and finetuning (bottom) setups on 9 downstream tasks for assessing dialog understanding. (DD++: DailyDialog++, B77: Banking77 task, R@k: Recall at k, MRR: Mean reciprocal rank). Our model consistently performs better than SOTA on all the tasks in both probing as well as finetuning setups.

For downstream tasks, we have two reasoning tasks based on the MuTual dataset (Cui et al., 2020), three classification tasks based on conversational intent detection (Casanueva et al., 2020), emotion detection (Welivita and Pu, 2020) and act classification (Stolcke et al., 1998), and four dialog evaluation tasks based on the DailyDialog++ dataset (DD++, Sai et al., 2020)<sup>7</sup>. Table 2 shows dataset details and metrics for these nine tasks. Both MuTual and DailyDialog++ datasets have an adversarial configuration for the respective tasks, which allows us to assess each of the evaluated models in adversarial settings also.

## 6 Results and Discussions

### 6.1 Pretraining DMI based Representations

During pretraining, we used “Validation MI” to evaluate model checkpoints. As the goal of our models is to learn a representation that captures maximum MI between the context and the response texts, this metric tracks how well the learned representation captures the mutual information between contexts and responses of unseen dialogs.

We use the validation split from Daily Dialog dataset as our validation set for evaluation the model during pretraining. It is not specific to a domain and, hence, covers a versatile range of topics. This set comprises 1,000 full conversations between two persons which on unrolling leads to 7,069 context-response (CR) pairs. We illustrate

<sup>7</sup>Note that DailyDialog++ is different from DailyDialog.

the variation in validation-MI metric against training steps in Fig. 3 in the Appendix.

### 6.2 Comparison of Representations on Downstream Task Performance

In this set of experiments, we probe/finetune the DMI models with various downstream tasks that require knowledge of many different types of dialog-understanding features. The results of our probing and finetuning experiments are shown in Table 3.

We have used two types of models as our baselines: generic pretrained models and dialog-specific pretrained models. RoBERTa, BERT, T5 (Raffel et al., 2019), GPT-2 are all trained on large corpora of generic web-crawled English text. Whereas, DialoGPT, DialogRPT, DEB and ConveRT models were trained on conversational data. For DialogRPT, we used “human-vs-rand” checkpoint released by authors. All models are 12-layer except Blender-small (8 layers), ConveRT (6 layers), DialogRPT (24 layers) and DMI\_Small (8 layers). We used the publicly available model checkpoints for all baselines, wherever possible. The ConveRT model’s checkpoint has been removed from Github<sup>8</sup> by its authors. Hence, it was only possible for us to MLP-probe the representations, without finetuning of the model, based on a cached version released by another user under a valid license<sup>9</sup>. Pretrained checkpoints for Meena, Con-

<sup>8</sup><https://github.com/PolyAI-LDN/polyai-models>

<sup>9</sup><https://github.com/davidalami/ConveRT>

Model	B77	SWDA	E-Intent		MuTual		MuTual Plus			DD++	DD++/adv	DD++/cross	DD++/full
	Acc.	Acc.	Acc.	R@1	R@2	MRR	R@1	R@2	MRR	Acc.	Acc.	Acc.	Acc.
DMI_Base	<b>93.93</b>	<b>74.50</b>	64.62	56.43	<b>79.91</b>	74.27	<b>52.14</b>	<b>75.06</b>	<b>71.09</b>	<b>91.03</b>	96.39	76.01	92.61
DMI_Base - Sym	93.28	72.69	<b>65.18</b>	<b>57.34</b>	77.88	<b>74.32</b>	48.08	72.69	68.60	90.94	<b>96.65</b>	<b>76.45</b>	<b>93.13</b>
DMI_Base - RoB	92.34	74.10	60.96	53.84	77.31	72.47	50.34	72.80	69.81	87.23	92.95	73.53	87.85
DMI_Base - Sym - RoB	91.59	73.55	60.71	54.06	75.40	72.24	47.97	71.45	68.24	86.79	92.96	70.29	87.13

Table 4: Ablation study results for the finetune setup for our base model on 9 downstream tasks. “-RoB” → No RoBERTa initialization. “-Sym” → Training with non-symmetric version of InfoNCE. (DD++: DailyDialog++, B77: Banking77 task, R@k: Recall at k, MRR: Mean reciprocal rank).

textPretrain and ConvFiT are not available, and hence we do not compare with them.

### Results in Probing Setup

We observe that, on average, DEB and ConveRT have good performance among the baselines. However, the RoBERTa model outperforms all other baselines on the MuTual task by a significant margin. In the MuTual Plus task, the DEB model outperforms other models in the R@1 and MRR metrics. ConveRT performs the best among all baselines on the other tasks. ConveRT’s loss function is also contrastive in nature and is similar to ours. This explains the model’s generally high strength across the tasks among all the baselines.

Our DMI\_Base beats ConveRT on all the tasks, and DMI\_Small beats the baseline on 7 out of 9 tasks. We believe DD++ tasks to be the most demanding ones with respect to context-level understanding. Here, all non-dialog baselines have a weaker performance, with DEB and ConveRT being the best of the bunch. These are also the tasks where our models excel the most, with both DMI-Small and DMI-Base beating all baselines with strong margins. DD++/cross is the most difficult among all four DD++ tasks. Here, the model is trained on random negative samples and tested on a dataset with human-curated adversarial negatives. Our DMI\_Base beats the best baseline on DD++/cross by 8.54 points. This shows the superior quality of context representations from our models.

### Results in Finetuning Setup

In the finetuning setup, on average, RoBERTa and DialogRPT have good performance among the baselines. DialogRPT performs well for DD++ tasks while Blender works well for the MuTual task. For all other tasks, RoBERTa is the best baseline, even outperforming models especially trained for dialog tasks (like DialogPT).

Similar to the probe setup, DMI\_Base beats baseline methods by significant margins. In general, finetune results are better than probe results across all models, as expected.

Our large-scale RoBERTa-initialized DMI\_Base model outperforms the best baseline for all tasks, by a considerable margin. Additionally, our DMI-based models are able to perform well uniformly across all tasks, unlike even baselines like DialogPT, DialogRPT and Blenderbot models which are explicitly trained on dialog data. This makes DMI the best overall model for dialog related tasks. Across multiple tasks, we show qualitative examples where our proposed DMI-based models provide accurate results, in the Appendix.

### 6.3 Ablations

We evaluate the importance of using RoBERTa based pretraining as well as the symmetric version of the InfoNCE loss in Table 4. We observe that RoBERTa based pretraining helps significantly across all tasks. The symmetric InfoNCE improves performance for SWDA and MuTual Plus tasks.

## 7 Conclusions and Future work

In this paper, we proposed the concept of Discourse Mutual Information (DMI) which is better suited for learning dialog-specific features in a self-supervised manner. Using the InfoMax principle we formulated a pretraining method for dialog-specific representation learning. Across 9 downstream dialog understanding tasks, our 12-layer model outperforms state-of-the-art methods. Further, we showed that on most of these tasks, even our 8-layer model outperforms standard 12-layer pretrained models. These experiments show the potential of the proposed DMI objective towards building dialog understanding models. We make the code publicly available<sup>10</sup>. Although we experimented only with dialog modeling in this paper, we believe that the proposed DMI objective is generic enough to be applied to any type of discourse in any domain. In the future, we would like to explore how to harness DMI representations for generative conversation modeling.

<sup>10</sup><https://anonymous.4open.science/r/2022-DMI-anonymous>



## 8 Ethical considerations

Like many other pretrained language representation models, the proposed model may also have learned patterns associated with exposure bias. Interpretability associated with the output is rather limited, hence users should use the outputs carefully. The proposed model ranks possible response candidates, and does not filter out any “problematic” candidates. Thus, for applications, where candidate responses could be problematic, (e.g., offensive, hateful, abusive, etc.), users should carefully filter them out before providing them as input to our model.

All the datasets used in this work are publicly available. We did not collect any new dataset as part of this work.

Banking77 Casanueva et al., 2020 has been obtained from <https://github.com/PolyAI-LDN/task-specific-datasets>. It is available under a Creative Commons Attribution 4.0 International license with details here<sup>11</sup>.

SWDA Stolcke et al., 1998: The dataset has been obtained from <http://compprag.christopherpotts.net/swda.html>. This work is licensed under a Creative Commons Attribution-NonCommercial-ShareAlike 3.0 Unported License.

E-Intent Welivita and Pu, 2020: The dataset was downloaded from <https://github.com/anuradha1992/EmpatheticIntents>. The original dataset is available at <https://github.com/facebookresearch/EmpatheticDialogues> which is under the Creative Commons Attribution 4.0 International license.

MuTual and MuTual-plus Cui et al., 2020: The datasets have been downloaded from <https://github.com/Nealclly/MuTual>. Licensing is unclear; the authors do not mention any license information or terms of use.

DailyDialog++ Sai et al., 2020: The dataset was downloaded from <https://github.com/iitmnlp/DailyDialog-plusplus>. The data is available under the MIT License.

rMax or Reddit-727M conversational-data Henderson et al., 2019: the dataset has been obtained from <https://github.com/PolyAI-LDN/conversational-datasets/tree/master/reddit>. The dataset is available under the Apache License Version 2.0.

[//github.com/PolyAI-LDN/conversational-datasets/tree/master/reddit](https://github.com/PolyAI-LDN/conversational-datasets/tree/master/reddit). The dataset is available under the Apache License Version 2.0.

## References

- Daniel Adiwardana, Minh-Thang Luong, David R So, Jamie Hall, Noah Fiedel, Romal Thoppilan, Zi Yang, Apoorv Kulshreshtha, Gaurav Nemade, Yifeng Lu, et al. 2020. Towards a human-like open-domain chatbot. *arXiv preprint arXiv:2001.09977*.
- Alexander A. Alemi, Ian Fischer, Joshua V. Dillon, and Kevin Murphy. 2017. **Deep variational information bottleneck**. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*. OpenReview.net.
- Mohamed Ishmael Belghazi, Aristide Baratin, Sai Rajeswar, Sherjil Ozair, Yoshua Bengio, R. Devon Hjelm, and Aaron C. Courville. 2018. **Mutual information neural estimation**. In *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholm, Sweden, July 10-15, 2018*, volume 80 of *Proceedings of Machine Learning Research*, pages 530–539. PMLR.
- Iñigo Casanueva, Tadas Temčinas, Daniela Gerz, Matthew Henderson, and Ivan Vulić. 2020. **Efficient intent detection with dual sentence encoders**. In *Proceedings of the 2nd Workshop on Natural Language Processing for Conversational AI*, pages 38–45, Online. Association for Computational Linguistics.
- Thomas M Cover. 1999. *Elements of information theory*. John Wiley & Sons.
- Leyang Cui, Yu Wu, Shujie Liu, Yue Zhang, and Ming Zhou. 2020. **MuTual: A dataset for multi-turn dialogue reasoning**. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1406–1416, Online. Association for Computational Linguistics.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. **BERT: Pre-training of deep bidirectional transformers for language understanding**. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Tianyu Gao, Xingcheng Yao, and Danqi Chen. 2021. **SimCSE: Simple contrastive learning of sentence embeddings**. In *EMNLP*, pages 6894–6910, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Xiang Gao, Yizhe Zhang, Michel Galley, Chris Brockett, and Bill Dolan. 2020. **Dialogue response ranking**

<sup>11</sup><https://github.com/PolyAI-LDN/task-specific-datasets/blob/master/LICENSE>

691	<a href="#">training with large-scale human feedback data</a> . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 386–395, Online. Association for Computational Linguistics.	
692		
693		
694		
695		
696	Beliz Gunel, Jingfei Du, Alexis Conneau, and Veselin Stoyanov. 2020. Supervised contrastive learning for pre-trained language model fine-tuning. In <i>International Conference on Learning Representations</i> .	
697		
698		
699		
700	Matthew Henderson, Paweł Budzianowski, Iñigo Casanueva, Sam Coope, Daniela Gerz, Girish Kumar, Nikola Mrkšić, Georgios Spithourakis, Pei-Hao Su, Ivan Vulić, and Tsung-Hsien Wen. 2019. <a href="#">A repository of conversational datasets</a> . In <i>Proceedings of the First Workshop on NLP for Conversational AI</i> , pages 1–10, Florence, Italy. Association for Computational Linguistics.	
701		
702		
703		
704		
705		
706		
707		
708	Matthew Henderson, Iñigo Casanueva, Nikola Mrkšić, Pei-Hao Su, Tsung-Hsien Wen, and Ivan Vulić. 2020. <a href="#">ConveRT: Efficient and accurate conversational representations from transformers</a> . In <i>EMNLP</i> , pages 2161–2174, Online. Association for Computational Linguistics.	
709		
710		
711		
712		
713		
714	R. Devon Hjelm, Alex Fedorov, Samuel Lavoie-Marchildon, Karan Grewal, Philip Bachman, Adam Trischler, and Yoshua Bengio. 2019. <a href="#">Learning deep representations by mutual information estimation and maximization</a> . In <i>7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019</i> . OpenReview.net.	
715		
716		
717		
718		
719		
720		
721	Ehsan Hosseini-Asl, Bryan McCann, Chien-Sheng Wu, Semih Yavuz, and Richard Socher. 2020. <a href="#">A simple language model for task-oriented dialogue</a> . In <i>Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual</i> .	
722		
723		
724		
725		
726		
727		
728	Lingpeng Kong, Cyprien de Masson d’Autume, Lei Yu, Wang Ling, Zihang Dai, and Dani Yogatama. 2020. <a href="#">A mutual information maximization perspective of language representation learning</a> . In <i>8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020</i> . OpenReview.net.	
729		
730		
731		
732		
733		
734		
735	Yanran Li, Hui Su, Xiaoyu Shen, Wenjie Li, Ziqiang Cao, and Shuzi Niu. 2017. <a href="#">DailyDialog: A manually labelled multi-turn dialogue dataset</a> . In <i>Proceedings of the Eighth International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 986–995, Taipei, Taiwan. Asian Federation of Natural Language Processing.	
736		
737		
738		
739		
740		
741		
742	Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. <i>arXiv preprint arXiv:1907.11692</i> .	
743		
744		
745		
746		
	Shikib Mehri, Evgeniia Razumovskaia, Tiancheng Zhao, and Maxine Eskenazi. 2019. <a href="#">Pretraining methods for dialog context representation learning</a> . In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 3836–3845, Florence, Italy. Association for Computational Linguistics.	747
		748
		749
		750
		751
		752
		753
	Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013. Efficient estimation of word representations in vector space. <i>arXiv preprint arXiv:1301.3781</i> .	754
		755
		756
		757
	Aaron van den Oord, Yazhe Li, and Oriol Vinyals. 2018. Representation learning with contrastive predictive coding. <i>arXiv preprint arXiv:1807.03748</i> .	758
		759
		760
	Matthew Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. 2018. <a href="#">Deep contextualized word representations</a> . In <i>Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)</i> , pages 2227–2237, New Orleans, Louisiana. Association for Computational Linguistics.	761
		762
		763
		764
		765
		766
		767
		768
		769
	Tiago Pimentel, Josef Valvoda, Rowan Hall Maudslay, Ran Zmigrod, Adina Williams, and Ryan Cotterell. 2020. <a href="#">Information-theoretic probing for linguistic structure</a> . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 4609–4622, Online. Association for Computational Linguistics.	770
		771
		772
		773
		774
		775
		776
	Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. <i>OpenAI blog</i> , 1(8):9.	777
		778
		779
		780
	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2019. Exploring the limits of transfer learning with a unified text-to-text transformer. <i>arXiv preprint arXiv:1910.10683</i> .	781
		782
		783
		784
		785
	Stephen Roller, Emily Dinan, Naman Goyal, Da Ju, Mary Williamson, Yinhan Liu, Jing Xu, Myle Ott, Eric Michael Smith, Y-Lan Boureau, and Jason Weston. 2021. <a href="#">Recipes for building an open-domain chatbot</a> . In <i>Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume</i> , pages 300–325, Online. Association for Computational Linguistics.	786
		787
		788
		789
		790
		791
		792
		793
	Ananya B. Sai, Akash Kumar Mohankumar, Siddhartha Arora, and Mitesh M. Khapra. 2020. <a href="#">Improving dialog evaluation with a multi-reference adversarial dataset and large scale pretraining</a> . <i>Transactions of the Association for Computational Linguistics</i> , 8:810–827.	794
		795
		796
		797
		798
		799
	Florian Schroff, Dmitry Kalenichenko, and James Philbin. 2015. Facenet: A unified embedding for face recognition and clustering. In <i>Proceedings of</i>	800
		801
		802

803 *the IEEE conference on computer vision and pattern*  
804 *recognition*, pages 815–823.

805 Jiaming Song and Stefano Ermon. 2020. [Understanding](#)  
806 [the limitations of variational mutual information esti-](#)  
807 [mators](#). In *8th International Conference on Learning*  
808 *Representations, ICLR 2020, Addis Ababa, Ethiopia,*  
809 *April 26-30, 2020*. OpenReview.net.

810 Andreas Stolcke, Elizabeth Shriberg, Rebecca Bates,  
811 Noah Coccaro, Daniel Jurafsky, Rachel Martin,  
812 Marie Meteer, Klaus Ries, Paul Taylor, Carol Van  
813 Ess-Dykema, et al. 1998. Dialog act modeling for  
814 conversational speech. In *AAAI Spring Symposium*  
815 *on Applying Machine Learning to Discourse Process-*  
816 *ing*, pages 98–105.

817 Naftali Tishby, Fernando C Pereira, and William Bialek.  
818 2000. The information bottleneck method. *arXiv*  
819 *preprint physics/0004057*.

820 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob  
821 Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz  
822 Kaiser, and Illia Polosukhin. 2017. [Attention is all](#)  
823 [you need](#). In *Advances in Neural Information Pro-*  
824 *cessing Systems 30: Annual Conference on Neural*  
825 *Information Processing Systems 2017, December 4-9,*  
826 *2017, Long Beach, CA, USA*, pages 5998–6008.

827 Ivan Vulić, Pei-Hao Su, Samuel Coope, Daniela  
828 Gerz, Paweł Budzianowski, Iñigo Casanueva, Nikola  
829 Mrkšić, and Tsung-Hsien Wen. 2021. [ConvFiT: Con-](#)  
830 [versational fine-tuning of pretrained language mod-](#)  
831 [els](#). In *EMNLP*, pages 1151–1168, Online and Punta  
832 Cana, Dominican Republic. Association for Compu-  
833 tational Linguistics.

834 Anuradha Welivita and Pearl Pu. 2020. [A taxonomy of](#)  
835 [empathetic response intents in human social conversa-](#)  
836 [tions](#). In *Proceedings of the 28th International Con-*  
837 *ference on Computational Linguistics*, pages 4886–  
838 4899, Barcelona, Spain (Online). International Com-  
839 mittee on Computational Linguistics.

840 Chien-Sheng Wu, Steven C.H. Hoi, Richard Socher, and  
841 Caiming Xiong. 2020. [TOD-BERT: Pre-trained natu-](#)  
842 [ral language understanding for task-oriented dialogue](#).  
843 In *EMNLP*, pages 917–929, Online. Association for  
844 Computational Linguistics.

845 Yi-Ting Yeh and Yun-Nung Chen. 2019. [QAInfomax:](#)  
846 [Learning robust question answering system by mutu-](#)  
847 [al information maximization](#). In *EMNLP-IJCNLP*,  
848 pages 3370–3375, Hong Kong, China. Association for  
849 Computational Linguistics.

850 Yizhe Zhang, Siqi Sun, Michel Galley, Yen-Chun Chen,  
851 Chris Brockett, Xiang Gao, Jianfeng Gao, Jingjing  
852 Liu, and Bill Dolan. 2020. [DIALOGPT : Large-scale](#)  
853 [generative pre-training for conversational response](#)  
854 [generation](#). In *Proceedings of the 58th Annual Meet-*  
855 *ing of the Association for Computational Linguistics:*  
856 *System Demonstrations*, pages 270–278, Online. As-  
857 sociation for Computational Linguistics.

## A Mutual Information Estimators

In this paper, we experiment with various different MI estimators, and found InfoNCE-S to be the best (both in terms of accuracy as well as training speed). The mathematical formulation of these estimators is provided below.

1. InfoNCE was proposed by Oord et al. (2018).

It connects to the mutual information value  $I(X; Y)$  as,

$$I(X; Y) \geq \log(N) - L_N$$

$$L_N = -\mathbb{E}_{P_{XY}} \left[ \log \frac{e^{f(x,y)}}{\sum_{y' \in Y} e^{f(x,y')}} \right]$$

2. MINE (Belghazi et al., 2018)

$$I(X; Y) \geq \sup_{\theta \in \Theta} [I_{\theta}^{(MINE)}(X; Y) =$$

$$\mathbb{E}_{P_{XY}} [T(x, y)] - \log \mathbb{E}_{P_X \times P_Y} [e^{T(x,y)}]]$$

3. JSD (Hjelm et al., 2019)

$$I_{\theta}^{(JSD)}(X; Y) = \mathbb{E}_{P_{XY}} [-sp(-T(x, y))] - \mathbb{E}_{P_X \times P_Y} [sp(T(x, y))]$$

4. SMILE (Song and Ermon, 2020)

$$I_{\theta}^{(smile)}(X; Y) = \mathbb{E}_{P_{XY}} [T(x, y)]$$

$$- \log \mathbb{E}_{P_X \times P_Y} [\text{clip}(e^{T(x,y)}, e^{-\tau}, e^{\tau})]$$

Use of the InfoMax objective for self-supervised learning has been more prevalent in the computer vision domain than in NLP. Although as Kong et al. (2020) have previously shown, many existing loss functions used for training NLP models can be derived directly from the InfoMax framework. Kong et al. (2020) had only focused on various language model objectives that focus on words given the surrounding context. The authors showed that this objective translates to maximizing mutual information between the context and the missing word within the context.

In dialog domain also, InfoMax-equivalent loss functions have been used. First, Henderson et al. (2020) used contrastive formulation of the response selection task as a pretraining objective for dialog representation. Other prior works on response selection models often used a binary-cross entropy loss for training. Both these loss functions are actually equivalent to various lower bound estimators



for mutual information. In the QAInfoMax model (Yeh and Chen, 2019), the authors used the Deep-InfoMax loss function (Hjelm et al., 2019) as a regularizer and showed that representations learned with or in-presence of an InfoMax regularizer are more resilient to adversarial attacks while maintaining the same level of task performance. We also observe the same effect in our DD++/cross experiments. This is because of the self-supervised yet task-specific nature of the loss function.

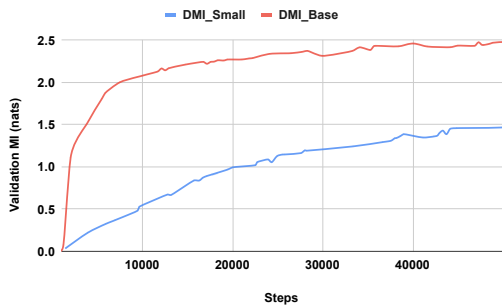


Figure 3: Validation-MI profile during pretraining

## B Response Retrieval Experiment

We wanted to investigate if the proposed model can rank good responses higher compared to more generic/bland ones. Hence to test against an extreme setting we simulate a response selection task for a very large pool using the test set of Daily Dialog (Li et al., 2017) dataset. We took all the  $\sim 7000$  responses from test set of the daily dialog dataset as the response pool. Next, for a few randomly selected context examples, we illustrate the top two ranked as well as ground truth responses for two full conversations in Tables 5 and 6. Of course, the ground truth response was removed from the pool for each context. The ranking of responses were done using the  $f(c, r)$  function from the trained DMI\_Base model. From the examples of response selection, we can observe that the model is able to both avoid bland responses and select responses that are relevant to the current context even from such a large pool. This shows the usefulness of dialog specific pretrained representation trained using the DMI objective.

## C Prediction Samples and Error Analysis

In Table 7, we show sample predictions from the DailyDialog++/cross task (Sai et al., 2020). As DialoGPT has the best performance in the probe setup, on this task among the baselines, we choose

it for error analysis. We randomly sampled 11 instances where the DialoGPT model made a mistake and observed the behavior of our DMI\_Base model on these samples. We see that our model correctly predicts for all 6 out of 6 negative samples and out of the 5 positive samples DMI\_Base predicts the label of 2 samples correctly (overall 8/11 correct predictions by our model). This shows that our model has a better understanding of the context and response inputs, which makes it robust against the adversarial negative samples. As can be seen in samples 2, 3, 5 and 6, the incorrect predictions by the DialoGPT model might have been caused by presence of common or similar meaning tokens (*cook, food; million; long; employee*) between context and response. This means that DialoGPT often relies on weak token-based cues for prediction.

For error analysis on the Empathetic-Intent (E-Intent) task (Welivita and Pu, 2020), we chose the ConveRT model as the baseline to compare against predictions from our DMI\_Base model. First, we randomly select 10 samples from the test set of the E-Intent task where the baseline ConveRT model makes a mistake. Then the predictions from the DMI\_Base model are observed on these 10 samples. The input utterances, true labels and the predictions made by the model are shown in Table 8. Out of these 10 samples, DMI\_Base is able to predict the labels for 6 instances correctly. We notice that though sample inputs often contain more than one emotion, the one denoted by the gold label is generally the primary one. Our model is able to capture this emotion correctly more often than the baseline, with such mixed-emotion samples.

Fig. 4 shows the confusion matrix for our DMI\_Base model for the Empathetic-Intent task. The accuracy is highest for afraid, acknowledging and questioning classes (each above 95%). Some of the most confusing pairs of classes are (annoyed, wishing), (anxious, apprehensive), (caring, confident), (content, grateful), (content, lonely).



Dialog	Rank-1 Response	Rank-2 Response
<b>S1:</b> Hey man , you wanna buy some weed ?	Yes , smoking or non-smoking ?	Sure , what ' s up ?
<b>S2:</b> Some what ?	My cigarettes .	Come on man ! I even got dope and acid ! Try some !
<b>S1:</b> Weed ! You know ? Pot , Ganja , Mary Jane some chronic !	I ' m sorry , Mrs . Lin . We don ' t allow smoking in the office .	Ha ! That ' s right . You used to smoke . I forgot !
<b>S2:</b> Oh , umm , no thanks .	Come on man ! I even got dope and acid ! Try some !	Come on , you can at least try a little , besides your cigarette .
<b>S1:</b> I also have blow if you prefer to do a few lines .	No , thanks . Oh , yes , could I have some juice for the mini bar ?	No , thanks . I ' Ve had enough . I ' ll have my bill , please .
<b>S2:</b> No , I am ok , really .	Great ! Would you prefer smoking or non-smoking ?	Come on , you can at least try a little , besides your cigarette .
<b>S1:</b> Come on man ! I even got dope and acid ! Try some !	I will , but let ' s go to the disco after dinner .	No , thanks . I ' ve had more than enough . In fact , I must be running along .
<b>S2:</b> Do you really have all of these drugs ? Where do you get them from ?	I found them at JC Pennies .	They are my personal effects .
<b>S1:</b> I got my connections ! Just tell me what you want and I ' ll even give you one ounce for free .	Hmm ... I ' ll tell you what . I ' ll lend you four thousand dollars , but you have to pay me back next week .	No , thanks . Oh , yes , could I have some juice for the mini bar ?
<b>S2:</b> Sounds good ! Let ' s see , I want .	Here you go . On any other day , it would cost me a fortune , but it ' s on special offer today .	All right , Maria . I ' ll give you until tomorrow at 4:00 to produce a satisfactory piece of work , but otherwise , you ' ll have to re-do it .
<b>S1:</b> Yeah ?	But of course ! Well , it ' s been great talking to you , but I have to get going .	Good . I was thinking that I ' d like to invite you to watch a movie . I can meet you at the cinema gate .

Table 5: Retrieval Example 1: Top two ranked responses from a large pool, as well as ground truth response for a conversation. Note that every line corresponds to one utterance in the conversation either from speaker S1 or speaker S2. **How to read the table:** For any context (all entries in first column upto any row), the ground truth response can be automatically obtained from the immediate next entry in first column. The response selected by the model in any cell  $(t, 1 \text{ or } 2)$  is for the true context from row 1 to row  $t$  in the first column.

Dialog	Rank-1 Response	Rank-2 Response
<b>S1:</b> Could I have my bill , please ?	Sure . Here is your receipt .	Sure . Your cash back is \$ 13 . And we ' ll bring out your fries in two minutes .
<b>S2:</b> Certainly , sir .	Thanks . Now can I make the full deposit ?	Thank you ! Would you like an aisle seat or a window seat ?
<b>S1:</b> I ' m afraid there ' s been a mistake .	I am really sorry too , maybe I can give you a call sometime .	Oh , I ' m sorry . However , if you could help me out , I ' ll double the pay for the hours worked .
<b>S2:</b> I ' m sorry , sir . What seems to be the trouble ?	Not much . I had to pay an unexpected bill , so I needed the money back .	Oh , nothing special . I ' m just a bit tired .
<b>S1:</b> I believe you have charged me twice for the same thing . Look , the figure of 6.5 dollar appears here , then again here .	One moment , please , sir . ... Here ' s your bill . Would you like to check and see if the amount is correct ?	Sir , I deleted the \$ 10 , but I had to add a \$ 2 service charge to your bill .

Table 6: Retrieval Example 2: Top two ranked responses from a large pool, as well as ground truth response for a conversation. Note that every line corresponds to one utterance in the conversation either from speaker S1 or speaker S2. **How to read the table:** For any context (all entries in first column upto any row), the ground truth response can be automatically obtained from the immediate next entry in first column. The response selected by the model in any cell  $(t, 1 \text{ or } 2)$  is for the true context from row 1 to row  $t$  in the first column.

ID	Context	Candidate Response	Gold	DialoGPT	DMI_Base
<i>DMI_Base predicts correctly</i>					
1	All right. I'll take it. __eou__ Do you like to use chopsticks __eou__ Yes, I like using chopsticks.	When you get closer, you see that each horizontal section is made up of two pieces that converge in a right angle.	0	1	0
2	And you'll have to sell your motorcycle. And your cameras. Right? __eou__ Maybe I'll cook once or twice a week. How is that?	I go to the temple twice a week so I prefer vegetarian food.	0	1	0
3	But I heard the box office rose up to 15 million in the first week. __eou__ Box office can't explain everything. I do not think it is cheerful or well-made. The plot is old and the female character is not pretty. __eou__ My sister has given me two tickets for tonight. It is called 'The life of Rose', a French movie.	I got 1 million views on my youtube channel in one week.	0	1	0
4	Glad you like it. By the way, is this your first time to China, Mr. White? __eou__ Yes, as a representative of IBM. I hope to conclude some business with you. __eou__ We also hope to expand our business with you.	May I know what and all process you have?	1	0	1
5	Good. I have to go right now. I really hope this meeting doesn't last too long. __eou__ They usually go on for ages. __eou__ I'll stop by if I have time later. Make sure everyone knows that we must stick to the deadlines.	I don't cut my hair because I really like to keep it long.	0	1	0
6	Of course. The main thing is that all our work must be completed on schedule. We even allow our employee to go home early if they finish their work early. __eou__ How often do you have meetings? __eou__ You should attend a department meeting every Monday morning. There are other meetings for people working together on certain projects. Department heads also attend an interdepartmental meeting each week.	In the newsletter, I gave employees column references this week.	0	1	0
7	Sounds interesting! That must be very convenient. __eou__ Yes, you're right. I can blog wherever and whenever I'm on the move. It's especially good when I'm on a business trip and my laptop happens to be away from me. __eou__ How can you do that?	I sank parents money into my business it is not convenient.	0	1	0
8	There is a wait right now to use the computers. __eou__ That's fine. __eou__ Would you please write your name on this list?	Sure, please give me a pen.	1	0	1
<i>DMI_Base predicts wrongly</i>					
9	How much cash would you like? __eou__ I want \$150. __eou__ Here's your \$150.	Well! I never forget your help.	1	0	0
10	I see, sir. This one is very good. __eou__ Is it? __eou__ You may rest assured. It sells well. __eou__ May I have a look at the introduction?	It has been recommended by top nutritionists.	1	0	0
11	Sir, tell us about your experience with Super Bulk-up. __eou__ Well, it's completely changed my life. __eou__ Tell us how.	The change is right in front of you, isn't it?	1	0	0

Table 7: Sample Predictions from the DD++/Cross task. In DD++/Cross, the models are trained using randomly sampled negatives and tested on curated adversarial negative samples. In each sample, the input context comprises the utterances, previous to the response, spoken by the two participants. Such utterances within a context are delimited by a special token “\_\_eou\_\_”.

ID	Input Utterance	Gold Label	ConveRT	DMI_Base
<i>DMI_Base predicts correctly</i>				
1	i feel very thankful for everything that i have, i live a really good life in my liking	grateful	content	grateful
2	I'm training a new girl at work. She is doing so good for her first week!	proud	confident	proud
3	It broke my heart today when I went to the grocery store and found out that they were out of Dean's French Onion Dip.	disappointed	devastated	disappointed
4	My wife's birthday is coming up. I got her a gift and the party planned out way ahead of time this year.	prepared	surprised	prepared
5	My friend helped me to pack	grateful	trusting	grateful
6	For two years now I've been walking with help of a walker, following a botched hip operation. Recently, at a physical therapy session, I was able to walk with a cane the length of the treatment room. I felt quite good about myself!	proud	caring	proud
<i>DMI_Base predicts wrongly</i>				
7	I was trying to plan my wedding by getting a caterer, and they kept blowing us off over and over again.	furious	disappointed	disappointed
8	Being a successful single mothr.	proud	content	content
9	We were over at our friend's house for a dinner and I was in the kitchen helping her cook. I had melted butter in a baking dish to make dessert, and I poured cold milk into it like the recipe said to do. It ended up cracking the dish. I felt bad. I offered to buy her a new one.	guilty	caring	ashamed
10	One time I had done really well in a class. I fully expected to get an A in it	anticipating	confident	disappointed

Table 8: Example Predictions on the Empathetic-Intent (E-Intent) task by ConveRT and our DMI\_Base model.

