

Measuring and Improving Semantic Diversity of Dialogue Generation

Anonymous ACL submission

Abstract

Response diversity has become an important criterion for evaluating the quality of open-domain dialogue generation models. However, current evaluation metrics for response diversity do not capture semantic diversity of generated responses, as they only consider lexical aspects of the responses. In this paper, we introduce a new automatic evaluation metric to measure the semantic diversity of generated responses. Through human evaluation, we demonstrate that our proposed metric highly correlates to human judgments on response diversity than existing lexical-level diversity metrics. Furthermore, motivated by the analysis of an existing dialogue dataset, we propose a simple yet effective learning method that improves the semantic diversity of generated responses through response re-weighting based on the semantic distribution of the training dataset. Through automatic and human evaluation, we show that our proposed learning method better improves both response diversity and coherency compared to other baseline methods.

1 Introduction

Open-domain dialogue generation (Sordani et al., 2015; Bordes et al., 2017) has greatly progressed with the development of large-scale pretrained language models (Radford et al.; Roller et al., 2021) in the last decade. However, although dialogue generation models can produce fluent responses for a given context, they are also known for frequently generating dull and uninformative generic responses (e.g., "I don't know"), degrading the interestingness of responses (Serban et al., 2016; Li et al., 2016a). To alleviate this problem, many studies (Zhao et al., 2017; Li et al., 2017a; Zhang et al., 2018) have been conducted to enhance the diversity of generated responses, and response diversity has become an important criterion for evaluating the

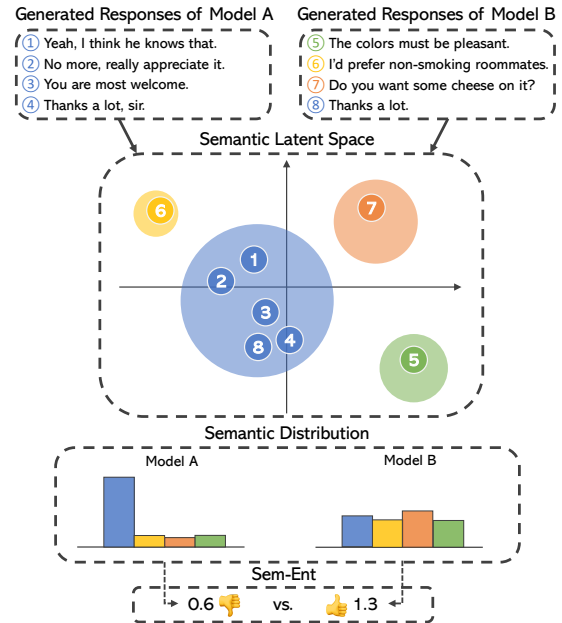


Figure 1: An illustration of measuring semantic diversity of generated responses. Although both Model A and Model B generate lexically diverse responses, we argue that the responses of Model B seem more varied in human perception because they are semantically diverse. Our proposed Sem-Ent measures semantic diversity based on the semantic distribution of generated responses.

quality of generated responses.¹

The current evaluation protocol employs lexical-level evaluation metrics such as *distinct-n* (Dist-*n*) (Li et al., 2016a) and *entropy-n* (Ent-*n*) (Serban et al., 2017; Zhang et al., 2018) to measure the diversity of generated responses. However, it is unclear whether lexical-level evaluation metrics can successfully capture the human judgment on response diversity. For instance, in Figure 1, responses generated by model A and model B both show high lexical diversity, but humans intuitively recognize that the responses of model B

¹According to recent survey papers (Ni et al., 2021; Liang and Li, 2021), more than thirty studies within five years have assessed dialogue generation models from the diversity perspective.

are more diverse. We argue that considering a *semantic diversity* of the generated responses is more important for capturing human judgment on response diversity. However, the lexical-level metrics cannot directly capture the semantic diversity since responses including similar words can have very different semantics, and responses with different words can have similar semantics (Yarats and Lewis, 2018). Nevertheless, most studies have conducted an evaluation with only the lexical-level evaluation metrics to measure the diversity of generated responses because there is no alternative metric to measure the semantic diversity.

To this end, we propose *Sem-Ent* (**Semantic-Entropy**), which is a new automatic evaluation metric for measuring the semantic diversity of generated responses. *Sem-Ent* first maps generated responses into a semantic latent space using a pre-trained language model (e.g., DialoGPT (Zhang et al., 2020) and BERT (Devlin et al., 2019)). Then, the metric measures the semantic diversity of generated responses by measuring how the responses are evenly distributed in the semantic latent space based on semantic clusters, as shown in Figure 1. Through human evaluation, we demonstrate that *Sem-Ent* is more highly correlated with human judgments on response diversity than existing lexical-level evaluation metrics. The human evaluation further shows that *Sem-Ent* highly correlates with human judgments about how they feel generated responses are interesting.

Furthermore, we observe that the semantic distribution of responses in the dialogue dataset is highly imbalanced. This imbalance leads the model to produce semantically less diverse responses. To address this problem, we propose a simple yet effective learning method of dialogue generation models. Our proposed method, *DRESS* (Diversifying RESponses Semantically), induces dialogue generation models to learn more about responses with rare semantics and learn less about responses with frequent semantics. From this, dialogue generation models could produce more semantically diverse responses. Experiments on two benchmark datasets demonstrate that *DRESS* shows substantially better semantic diversity compared to state-of-the-art baseline methods, along with the gain in response coherency. Interestingly, *DRESS* achieves better performance in evaluation metrics for lexical-level diversity than baselines even though it focuses on improving the semantic diversity of generated re-

sponses. Moreover, human evaluation results also affirm the effectiveness of *DRESS*, where *DRESS* outperforms all baseline methods in terms of appropriateness and informativeness of generated responses.

Our Contributions: (1) A new automatic evaluation metric for measuring semantic diversity (*Sem-Ent*), which is highly correlated with human judgment on response diversity. (2) A simple yet effective learning method of dialogue generation models (*DRESS*) for improving the semantic diversity of generated responses. (3) Experiments on two benchmark datasets, showing that *DRESS* outperforms the baseline methods in both semantic diversity and lexical-level diversity. (4) A Python library² of *Sem-Ent*, contributing to the community of open-domain dialogue generation.

2 Related Work

2.1 Open-domain Dialogue Models for Enhancing Response Diversity

Since generating dull and uninformative responses is a well-known and essential problem in open-domain dialogue (Vinyals and Le, 2015; Li et al., 2016a), numerous lines of works have been proposed to address this issue. Li et al. (2016a) replace the standard maximum likelihood objective into maximum mutual information objective to penalize generic responses. This new objective function has been continuously adopted in subsequent works to increase the specificity and diversity of generated responses (Li et al., 2016c; Zhang et al., 2018, 2020). Another line of work improves diversity by modeling the one-to-many relationship of open-domain dialogue using latent variables to generate multiple and diverse responses (Serban et al., 2017; Zhao et al., 2017; Bao et al., 2020a,b; Chen et al., 2019; Zhang et al., 2019; Gao et al., 2019). Some methods selectively penalize frequent responses by removing them from the training dataset (Csáky et al., 2019) or applying negative training to frequent responses (He and Glass, 2020). Using different decoding algorithms can improve the response diversity; Li et al. (2016b) and Vijayakumar et al. (2018) directly modify the beam search algorithm to promote the response diversity. Sampling-based decoding algorithms such as top-k sampling (Fan et al., 2018) and nucleus sampling (Holtzman et al., 2019) are known to improve the diversity of generated responses. Wang et al. (2021) diversify re-

²Link will be released after publication.

sponses by adaptively modifying the target token distribution with a lightweight decoder to prevent the model from being over-confident.

2.2 Metrics for Capturing Response Diversity

Response diversity metrics for open-domain dialogue generation models can mainly be categorized into two groups. Referenced metrics (Zhao et al., 2017; Gao et al., 2019) use the reference responses provided by human annotators to capture the response diversity by computing a recall value based on various similarity metrics such as BLEU and embedding similarity. On the other hand, unreferenced metrics measure the response diversity without the use of reference responses generated by human annotators. Therefore, unreferenced metrics are more widely adopted than referenced metrics because they can measure response diversity even in the absence of reference responses. Dist- n (Li et al., 2016a) measures the response diversity with the fraction of distinct n -grams over possible n -grams in all generated responses. Ent- n metric (Serban et al., 2017; Zhang et al., 2018) is suggested to improve the Dist- n metric by taking the frequency difference of n -grams into account. LF (Li et al., 2019) calculates the frequency of low-frequency words in generated responses as the response diversity. Our work focuses on introducing a semantic diversity metric that alleviates the limitation of the aforementioned unreferenced diversity metrics of considering only lexical aspect of generated responses.

3 Measuring Semantic Diversity

3.1 Sem-Ent

Let $\mathcal{D} = \{(c_1, r_1), (c_2, r_2), \dots, (c_m, r_m)\}$ denote a training dataset consisting of m dialogues where c_i and r_i denote the context and its response of the i -th dialogue, respectively. Dialogue generation is to generate a response r for a given context c .

We are motivated by recent empirical observations that responses can be clustered by the semantic similarity between the responses (Ko et al., 2020; Gao et al., 2020). By following Csáky et al. (2019); Pillutla et al. (2021), we cluster responses in \mathcal{D} by utilizing a pretrained language model. Here, we select DialoGPT (Zhang et al., 2020) as the language model. Each response r_i in \mathcal{D} is turned into a semantic representation $e(r_i)$ by the language model, and then k semantic clusters are formed from the semantic representations by the

k -means algorithm (Lloyd, 1982). Let \mathcal{C} denote a set of the obtained k semantic clusters.

Consider a test dataset $\tilde{\mathcal{D}} = \{(\tilde{c}_1, \tilde{r}_1), \dots, (\tilde{c}_n, \tilde{r}_n)\}$ consisting of n dialogues. During evaluation, a dialogue generation model M generates responses $\mathcal{R}^M = \{r_1^M, \dots, r_n^M\}$ for the contexts $\{\tilde{c}_1, \dots, \tilde{c}_n\}$ in $\tilde{\mathcal{D}}$, respectively. Sem-Ent measures the semantic diversity of \mathcal{R}^M generated by the model M . To compute Sem-Ent, we require a semantic distribution $P(\mathcal{R}^M)$, but there is no direct way to obtain the exact distribution. Thus, we approximate the semantic distribution $P(\mathcal{R}^M)$ using a distribution $\tilde{P}(\mathcal{C}) = [\tilde{p}(1); \dots; \tilde{p}(k)]$ of the semantic clusters \mathcal{C} as follows:

$$\tilde{p}(j) = \frac{1}{n} \sum_{i=1}^n \mathbb{I}(\phi_{\mathcal{C}}(e(r_i^M)) = j), \quad (1)$$

where $\phi_{\mathcal{C}}(x) \in \{1, \dots, k\}$ is a cluster mapping function that returns the cluster id of x from \mathcal{C} . $\tilde{p}(j)$ is the probability of the j -th cluster, indicating how many generated responses are assigned to the j -th semantic cluster.

Sem-Ent is an entropy of $\tilde{P}(\mathcal{C})$, which is calculated with $\tilde{P}(\mathcal{C})$ approximating the semantic distribution of \mathcal{R}^M as follows:

$$\text{Sem-Ent}(\mathcal{R}^M) = - \sum_{j=1}^k \tilde{p}(j) \cdot \log \tilde{p}(j). \quad (2)$$

Interpretation of Sem-Ent is quite straightforward: Sem-Ent gets lower when the semantic distribution gets more imbalanced, i.e., when models generate responses belonging to only several specific semantic clusters. Conversely, Sem-Ent gets the highest value of $\log k$ when generated responses are uniformly distributed to each semantic cluster.

3.2 Correlation with Human judgment

We conduct a human evaluation to demonstrate that Sem-Ent successfully captures human judgments on response diversity.

Experimental Setup. We use a similar experimental setup to that of Pillutla et al. (2021) for analyzing the correlation between response diversity metrics and human judgment. We prepare eight inference settings from two generation models (Blender-90M (Roller et al., 2021) and BART-large (Lewis et al., 2020)) and four decoding algorithms (greedy, beam, top- k sampling, and nucleus sampling). The generation models are fine-tuned on DailyDialog (Li et al., 2017b) dataset that consists of daily conversations about various topics.

Metric	Correlation	Dist-3	Ent-3	LF	MAUVE	Sem-Ent
Diversity/BT	Pearson	0.348 (0.399)	0.702 (0.052)	-0.232 (0.580)	0.134 (0.750)	0.810 (0.015)
	Spearman	0.381 (0.352)	0.667 (0.071)	0.000 (1.000)	0.547 (0.160)	0.762 (0.028)
Interesting/BT	Pearson	0.261 (0.533)	0.671 (0.068)	-0.260 (0.533)	0.098 (0.817)	0.789 (0.020)
	Spearman	0.381 (0.352)	0.714 (0.047)	0.048 (0.911)	0.523 (0.182)	0.667 (0.020)

Table 1: Correlation of various diversity measures with human judgments. "BT" denotes the Bradley-Terry score for a pair-wise human evaluation and the value inside the parenthesis indicates p-value.

Then, each inference setting is paired with other settings, which gives a total of 28 ($8C_2$) pairs of settings.

For every round, a human annotator is assigned to a set that contains ten independent contexts c_1, \dots, c_{10} and two sides of responses r_{1a}, \dots, r_{10a} and r_{1b}, \dots, r_{10b} generated with two different settings. The annotator is asked to select which side of responses is better in two criteria; whether (1) shows more diversity and (2) shows more interesting and creative responses, using a 5-point Likert scale. We obtain 25 preference ratings for each pair of inference settings. These annotation results are converted into each setting’s score by using the Bradley-Terry model (Marden, 1996) fitted by pair-wise annotations. We measure the correlation between the Bradley-Terry score and diversity metrics to check how each metric correlates with the human judgment on each criterion. More details about human evaluation are included in Appendix. **Baseline Metrics.** We compare Sem-Ent with existing lexical-level response diversity metrics: Dist- n (Li et al., 2016a), Ent- n (Serban et al., 2017; Zhang et al., 2018) and LF (Li et al., 2019). We also include recently proposed MAUVE (Pillutla et al., 2021) as a baseline metric. MAUVE shares some properties with Sem-Ent such that it evaluates the distributional property of generated responses with semantic latent representations. However, it is designed to measure the divergence of generated responses from human responses, not for directly measuring response diversity. We compare Sem-Ent to MAUVE to verify that our Sem-Ent is more suitable for measuring the response diversity in open-domain dialogue generation.

Results. Table 1 shows the correlation between the human judgments and the different diversity metrics in terms of Pearson and Spearman rank correlation. Our Sem-Ent shows the highest correlation (on both Pearson correlation and Spearman correlation) with human judgment on response diversity compared to other evaluation metrics with a significant margin. Especially, Dist- n , the most com-

monly used metric for response diversity, shows a much lower correlation (0.348) compared to Sem-Ent (0.810). These results support that Sem-Ent is a good surrogate for measuring human judgment on response diversity and strongly suggest that analyzing the semantic diversity of generated responses is crucial for capturing human perception of response diversity. Moreover, MAUVE shows a lower correlation with human judgment on response diversity. This result implies that a closer gap between human responses and generated responses does not always indicate that generated responses are diverse since human responses contain many dull responses frequently (also studied in Section 4.1 and by Csáky et al. (2019)).

We also observe that Sem-Ent shows a high correlation with human judgment on interestingness; Sem-Ent has a similar correlation to Ent- n and shows a substantially higher correlation than Dist- n , LF, and MAUVE. We believe that the strong correlation of Sem-Ent with human judgment on response diversity leads to a high correlation with a closely related model property, interestingness.

In Section 6, we further justify that Sem-Ent is robust to a choice of configurations used for the metric such as a choice of the language model for extracting semantic representations of responses and a number of clusters k .

4 DRESS: Diversifying RESPONSES Semantically

4.1 Diagnosing the Semantic Distribution of Dialogue Dataset

As shown in Section 3.2, semantic distribution of responses provides a crucial clue for understanding the diversity of the responses. Therefore, we analyze the semantic distribution of the responses in the training dataset. Figure 2 depicts the semantic distribution $\tilde{P}(\mathcal{R})$ of the responses $\mathcal{R} = \{r_1, r_2, \dots, r_m\}$ in the training data of DailyDialog dataset. As shown in the figure, the semantic distribution of the training dataset \mathcal{D} is highly skewed – almost half of the responses fall into the

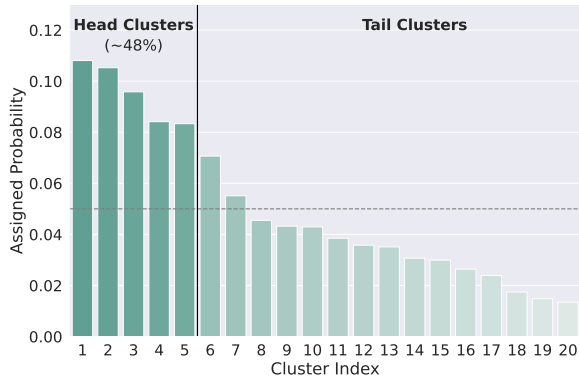


Figure 2: Semantic distribution of the responses in the train split of DailyDialog. Clusters are sorted in the descending order of the assigned probabilities. The dashed line indicates the uniformly distributed probability, 0.05.

Index	Responses
2	<ul style="list-style-type: none"> • Yeah . I know . • Thank you . • You are most welcome . • No more , thank you very much . • Not yet .
13	<ul style="list-style-type: none"> • that sounds great . Do you know if there are any vegetable dishes that are spicy ? • Do you want cheese on it ? • I agree . The colors must be soft and pleasant . You should feel comfortable when you cook our dinners .
18	<ul style="list-style-type: none"> • I bought a new mattress and some fresh bed-clothes . I also bought a new dressing table and a new bedside table . • I ' d prefer non-smoking roommates , but I guess I ' ll have to take what I can get ! • A single room with a front view is 100 dollars per night , one with a rear is 80 dollars .

Table 2: Response examples of the semantic clusters. *Index* column indicates the Cluster Index in Figure 2.

top five frequent clusters (head clusters). Moreover, the frequent clusters tend to contain more generic and dull responses compared to infrequent clusters (tail clusters), as illustrated in Table 2. Contrarily, responses in the infrequent clusters have a wider variety of topics, intents, and diverse vocabularies. Since the training data is skewed towards semantically generic and dull responses, naively training with this data will lead to a low semantic diversity of generated responses.

4.2 Improving Semantic Diversity with DRESS

We introduce a simple yet effective learning method of generation models for improving semantic diversity, DRESS, which addresses the problem of the imbalanced semantic distribution by re-

weighting the instances in the training dataset. The purpose of DRESS is simple: inducing generation models to learn more about responses in the infrequent semantic clusters and contrarily learn less about responses in the frequent semantic clusters. To this end, DRESS modifies the learning objective into the weighted loss function and applies Negative Training (He and Glass, 2020; Li et al., 2020) to the modified objective.

A conventional dialogue generation model is trained by optimizing an NLL (negative log-likelihood) objective as follows:

$$L_{NLL}(D) = - \sum_{i=1}^m \log p_{\theta}(r_i|c_i), \quad (3)$$

where θ indicates parameters of dialogue generation models. Instead of using vanilla NLL objective, we propose to utilize weighted NLL objective in DRESS using weight of responses $w(r_i)$:

$$L_{DRESS}(D) = - \sum_{i=1}^m w(r_i) \cdot \log p_{\theta}(r_i|c_i). \quad (4)$$

The goal of weighted NLL objective is to assign smaller weights to the responses in frequent semantic clusters and assign bigger weights to responses in infrequent semantic clusters to balance the semantic distribution. To meet this condition, the weighting function $w(r)$ should satisfy the constraint: if $\tilde{p}(\phi_c(e(r_i))) \leq \tilde{p}(\phi_c(e(r_j)))$, then $w(r_i) \geq w(r_j)$. Inspired by focal loss (Lin et al., 2017) which is used in the long-tail classification problem (Liu et al., 2019b; Hong et al., 2021), we calculate w as follows:

$$w(r) = (1 - \tilde{p}(\phi_c(e(r))))^{\gamma}, \quad (5)$$

where γ is a hyperparameter for controlling a degree of re-weighting (higher γ means more intense re-weighting).

Moreover, to penalize responses in frequent semantic clusters intensively, we utilize Negative Training (He and Glass, 2020; Li et al., 2020) jointly with the weighted objective function. For every epoch, the model generates responses to each given context. If generated responses are included in head clusters (here, the assigned probability of clusters is bigger than 0.1), then those generated responses are assumed as negative examples, i.e., assigning $w(r) = -1$.

5 Experiments

5.1 Experimental Setup

We conduct experiments to demonstrate that the proposed learning method successfully improves response diversity.

Dataset. We conduct experiments on two English open-domain dialogue datasets: DailyDialog and OpenSubtitles (Lison and Tiedemann, 2016). DailyDialog consists of 13K dialogues which includes 87K context-response pairs, and we split the dialogues into train/valid/test sets in 8:1:1. The test set of DailyDialog contains 6.7K context-response pairs. OpenSubtitles is a large corpus containing movie scripts, and we use the version released in 2018 with 100K context-response pairs for the training and validation set each. We get rid of context-response pairs whose response is shorter than five words from the original test set and randomly sample 10K pairs as test data.

Automated Metrics. As the goal of diversity-promoting dialogue generation models is to generate diverse responses without hurting the coherency of responses, we focus on two criteria: response diversity and coherency. For measuring response diversity, we use both lexical-level diversity metrics (Dist- n , Ent- n , and LF) and a semantic diversity metric (Sem-Ent, $k = 20$). For measuring response coherency, we employ MaUde (Sinha et al., 2020), an unreferenced dialogue response evaluation metric that shows a high correlation with human judgments on the fluency of responses.

Human Evaluation. We further conduct a pairwise comparison through the human evaluation for evaluating generated responses since automatic evaluations are sometimes not trustworthy. We use Amazon Mechanical Turk to collect the annotations. Each annotator evaluates which model is better in terms of *Appropriateness* for measuring response coherency and *Informativeness* for evaluating whether the given response has meaningful information relevant to its given context. We collect annotations for 50 test cases per each model pair, and three annotators rate each test case to improve the robustness of the evaluation result. More details about evaluation protocol (e.g., interface for collecting annotation) are shown in Appendix.

5.2 Baseline Methods

MMI (Li et al., 2016a) increases response diversity by maximizing the mutual information between context and response rather than maximizing the

likelihood as in conventional dialogue models. We utilize the MMI-antiLM as our MMI baseline.

CVAE (Zhao et al., 2017) is a representative model among dialogue generation models that utilize latent variables to increase response diversity. CVAE builds the response generation process as a conditional variational auto-encoder of a response with dialogue context as a condition.

EDF (Entropy-based Data Filtering) (Csáky et al., 2019) enhances response diversity by filtering out context-response pairs that increase one-to-many or many-to-one problems in the training dataset. We use target side entropy to filter the pairs.

NT (Negative Training) (He and Glass, 2020) directly penalizes the generation of generic responses by applying reverse direction gradient for the losses of the generic responses, leading to maximizing the loss rather than minimizing it.

AdaLabel (Wang et al., 2021) alleviates the overconfidence problem of generation models to improve response diversity by dynamically smoothing the target token distribution with an auxiliary lightweight decoder.

5.3 Implementation Details

We take two Transformer-based sequence-to-sequence models: Blender-90M (Roller et al., 2021) and BART-large (Lewis et al., 2020) as the underlying generation models to demonstrate that our method widely works well on different architectures. For DRESS, we set $\gamma = 30$ and the number of clusters $k = 20$ in our whole experiments unless otherwise specified. All models use greedy decoding strategy, and we utilize both blocking repeated n -grams (Paulus et al., 2017) ($n = 3$) within the generated response and the input sequence to prevent models from repeating subsequences. Moreover, we release our implementation code³ publicly to help researchers reproduce the result.

6 Results and Analysis

6.1 Evaluation Results

Table 3 shows the automatic evaluation results. Overall, DRESS achieves the best performance in both semantic and lexical-level response diversity while showing high response fluency for most of the experimental setups. To be more specific, as shown in the table, DRESS shows a substantially higher semantic diversity (Sem-Ent) than all other baseline models in every experimental setup.

³Link will be released after publication.

Backbone	Method	Dist-1	Dist-2	Dist-3	Ent-1	Ent-2	Ent-3	LF	MaUdE	Sem-Ent
Blender-90M (DailyDialog)	Vanilla	0.0453	0.2103	0.3881	7.1322	10.7502	12.3950	0.2234	0.8489	2.5486
	MMI	0.0349	0.1677	0.3069	7.0730	10.3806	11.9808	0.2155	0.8208	2.5784
	CVAE	0.0471	<u>0.2389</u>	<u>0.4459</u>	7.4074	11.2797	12.9969	0.2449	0.8552	2.6261
	EDF	<u>0.0473</u>	0.2271	0.4226	7.2888	11.0283	12.7132	0.2402	<u>0.8593</u>	2.5872
	NT	0.0475	0.2351	0.4422	7.3994	11.2561	13.0111	0.2467	0.8597	2.6434
	AdaLabel	0.0377	0.1982	0.3915	7.1546	10.8772	12.6829	0.2158	0.8443	2.6038
	DRESS(-NT)	0.0445	0.2295	0.4360	7.4560	<u>11.3273</u>	<u>13.1028</u>	<u>0.2474</u>	0.8460	<u>2.7576</u>
DRESS	0.0460	0.2404	0.4571	7.5468	11.5094	13.3060	0.2576	0.8575	2.7819	
BART-large (DailyDialog)	Vanilla	0.0462	0.2168	0.4056	7.3913	11.2075	12.8648	0.2593	0.8854	2.4251
	MMI	0.0497	0.2329	0.4355	7.4748	11.4060	13.0898	0.2623	0.8787	2.4646
	CVAE	0.0429	0.2416	0.5117	7.2728	11.2968	13.1643	0.2558	0.8744	2.4215
	EDF	0.0597	0.2926	0.5355	7.9606	12.1776	13.8786	0.3036	0.8918	2.5842
	NT	<u>0.0571</u>	<u>0.2919</u>	0.5424	8.0267	12.3098	14.0577	<u>0.3070</u>	0.9024	2.6690
	AdaLabel	0.0482	0.2573	0.5136	7.9152	12.0968	13.9496	0.2936	0.8947	2.6336
	DRESS(-NT)	0.0554	0.2909	<u>0.5448</u>	8.1722	<u>12.5195</u>	<u>14.3244</u>	0.3079	0.9192	<u>2.8444</u>
DRESS	0.0547	0.2906	0.5504	8.1821	12.5533	14.3890	0.3052	<u>0.9153</u>	2.8548	
Blender-90M (OpenSubtitles)	Vanilla	0.0373	0.1550	0.2698	6.5882	9.5097	10.7983	0.1758	0.8459	2.4702
	MMI	0.0426	0.1660	0.2755	6.4854	9.2276	10.3364	0.2005	0.8721	2.4469
	CVAE	0.0393	0.1804	0.3398	7.0092	10.5135	11.8959	0.2073	0.9214	2.5726
	EDF	0.0476	0.2019	0.3536	7.0189	10.3899	11.8036	0.2161	0.8777	2.5738
	NT	<u>0.0504</u>	<u>0.2216</u>	<u>0.3969</u>	<u>7.3734</u>	<u>11.0928</u>	<u>12.6594</u>	<u>0.2480</u>	0.8944	2.7049
	AdaLabel	0.0431	0.1913	0.3573	7.0306	10.5280	12.0680	0.2063	0.8708	2.6407
	DRESS(-NT)	0.0499	0.2178	0.3817	7.3316	10.8422	12.2530	0.2308	0.8927	<u>2.7114</u>
DRESS	0.0524	0.2351	0.4180	7.5113	11.2355	12.7612	0.2612	<u>0.9041</u>	2.7654	
BART-large (OpenSubtitles)	Vanilla	0.0262	0.1028	0.1806	5.8507	8.2064	9.2760	0.1532	0.7803	2.2043
	MMI	0.0275	0.1094	0.1923	6.0557	8.5303	9.6961	0.1595	0.8067	2.1626
	CVAE	0.0226	0.1460	0.3495	6.2232	9.7304	11.4593	0.1507	0.8600	2.3005
	EDF	0.0474	<u>0.2056</u>	<u>0.3572</u>	7.0338	10.5464	11.9977	0.2209	0.8558	2.5346
	NT	0.0228	0.0948	0.1594	5.5542	8.2025	9.6915	0.1165	0.8298	2.6368
	AdaLabel	0.0381	0.1772	0.3316	7.0306	10.5667	12.0747	0.2030	<u>0.8647</u>	2.5652
	DRESS(-NT)	0.0456	0.2006	0.3509	<u>7.1669</u>	<u>10.6915</u>	<u>12.1509</u>	<u>0.2220</u>	0.8618	<u>2.6620</u>
DRESS	<u>0.0472</u>	0.2178	0.3890	7.4656	11.2761	12.8601	0.2322	0.8873	2.7406	

Table 3: Automatic evaluation results in terms of various diversity metrics (Dist- n , Ent- n , LF, and Sem-Ent) and coherency metric (an average MaUdE of generated responses). **Bolded** value indicates the best result and underlined value indicates the runner-up among the results. DRESS(-NT) indicates the variant version of DRESS that only utilizes the weighted NLL without Negative Training.

Figure 3 illustrates the detailed semantic distribution of the generated responses. While the Vanilla model shows a high probability on the head semantic clusters (e.g., Cluster 1, 2, 4) and low probability on the tail semantic clusters (e.g., Cluster 13~20), DRESS effectively reduces the probabilities of the head semantic cluster and boosts probabilities of the tail clusters. It is quite intriguing that DRESS also achieves better performance in lexical-level response diversity (Dist- n , Ent- n , and LF). Furthermore, MaUdE results indicate that DRESS preserves better response coherency compared to other baseline methods.

Apart from automatic evaluation, we further compare DRESS with baseline methods in pairwise human evaluation to verify the effectiveness of DRESS. Table 4 shows the evaluation results, showing clear improvements in terms of appropri-

ateness and informativeness from using DRESS.

6.2 Analysing DRESS

Changing Hyperparameters of DRESS. We examine how the automatic results change when we vary the hyperparameters of DRESS: γ in Equation 4.2 and the number of clusters k . Table 5 shows the results about the effect of the hyperparameters. We find that increasing γ induces models to produce more diverse responses, which can be shown by improvement in Dist-3, Ent-3, and Sem-Ent. We also observe that decreasing k induces the models to generate more diverse responses. However, MaUdE gets degraded while response diversity improves, which implies a trade-off between response diversity and coherence.

Ablation Study. To verify the effect of our weighted NLL, we conduct an ablation study. In

Comparison (A vs. B)	Appropriateness			Informativeness		
	A wins (%)	B wins (%)	Tie (%)	A wins (%)	B wins (%)	Tie (%)
Ours vs Vanilla	35.3	24.7	40.0	36.0	28.0	36.0
Ours vs MMI	40.0	34.7	25.3	40.7	36.0	23.3
Ours vs CVAE	44.7	30.0	25.3	36.7	36.0	27.3
Ours vs EDF	35.3	24.7	40.0	32.7	23.3	44.0
Ours vs NT	28.0	25.3	46.7	37.3	26.0	36.7
Ours vs AdaLabel	28.7	24.0	47.3	32.7	31.3	36.0

Table 4: Human pair-wise comparison results in terms of appropriateness and informativeness of generated responses. The evaluation is conducted on the test set of DailyDialog with Blender-90M using greedy decoding.

Config	Dist-3	Ent-3	MaUdE	Sem-Ent
$\gamma = 1.0$	0.4333	12.8968	0.8570	2.6233
$\gamma = 5.0$	0.4400	12.9989	0.8593	2.6551
$\gamma = 10.0$	0.4410	13.0670	0.8583	2.6959
$\gamma = 30.0$	0.4571	13.3060	0.8575	2.7819
$\gamma = 100.0$	0.4625	13.5839	0.8436	2.8444
$k = 10$	0.4748	13.7596	0.8390	2.8451
$k = 20$	0.4571	13.3060	0.8575	2.7819
$k = 50$	0.4318	13.0001	0.8513	2.7009
$k = 100$	0.4311	12.8857	0.8637	2.6258

Table 5: Analysing the effect of hyperparameters, γ and k . When changing γ , we fix k to 20. When changing k , we fix γ to 30.0.

Table 3, DRESS(-NT) indicates the variant of DRESS without Negative Training and only utilizes weighted NLL. DRESS(-NT) shows a slight degradation in Sem-Ent compared to DRESS. Nonetheless, DRESS(-NT) achieves better performance in Sem-Ent than other baseline methods excluding DRESS. Moreover, DRESS(-NT) also shows a higher lexical-level diversity than other baseline methods, along with high MaUdE scores.

6.3 Robustness of Sem-Ent on the Choice of Configurations

In this section, we examine the robustness of Sem-Ent changing the configurations used for calculating the metric. Several configurations can be changed in Sem-Ent, including the choice of language models for mapping responses r into a semantic representation $e(r)$ and the number of clusters k for the k -means algorithm. Varying the configurations, we compute Sem-Ent on responses generated by Blender-90M for the test set of DailyDialog with all methods (in Table 3). We then measure the Spearman correlation between the computed Sem-Ent of different configurations.

For the choice of language models, we compare three variants: DialogGPT, RoBERTa (Liu et al., 2019a), and GPT2-large (Radford et al.). The average Spearman correlation between the pairs of

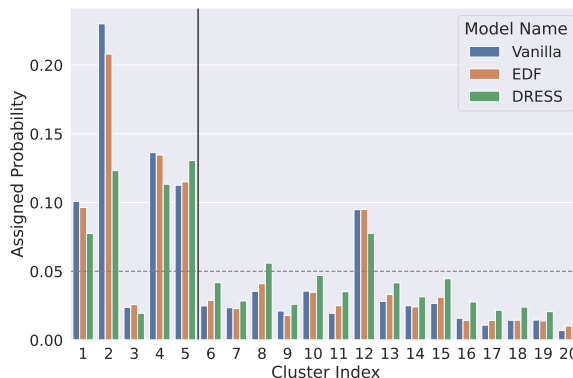


Figure 3: Probability distribution of the responses generated by Vanilla, EDF and DRESS. The dashed line indicates the uniformly distributed probability, 0.05.

these three variants (3 pairs) is 0.8809. For the number of clusters, we vary the number k with values in $\{10, 20, 50, 100\}$ and compare the Sem-Ent rankings. The average Spearman correlation between these configurations (6 pairs) is 0.9821. High correlations show that Sem-Ent produces similar rankings of different models regardless of different configurations, indicating that Sem-Ent is a robust metric for calculating response diversity.

7 Conclusion

In this work, we argue that semantic diversity is overlooked while measuring response diversity of dialogue generation; thus, we present a new automatic evaluation metric, Sem-Ent, which can measure the semantic diversity of generated responses. Sem-Ent correlates with human judgments on response diversity more than other automatic diversity metrics and also shows a high correlation with human judgments in interestingness. Moreover, we introduce a new learning method, DRESS, to improve the semantic diversity of dialogue generation. Evaluation results show that DRESS improves both the semantic diversity and lexical-level diversity of dialogue generation, along with the gain in response coherency.

Ethical Considerations

Dialogue generation models can reveal some biases and toxicities from their responses since these models leverage large-scale web-crawled data for pretraining. This is a common consideration for works related to dialogue generation. Moreover, while our paper focuses on diversifying responses in semantic viewpoint, the model may unintentionally learn about offensive words while diversifying responses. We believe it will be meaningful to reduce potential harmful responses considering semantics in future work.

References

- Siqi Bao, Huang He, Fan Wang, Hua Wu, and Haifeng Wang. 2020a. Plato: Pre-trained dialogue generation model with discrete latent variable. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 85–96.
- Siqi Bao, Huang He, Fan Wang, Hua Wu, Haifeng Wang, Wenquan Wu, Zhen Guo, Zhibin Liu, and Xinchao Xu. 2020b. Plato-2: Towards building an open-domain chatbot via curriculum learning. *arXiv preprint arXiv:2006.16779*.
- Antoine Bordes, Y-Lan Boureau, and Jason Weston. 2017. [Learning end-to-end goal-oriented dialog](#). In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*. OpenReview.net.
- Chaotao Chen, Jinhua Peng, Fan Wang, Jun Xu, and Hua Wu. 2019. Generating multiple diverse responses with multi-mapping and posterior mapping selection. *arXiv preprint arXiv:1906.01781*.
- Richárd Csáky, Patrik Purgai, and Gábor Recski. 2019. Improving neural conversational models with entropy-based data filtering. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 5650–5669.
- 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.
- Angela Fan, Mike Lewis, and Yann Dauphin. 2018. Hierarchical neural story generation. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 889–898.

- Xiang Gao, Sungjin Lee, Yizhe Zhang, Chris Brockett, Michel Galley, Jianfeng Gao, and Bill Dolan. 2019. Jointly optimizing diversity and relevance in neural response generation. In *Proceedings of NAACL-HLT*, pages 1229–1238.
- Yifan Gao, Piji Li, Wei Bi, Xiaojiang Liu, Michael Lyu, and Irwin King. 2020. Dialogue generation on infrequent sentence functions via structured meta-learning. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: Findings*, pages 431–440.
- Seungju Han, Beomsu Kim, Seokjun Seo, Enkhbayar Erdenee, and Buru Chang. 2021. Understanding and improving the exemplar-based generation for open-domain conversation. *arXiv preprint arXiv:2112.06723*.
- Tianxing He and James Glass. 2020. Negative training for neural dialogue response generation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 2044–2058.
- Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. 2019. The curious case of neural text degeneration. In *International Conference on Learning Representations*.
- Youngkyu Hong, Seungju Han, Kwanghee Choi, Seokjun Seo, Beomsu Kim, and Buru Chang. 2021. Disentangling label distribution for long-tailed visual recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6626–6636.
- Beomsu Kim, Seokjun Seo, Seungju Han, Enkhbayar Erdenee, and Buru Chang. 2021. Distilling the knowledge of large-scale generative models into retrieval models for efficient open-domain conversation. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 3357–3373.
- Wei-Jen Ko, Avik Ray, Yilin Shen, and Hongxia Jin. 2020. Generating dialogue responses from a semantic latent space. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 4339–4349.
- Iliia Kulikov, Alexander Miller, Kyunghyun Cho, and Jason Weston. 2019. Importance of search and evaluation strategies in neural dialogue modeling. In *Proceedings of the 12th International Conference on Natural Language Generation*, pages 76–87.
- Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7871–7880.

681	Jiwei Li, Michel Galley, Chris Brockett, Jianfeng Gao, and William B Dolan. 2016a. A diversity-promoting objective function for neural conversation models. In <i>Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 110–119.	Ziwei Liu, Zhongqi Miao, Xiaohang Zhan, Jiayun Wang, Boqing Gong, and Stella X Yu. 2019b. Large-scale long-tailed recognition in an open world. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pages 2537–2546.	736 737 738 739 740
688	Jiwei Li, Will Monroe, and Dan Jurafsky. 2016b. A simple, fast diverse decoding algorithm for neural generation. <i>arXiv preprint arXiv:1611.08562</i> .	Stuart Lloyd. 1982. Least squares quantization in pcm. <i>IEEE transactions on information theory</i> , 28(2):129–137.	741 742 743
689		Edward Loper and Steven Bird. 2002. Nltk: The natural language toolkit. In <i>Proceedings of the ACL-02 Workshop on Effective Tools and Methodologies for Teaching Natural Language Processing and Computational Linguistics</i> , pages 63–70.	744 745 746 747 748
691	Jiwei Li, Will Monroe, Alan Ritter, Dan Jurafsky, Michel Galley, and Jianfeng Gao. 2016c. Deep reinforcement learning for dialogue generation. In <i>Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing</i> , pages 1192–1202.	John I Marden. 1996. <i>Analyzing and Modeling Rank Data</i> . CRC Press.	749 750
696		Alexander Miller, Will Feng, Dhruv Batra, Antoine Bordes, Adam Fisch, Jiasen Lu, Devi Parikh, and Jason Weston. 2017. Parlai: A dialog research software platform. In <i>Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing: System Demonstrations</i> , pages 79–84.	751 752 753 754 755 756
697	Jiwei Li, Will Monroe, Tianlin Shi, Sébastien Jean, Alan Ritter, and Dan Jurafsky. 2017a. Adversarial learning for neural dialogue generation. In <i>Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing</i> , pages 2157–2169.	Jinjie Ni, Tom Young, Vlad Pandelea, Fuzhao Xue, Vinay Adiga, and Erik Cambria. 2021. Recent advances in deep learning based dialogue systems: A systematic survey. <i>arXiv preprint arXiv:2105.04387</i> .	757 758 759 760
702	Margaret Li, Stephen Roller, Iliia Kulikov, Sean Welleck, Y-Lan Boureau, Kyunghyun Cho, and Jason Weston. 2020. Don’t say that! making inconsistent dialogue unlikely with unlikelihood training. In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 4715–4728.	Romain Paulus, Caiming Xiong, and Richard Socher. 2017. A deep reinforced model for abstractive summarization. <i>arXiv preprint arXiv:1705.04304</i> .	761 762 763
703		Krishna Pillutla, Swabha Swayamdipta, Rowan Zellers, John Thickstun, Sean Welleck, Yejin Choi, and Zaid Harchaoui. 2021. Mauve: Measuring the gap between neural text and human text using divergence frontiers. <i>Advances in Neural Information Processing Systems</i> , 34.	764 765 766 767 768 769
704		Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners.	770 771 772
705		Stephen Roller, Emily Dinan, Naman Goyal, Da Ju, Mary Williamson, Yinhan Liu, Jing Xu, Myle Ott, Eric Michael Smith, Y-Lan Boureau, et al. 2021. Recipes for building an open-domain chatbot. In <i>Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume</i> , pages 300–325.	773 774 775 776 777 778 779
706		Iulian Serban, Alessandro Sordoni, Yoshua Bengio, Aaron Courville, and Joelle Pineau. 2016. Building end-to-end dialogue systems using generative hierarchical neural network models. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 30.	780 781 782 783 784 785
707		Iulian Serban, Alessandro Sordoni, Ryan Lowe, Laurent Charlin, Joelle Pineau, Aaron Courville, and Yoshua Bengio. 2017. A hierarchical latent variable encoder-decoder model for generating dialogues. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 31.	786 787 788 789 790 791
708	Yanran Li, Hui Su, Xiaoyu Shen, Wenjie Li, Ziqiang Cao, and Shuzi Niu. 2017b. Dailydialog: A manually labelled multi-turn dialogue dataset. In <i>Proceedings of the Eighth International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 986–995.		
709			
710			
711			
712			
713			
714	Zuchao Li, Rui Wang, Kehai Chen, Masso Utiyama, Eiichiro Sumita, Zhuosheng Zhang, and Hai Zhao. 2019. Data-dependent gaussian prior objective for language generation. In <i>International Conference on Learning Representations</i> .		
715			
716			
717			
718			
719	Hongru Liang and Huaqing Li. 2021. Towards standard criteria for human evaluation of chatbots: A survey. <i>arXiv preprint arXiv:2105.11197</i> .		
720			
721			
722	Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. 2017. Focal loss for dense object detection. In <i>Proceedings of the IEEE international conference on computer vision</i> , pages 2980–2988.		
723			
724			
725			
726	Pierre Lison and Jörg Tiedemann. 2016. Opensubtitles2016: Extracting large parallel corpora from movie and tv subtitles. In <i>Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC’16)</i> , pages 923–929.		
727			
728			
729			
730			
731	Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019a. Roberta: A robustly optimized bert pretraining approach. <i>arXiv preprint arXiv:1907.11692</i> .		
732			
733			
734			
735			

792	Koustuv Sinha, Prasanna Parthasarathi, Jasmine Wang,	Yuchi Zhang, Yongliang Wang, Liping Zhang, Zhiqiang	847
793	Ryan Lowe, William L Hamilton, and Joelle Pineau.	Zhang, and Kun Gai. 2019. Improve diverse text	848
794	2020. Learning an unreferenced metric for online	generation by self labeling conditional variational	849
795	dialogue evaluation. In <i>Proceedings of the 58th An-</i>	auto encoder. In <i>ICASSP 2019-2019 IEEE Interna-</i>	850
796	<i>tional Meeting of the Association for Computational</i>	<i>tional Conference on Acoustics, Speech and Signal</i>	851
797	<i>Linguistics</i> , pages 2430–2441.	<i>Processing (ICASSP)</i> , pages 2767–2771. IEEE.	852
798	Alessandro Sordoni, Michel Galley, Michael Auli, Chris	Tiancheng Zhao, Ran Zhao, and Maxine Eskenazi. 2017.	853
799	Brockett, Yangfeng Ji, Margaret Mitchell, Jian-Yun	Learning discourse-level diversity for neural dialog	854
800	Nie, Jianfeng Gao, and William B Dolan. 2015. A	models using conditional variational autoencoders.	855
801	neural network approach to context-sensitive genera-	In <i>Proceedings of the 55th Annual Meeting of the</i>	856
802	tion of conversational responses. In <i>Proceedings of</i>	<i>Association for Computational Linguistics (Volume</i>	857
803	<i>the 2015 Conference of the North American Chap-</i>	<i>1: Long Papers)</i> , pages 654–664.	858
804	<i>ter of the Association for Computational Linguistics:</i>		
805	<i>Human Language Technologies</i> , pages 196–205.		
806	Ashwin K Vijayakumar, Michael Cogswell, Ram-		
807	prasaath R Selvaraju, Qing Sun, Stefan Lee, David		
808	Crandall, and Dhruv Batra. 2018. Diverse beam		
809	search for improved description of complex scenes.		
810	In <i>Thirty-Second AAAI Conference on Artificial In-</i>		
811	<i>telligence</i> .		
812	Oriol Vinyals and Quoc V Le. 2015. A neural conversa-		
813	tional model. <i>arXiv preprint arXiv:1506.05869</i> .		
814	Yida Wang, Yinhe Zheng, Yong Jiang, and Minlie		
815	Huang. 2021. Diversifying dialog generation via		
816	adaptive label smoothing. <i>arXiv e-prints</i> , pages		
817	arXiv–2105.		
818	Thomas Wolf, Julien Chaumond, Lysandre Debut, Vic-		
819	tor Sanh, Clement Delangue, Anthony Moi, Pier-		
820	ric Cistac, Morgan Funtowicz, Joe Davison, Sam		
821	Shleifer, et al. 2020. Transformers: State-of-the-		
822	art natural language processing. In <i>Proceedings of</i>		
823	<i>the 2020 Conference on Empirical Methods in Nat-</i>		
824	<i>ural Language Processing: System Demonstrations</i> ,		
825	pages 38–45.		
826	Ka Wong, Praveen Paritosh, and Lora Aroyo. 2021.		
827	Cross-replication reliability—an empirical approach		
828	to interpreting inter-rater reliability. <i>arXiv preprint</i>		
829	<i>arXiv:2106.07393</i> .		
830	Denis Yarats and Mike Lewis. 2018. Hierarchical text		
831	generation and planning for strategic dialogue. In <i>In-</i>		
832	<i>ternational Conference on Machine Learning</i> , pages		
833	5591–5599. PMLR.		
834	Yizhe Zhang, Michel Galley, Jianfeng Gao, Zhe Gan,		
835	Xiujun Li, Chris Brockett, and Bill Dolan. 2018.		
836	Generating informative and diverse conversational		
837	responses via adversarial information maximization.		
838	<i>Advances in Neural Information Processing Systems</i> ,		
839	31:1810–1820.		
840	Yizhe Zhang, Siqi Sun, Michel Galley, Yen-Chun Chen,		
841	Chris Brockett, Xiang Gao, Jianfeng Gao, Jingjing		
842	Liu, and William B Dolan. 2020. Dialogpt: Large-		
843	scale generative pre-training for conversational re-		
844	sponse generation. In <i>Proceedings of the 58th An-</i>		
845	<i>tional Meeting of the Association for Computational</i>		
846	<i>Linguistics: System Demonstrations</i> , pages 270–278.		

A Appendix

A.1 Descriptive Statistics about Results

Confidence Interval of MaUdE In Table 3 and Table 5 from the main paper, we report the average MaUdE score of responses generated by each method. To provide descriptive statistics of evaluation, here we provide a 95% confidence interval of MaUdE in Table 6 and Table 7. Note that we only report confidence intervals of MaUdE since other diversity metrics (Dist- n , Ent- n , LF, Sem-Ent) return a single value from a set of responses, thus can not calculate the confidence interval.

Inter-Rater Reliability of Pairwise Human Evaluation We calculate a Fleiss’ Kappa for pairwise human evaluation results to measure the annotation variance. We find that Fleiss’ Kappas are 0.09 and 0.04 for appropriateness and informativeness, respectively. Although these values are not high, as Kulikov et al. (2019) and Wong et al. (2021) show on their paper, inter-rater reliability of annotation results using crowd-sourced annotators (such as our case, using Amazon Mechanical Turk) can be low since annotators show high cultural and training variances, especially when the task is subjective as our case. Note that 64 annotators participated in our human evaluation. Also, we limited the number of maximum annotations that a single annotator can be assigned to reduce the bias, which might have increased inter-rater diversity.

A.2 Human Evaluation Protocol

Evaluation for Comparing Metrics We use Amazon Mechanical Turk for collecting assessments, and Figure 4 shows the instructions and the interface for the human evaluation. We mitigate the bias from the annotator by setting a maximum number of annotations per worker as 20 and randomly shuffling the order of the model and the corresponding response. Since our task does not require particular expertise in linguistics, we open the evaluations to non-experts. Nonetheless, to control the annotation quality, we only allow the annotators who satisfy the following requirements: (1) HITs approval rate greater than 95%, (2) Location is one of Australia, Canada, New Zealand, United Kingdom, and the United States, (3) Lifetime number of HITs approved greater than 1000, following Kim et al. (2021); Han et al. (2021). We estimated that each HITs takes around 1.5 minutes on average (87 seconds per each HIT estimated by the 85th percentile of response times) and set the payment to USD 16

per hour. Therefore, annotators are paid USD 0.40 per HITs.

Evaluation for Comparing Methods As we described above, we also use Amazon Mechanical Turk, and we use the same setting to mitigate the bias and control the annotation quality. Figure 5 shows the instructions and the interface for the human evaluation. Here, annotators are paid USD 0.25 per HITs as we estimated that each HITs takes around 1.4 minutes on average (84 seconds per HITs estimated by the 85th percentile of response times) and set the payment to USD 10.7 per hour since the difficulty of the task is easier than above.

A.3 Evaluation Details

Bradley-Terry Model We use the Bradley-Terry model from pairwise human evaluation results to obtain the ranking of the models. Given parameters $\theta_1, \dots, \theta_n$, for two items i and j , the probability of the outcome $i \succ j$ is $p(i \succ j) = e^{\theta_i} / (e^{\theta_i} + e^{\theta_j})$. For more details about the Bradley-Terry model, please refer to [choix manual](#).

Calculating Dist- n , Ent- n , LF We use *NLTK* (Loper and Bird, 2002) package while calculating Dist- n , Ent- n , and LF, particularly for tokenizing sentence and preparing n -grams. When calculating Low-Frequency Token Ratio (LF), we choose words with an occurrence count less than 100 in each dataset.

Number of Experiments We run an experiment only once since our evaluation requires a human evaluation which requires an extra annotation budget.

A.4 Additional Examples of the Semantic Clusters

We provide additional response examples of the semantic clusters in DailyDialog dataset in Table 8.

A.5 Analysis of the Distribution of Generated Responses

Figure 6 illustrates the cumulative semantic probability distributions of the generated responses. DRESS clearly shows the most similar cumulative distribution to that of uniform distribution, which is a distribution that achieves the highest Sem-Ent value. Moreover, DRESS dramatically reduces the distribution of head clusters containing generic responses compared to other baseline methods and conversely enlarges the distribution of tail clusters.

Instructions

Task Info:
 We are studying how good AI models are at generating text on the internet. You are given a multiple dialogue contexts for each model, as well as and two responses from model A and B. These responses are written by an AI. You must choose (a) which of two responses are more diverse, (b) which of two responses is more interesting

Guidelines:

- There are five choices for each question: Definitely A/B, Slightly A/B, or Tie. Please use the "Tie" option extremely sparingly! (No more than one in every ten pairs should be chosen as a tie along any of the three questions).
- The questions can have different answers! Some text is very creative or interesting, but it doesn't quite fit the context or make sense.
- Try to focus on quality over quantity. The text can be long but contain rambly gibberish.
- Please do your best, some of these are pretty challenging!
- Answering each question should take around 1.5 minutes on average, as per our estimation. We have calibrated the pay to be \$16 per hour with this speed.

Example responses that Model A generates:

- Context: \${context_a_0}
- Response: \${resp_a_0}

- Context: \${context_a_1}
- Response: \${resp_a_1}

- Context: \${context_a_2}
- Response: \${resp_a_2}

- Context: \${context_a_3}
- Response: \${resp_a_3}

- Context: \${context_a_4}
- Response: \${resp_a_4}

- Context: \${context_a_5}
- Response: \${resp_a_5}

- Context: \${context_a_6}
- Response: \${resp_a_6}

- Context: \${context_a_7}
- Response: \${resp_a_7}

- Context: \${context_a_8}
- Response: \${resp_a_8}

- Context: \${context_a_9}
- Response: \${resp_a_9}

Example responses that Model B generates:

- Context: \${context_b_0}
- Response: \${resp_b_0}

- Context: \${context_b_1}
- Response: \${resp_b_1}

- Context: \${context_b_2}
- Response: \${resp_b_2}

- Context: \${context_b_3}
- Response: \${resp_b_3}

- Context: \${context_b_4}
- Response: \${resp_b_4}

- Context: \${context_b_5}
- Response: \${resp_b_5}

- Context: \${context_b_6}
- Response: \${resp_b_6}

- Context: \${context_b_7}
- Response: \${resp_b_7}

- Context: \${context_b_8}
- Response: \${resp_b_8}

- Context: \${context_b_9}
- Response: \${resp_b_9}

Which model generates more diverse responses, given the context?

(select one)
▼

Which model generates more interesting and creative, given the context?

(select one)
▼

Figure 4: The interface of human evaluation for assessing how responses are (a) diverse, (b) more interesting and creative.

956 A.6 Limitations of our Work

957 In this section, we discuss the potential limitations
 958 of our methods and the experimental procedure. To
 959 start with, our proposed diversity metric Sem-Ent
 960 requires a pre-trained language model to calculate
 961 the result. This indicates that it requires relatively
 962 heavier computational resources to calculate Sem-
 963 Ent compared to other lexical-based diversity met-
 964 rics such as Dist-*n* and Ent-*n*. Moreover, extend-
 965 ing Sem-Ent to other languages or other domains
 966 could be problematic if no high-quality pre-trained
 967 language model is available on that language or
 968 domain.

969 In terms of the experimental procedure, we per-

970 formed the experiment once rather than running it
 971 multiple times with different seeds. Since our eval-
 972 uation process incorporates a human annotation,
 973 which requires a payment to human annotators, we
 974 were not able to perform multiple sets of experi-
 975 ments due to the limitation on budget. In the same
 976 perspective, we were not able to obtain a sufficient
 977 number of annotations to acquire statistically sig-
 978 nificant results for every pairwise comparison. We
 979 run an experiment only once since our evaluation
 980 requires a human evaluation which requires an ex-
 981 tra annotation budget. Also, we only experimented
 982 with the English dialogue dataset (DailyDialog and
 983 English portion of the OpenSubtitles). Therefore
 984 our results do not necessarily guarantee the same

Instructions

Given the dialogue context, you need to choose a better response between two responses, A and B, with the criteria of **appropriateness** and **informativeness**.

Appropriateness is a metric for evaluating whether **the given response is fluent, logical, and appropriate to its given context**.

Informativeness is a metric for evaluating whether **the given response has meaningful information relevant to its given context**.

Dialogue #1
 Sure . It's probably a skiing show . Adam told me there are many skiing clubs whose members put on shows to attract more visitors .

Candidate #1
 A: yes , i know that the olympic games are often held in places where people ski regularly .
 vs.
 B: i ' ve never heard of such a thing . what are they ?

Select a response with better appropriateness. (Fluent, logical and appropriate to given context)

(select one) ▼

Select a responses with better informativeness. (Meaningful, Specific to given context, Informative)

(select one) ▼

Figure 5: The interface of pairwise human evaluation for appropriateness and informativeness.

985 result in other languages rather than English.

986 Also, we’d like to clarify that our proposed met- 1004
 987 ric, Sem-Ent, only focuses on measuring generated 1005
 988 responses’ diversity and does not consider the re- 1006
 989 sponse coherency. Although this is our intention 1007
 990 since we aimed to build an unreferenced diversity 1008
 991 metric, this limitation yields a drawback that Sem- 1009
 992 Ent should always be jointly used with another 1010
 993 metric that measures the response coherency (e.g., 1011
 994 MaUdE). Expanding Sem-Ent to consider the co- 1012
 995 herency with an input context will be an intriguing 1013
 996 future direction for our research. 1014

997 **A.7 Further Implementation Details** 1015

998 **Training Models** All of our experiments are done 1016
 999 using the ParlAI (Miller et al., 2017) framework. 1017
 1000 We leverage model weights of Blender-90M and 1018
 1001 BART-large from ParlAI. Blender-90M is pre- 1019
 1002 trained on Reddit corpus, and BART-large is pre- 1020

1003 trained jointly on Wikipedia and Toronto Books. 1004
 Note that Blender-90M has 90M parameters, and 1005
 BART-large consists of 400M parameters. All base- 1006
 lines and DRESS use the initial learning rate of 1007
 $7e - 6$ with Adam optimizer, except CVAE for 1008
 Blender-90M trained on DailyDialog using $2e - 5$, 1009
 MMI for Blender-90M trained on OpenSubtitles 1010
 using $1e - 6$, and CVAE for Blender-90M trained 1011
 on OpenSubtitles using $1e - 5$. We search for the 1012
 appropriate learning rate for those exceptions since 1013
 those exceptions are not stable enough to train the 1014
 model. We use a learning rate scheduler that re- 1015
 duces its learning rate by multiplying 0.5 when 1016
 the loss has stopped decreasing. All Blender-90M 1017
 models and all BART-large models are trained us- 1018
 ing batch size of 32 and 16 on single A100 GPU, 1019
 respectively. Training a single model takes less 1020
 than a day with these configurations. 1021

Language Model for Calculating Sem-Ent In 1021

Backbone	Method	MaUdE (\pm 95% CI)	
Blender-90M (DailyDialog)	Vanilla	0.8489 \pm 0.005	
	MMI	0.8208 \pm 0.005	
	CVAE	0.8552 \pm 0.005	
	EDF	0.8593 \pm 0.005	
	NT	0.8597 \pm 0.005	
	AdaLabel	0.8443 \pm 0.005	
	DRESS(-NT)	0.8460 \pm 0.005	
	DRESS	0.8575 \pm 0.005	
	BART-large (DailyDialog)	Vanilla	0.8854 \pm 0.005
		MMI	0.8787 \pm 0.005
CVAE		0.8744 \pm 0.005	
EDF		0.8918 \pm 0.004	
NT		0.9024 \pm 0.004	
AdaLabel		0.8947 \pm 0.004	
DRESS(-NT)		0.9192 \pm 0.003	
DRESS		<u>0.9153</u> \pm 0.003	
Blender-90M (OpenSubtitles)		Vanilla	0.8459 \pm 0.004
		MMI	0.8721 \pm 0.004
	CVAE	0.9214 \pm 0.003	
	EDF	0.8777 \pm 0.004	
	NT	0.8944 \pm 0.003	
	AdaLabel	0.8708 \pm 0.004	
	DRESS(-NT)	0.8927 \pm 0.003	
	DRESS	<u>0.9041</u> \pm 0.003	
	BART-large (OpenSubtitles)	Vanilla	0.7803 \pm 0.005
		MMI	0.8067 \pm 0.005
CVAE		0.8600 \pm 0.004	
EDF		0.8558 \pm 0.004	
NT		0.8298 \pm 0.005	
AdaLabel		<u>0.8647</u> \pm 0.004	
DRESS(-NT)		0.8618 \pm 0.004	
DRESS		0.8873 \pm 0.003	

Table 6: MaUdE with a 95% confidence interval when automatically evaluating various methods.

1022 this work, we test three language models to obtain embeddings from the response: DialoGPT, 1023 RoBERTa, and GPT2-large. For reproducibility, 1024 we utilize model weights which are publicly 1025 opened on HuggingFace Transformers (Wolf 1026 et al., 2020): microsoft/DialoGPT-large, 1027 roberta-base, and gpt2-large for DialoGPT, 1028 RoBERTa, and GPT2-large, respectively. 1029 **Software and Hardware** We use Python 3.8, Py- 1030 Torch 1.9.0 (py3.8_cuda11.1_cudnn8.0.5_0), Hug- 1031 gingFace Transformers 4.6.1, and ParlAI 1.3.0. All 1032 the experiments are done using NVIDIA A100- 1033 40GB GPUs, along with AMD EPYC 7742 64- 1034 Core Processors. 1035 **License** The DailyDialog dataset has CC-BY-NC- 1036 SA 4.0 license. OpenSubtitles dataset does not 1037 specify the license on the dataset. For the pre- 1038 trained models, DialoGPT, RoBERTa, and GPT-2 1039 large is all released with the MIT license. Since 1040 CC-BY-NC-SA 4.0 and MIT license both allow the 1041 utilization of the resource for research purposes, 1042

the use of these scientific artifacts in this work is 1043 valid. 1044

Config	MaUdE (\pm 95% CI)
$\gamma = 1.0$	0.8570 ± 0.004
$\gamma = 5.0$	0.8593 ± 0.004
$\gamma = 10.0$	0.8583 ± 0.004
$\gamma = 30.0$	0.8575 ± 0.004
$\gamma = 100.0$	0.8436 ± 0.004
$k = 10$	0.8390 ± 0.004
$k = 20$	0.8575 ± 0.004
$k = 50$	0.8513 ± 0.004
$k = 100$	0.8637 ± 0.004

Table 7: MaUdE with a 95% confidence interval when analysing the effect of hyperparameters, γ and k .

Index	Responses
1	<ul style="list-style-type: none"> • I'm going to the store . • Oh , yes . Hi , how are you ? • All right . Hop in , please . • I am , sir . • No problem . I 'll wait for your call .
2	<ul style="list-style-type: none"> • Yeah . I know . • Thank you . • You are most welcome . • No more , thank you very much . • Not yet .
18	<ul style="list-style-type: none"> • I bought a new mattress and some fresh bedclothes . I also bought a new dressing table and a new bedside table . • I ' d prefer non-smoking roommates , but I guess I ' ll have to take what I can get ! • A single room with a front view is 100 dollars per night , one with a rear is 80 dollars .
19	<ul style="list-style-type: none"> • Yes . Will you also make copies and file them using both methods ? • you should probably call the IT department and have them check your computer for virus . • I see . Well , can I have a look at your phone ? Unfortunately , this phone can ' t be used in the US . it ' s not compatible with our 3G network .
20	<ul style="list-style-type: none"> • A driver ' s license or something showing that you live in this city . • I want to change a new car . I like Honda best , especially the red one . But it is too expensive . • We use a vacuum cleaner that removes all the dirt , and we throw away all of the trash that we can find .

Table 8: Additional response examples of the semantic clusters of DailyDialog dataset. *Index* column indicates the Cluster Index in Figure 2.

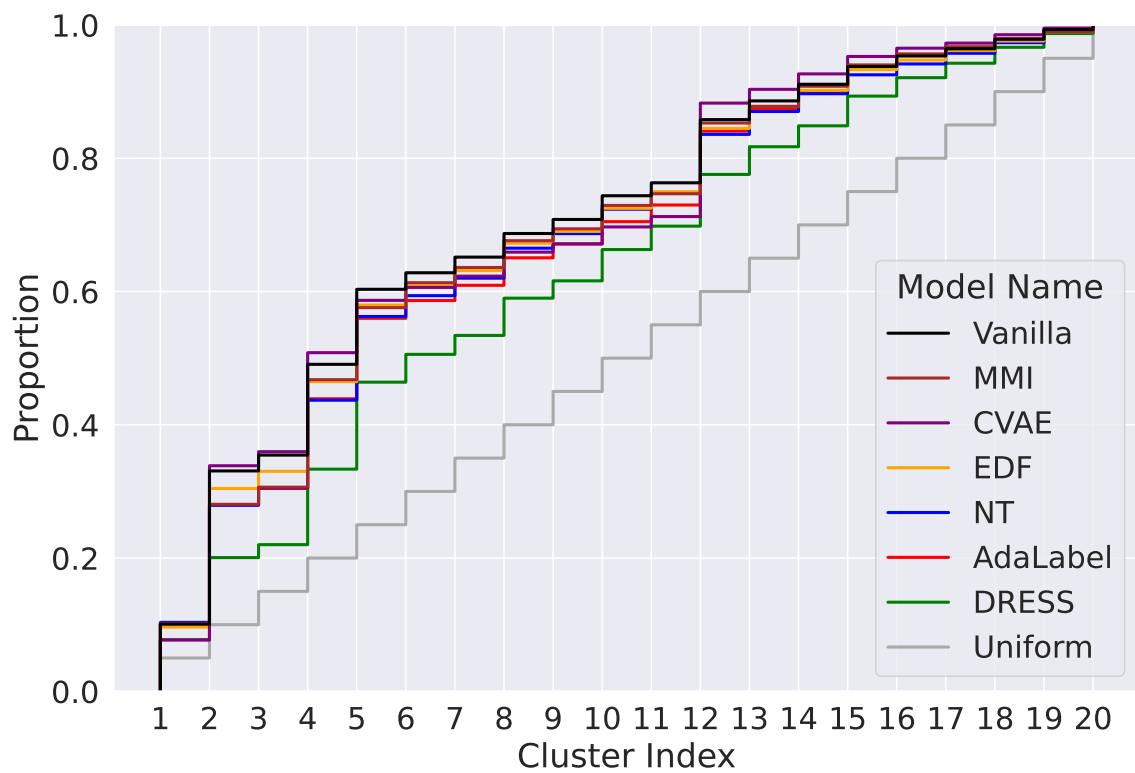


Figure 6: Cumulative probability distribution of the responses generated by different methods. *Uniform* illustrates the case of uniform cluster distribution.