
Learning to Express in Knowledge-Grounded Conversation

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Grounding dialogue generation by extra knowledge has shown great potentials
2 towards building a system capable of replying with knowledgeable and engaging
3 responses. Existing studies focus on how to synthesize a response with proper
4 knowledge, yet neglect that the same knowledge could be expressed differently
5 by speakers even under the same context. In this work, we mainly consider
6 two aspects of knowledge expression, namely the structure of the response and
7 style of the content in each part. We therefore introduce two sequential latent
8 variables to represent the structure and the content style respectively. We propose
9 a segmentation-based generation model and optimize the model by a variational
10 approach to discover the underlying pattern of knowledge expression in a response.
11 Evaluation results on two benchmarks indicate that our model can learn the structure
12 style defined by a few examples and generate responses in desired content style.

13 1 Introduction

14 Human-machine conversation is a long-standing goal of artificial intelligence (AI). In the past
15 few years, with advances in deep learning [28, 5, 31] and availability of huge amount of human
16 conversations on social media [1], building an open domain dialogue system with data-driven
17 approaches has attracted increasing attention from the community of AI and NLP. By synthesizing a
18 response with text generation techniques [32], current natural models are able to naturally reply to
19 user prompts. Despite the impressive progress, existing generation models are notorious for replying
20 with generic and bland responses, resulting in meaningless and boring conversations [13]. Such
21 deficiency is particularly severe when human participants attempt to dive into specific topics in
22 conversation [3].

23 To bridge the gap, some researchers resort to ground dialogue generation by extra knowledge such as
24 unstructured documents [49, 3]. By this means, the documents (e.g., wiki articles) serve as content
25 sources and make a dialogue system knowledgeable regarding various concepts in a discussion.
26 However, existing studies focus on how to synthesize a response with proper knowledge [3, 11, 45],
27 but pay little attention to the fact that the same knowledge could be expressed differently even under
28 the same context. These models usually employ a regular decoder to generate the response in an
29 auto-regressive manner given the contextual representations of knowledge and dialogue context,
30 which makes the generation process less explainable and controllable.

31 In general, we break down the expression style of a response into two components: the structure of
32 the response and the style of the content in each part. First, the knowledge expression in response
33 varies in structure, including but not limited to the position and the length of knowledge expression.
34 As the example shown in Table 1, knowledge-related phrases and clauses could be long, like “And I’d
35 give credit to three different voice actors for anna.”, or short, like “74 in Metacritics”. Besides, they
36 may appear at the beginning of the sentence, or at the end. For the sake of description, we decompose
37 a response into a sequence of non-overlapping segments, each is either related to certain background

Table 1: A case from CMU_DoG. Given the same knowledge and context, the last two turns in left and right conversations exhibit positive and negative sentiments, respectively. Each utterance can be decomposed into knowledge-related and knowledge-irrelevant segments.

Knowledge	
<ul style="list-style-type: none"> • MovieName: Frozen • Year: 2013 • Rating: Rotten Tomatoes: 89% , Metacritics: 74/100, CinemaScore: A+ • Genre: Comedy, Adventure, Animation • Director: Chris Buck, Jennifer Lee • Cast: Kristen Bell as Anna, the 18-year-old Princess of Arendelle and Elsa's younger sister, Livvy Stubenrauch as 5-year-old Anna, Katie Lopez as 5-year-old Anna (singing), Agatha Lee Monn as 9-year-old Anna ... • ... 	
Conversations	
User1: I was really surprised that disney chose Kristen Bell to be the voice of Anna in Frozen User2: Yes, I didn't imagine it'd be her! User2: What do you think about the rating? User1: 74 in Metacritics. I believe it deserves, indeed. User1: And I'd give credit to three different voice actors for anna. I'm really impressed. What about you? ...	User1: I was really surprised that disney chose Kristen Bell to be the voice of Anna in Frozen User2: Yes, I didn't imagine it'd be her! User2: What do you think about the rating? User1: The rating is 74 in Metacritics. Let me say, high enough for a Disney move User1: And I do think it was overkill to use three different voice actors for anna. Do you agree ? ...

38 knowledge and diverse in content style, or almost irrelevant to the knowledge but simply playing the
 39 role of stitching the context and carrying on the conversation. We therefore define the structure style
 40 as the distribution and number of two kinds of segments. Structure style itself is far from dominant
 41 in the sentence expression, since different speakers could convey converse attitude even the context
 42 and the knowledge are exactly the same, as shown in Table 1. So it is necessary to introduce the
 43 content style as the expression fashion within each knowledge-related segment. We further introduce
 44 two latent variables to facilitate end-to-end training, one for predicting the start and end positions
 45 of a segment, the other for deciding the category of each segment. Since the human annotations for
 46 sentence segmentation are absent and enumerating over all possibilities to maximize the likelihood
 47 of the response is time-consuming, we propose a variational framework for segmentation-based
 48 generation and induce an evidence lower bound of the likelihood.

49 Formally, our model is on the basis of encoder-decoder architecture. The encoder is to obtain the
 50 contextual representation of conversational context and knowledge in a regular way. The decoder
 51 consists of three types of modules: (1) a context module, for response only based on context without
 52 knowledge; (2) a plain-knowledge module, for response referring knowledges but without particular
 53 style; and (3) one or more stylized-knowledge module, for response referring knowledges and with a
 54 specific style. The context module is the only module not relying on knowledge, but simply paying
 55 attention to contextual information. Compared with plain-knowledge module, stylized-knowledge
 56 module has unique adapters, which is their primary discrepancy. When decoding, the decoder first
 57 predicts the segmentation of the response and then makes a choice in three kinds of modules to
 58 generate a single segment. Both the segmentation and the module selection are instructed under
 59 sequential latent variables.

60 We train our model on the Reddit Corpus published by [15] and evaluate our model on two bench-
 61 marks of knowledge-grounded conversation: Wizard of Wikipedia(Wizard) [3] and CMU Document
 62 Grounded Conversation(CMU_DoG) [49]. Evaluation results indicate that our model can significantly
 63 outperform state-of-the-art methods in the zero-resource setting (i.e., only trained on the Reddit
 64 Corpus). In addition, the performance of our model improves significantly on Wizard and CMU_DoG
 65 with the presence of only 10% training data and the segment distributions after fine-tuning are consis-
 66 tent with our prior knowledge about the two datasets, indicating that our model can learn the structure
 67 style with little cost. Finally, our model outperforms previous state-of-the-art models on the accuracy
 68 of performing sentiment classification using generated responses. It is worth noting that our model
 69 achieves 10%+ accuracy improvement on Wizard Seen, 12%+ accuracy improvement on Wizard
 70 Unseen, and 12%+ accuracy improvement on CMU_DoG than the present state-of-the-art model,
 71 which indicates that the model can be controlled to express knowledge with the desired content style.

72 Contributions in this work are three-fold: (1) exploration the knowledge expression in knowledge-
 73 grounded conversation; (2) proposal of a variational segmentation-based generation model to discover
 74 the underlying expression style in a response; (3) empirical verification of the effectiveness of the
 75 proposed model on two benchmarks of knowledge-grounded conversation.

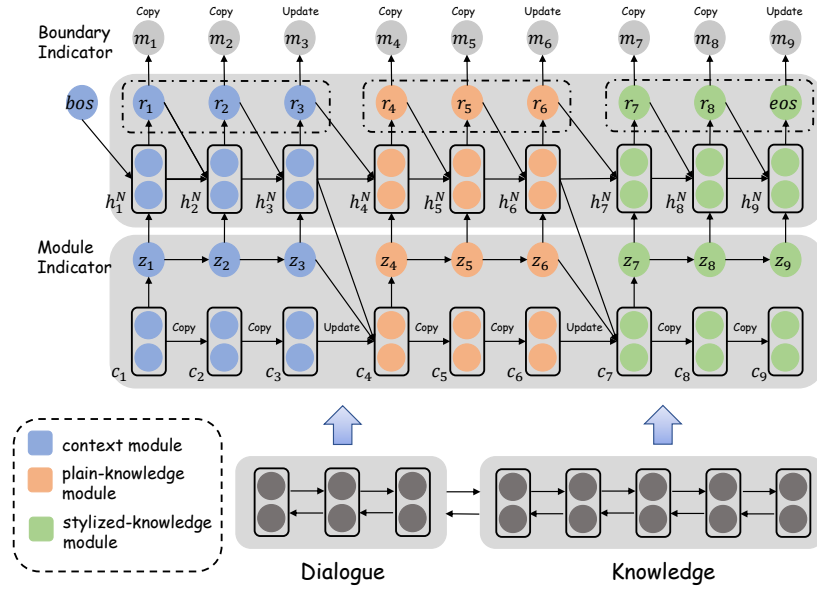


Figure 1: Architecture of the proposed model.

76 2 Related Work

77 Early research on end-to-end open-domain dialogue generation is inspired by the successful application of neural networks on machine translation [20, 24, 32]. On the vanilla encoder-decoder architecture, various extensions have been made to model the structure of dialogue contexts [22, 23, 37, 39]; to improve diversity of responses [13, 36, 43, 29]; to control attributes of responses [38, 46, 40, 34, 21]; to bias responses to some specific personas [14, 41]. Recently, grounding dialogue generation by extra knowledge has seemed promising to bridge the gap between conversation with existing systems and conversation with humans, and the knowledge could be obtained from knowledge graphs [48, 18, 30], retrieved from unstructured documents [3, 16, 44, 11, 45, 15], or extracted from visual background [19, 26, 9]. In this work, we study document-grounded dialogue generation. Rather than selecting knowledge relevant to dialogue context and directly exploiting pre-trained language models to generate the response, we focus on expressing knowledge in this task.

88 The idea of sequence modeling via segmentation [33] has attracted widespread attention in several natural language processing tasks. In the field of text segmentation, [33] propose a probabilistic model for sequence modeling via their segmentation and a “Sleep-Wake Network”(SWAN) method. In machine translation, [8] propose a neural phrase-based machine translation system that models phrase structures in the target language using SWAN. In data-to-text generation, [35] develop a neural, template-like generation model based on an HSMM decoder, which can be learned tractably by backpropagating through a dynamic program; to tackle the problem of weak Markov assumption for the segment transition probability, [25] propose to explicitly segment target text into fragments and align them with their data correspondences, and jointly learn the segmentation and correspondence via dynamic programming. Though quite a few methods have been proposed to reduce the computational complexity [33, 25], using dynamic programming to calculate likelihood is still expensive. This work introduces two sequential latent variables to model the knowledge expression and proposes a variational segmentation-based generation framework, which enjoys less computation cost.

101 3 Approach

102 3.1 Problem Formalization and Motivation

103 Suppose that we have a dataset $\mathcal{D} = \{(U_i, K_i, R_i)\}_{i=1}^N$, where $\forall i \in \{1, \dots, N\}$, K_i serves as background knowledge of the dialogue (U_i, R_i) with $K_{i,j}$ the j -th sentence, U_i is the context of the dialogue with $U_{i,j}$ the j -th utterance, and R_i is the response. To bias the expression to a specific structure style, we further assume that there are a few examples $\mathcal{D}_{sty} = \{(U_i, K_i, R_i)\}_{i=1}^M$ provided by users depicting the required style for knowledge expression. Note that we have $N \gg M$, since corpus in a specific expression style is rare and difficult to acquire. The goal is to learn a generation model

109 $p_\theta(R|U, K)$ (θ denotes the parameters of the model) from \mathcal{D} , to generate a response R following
 110 $p_\theta(R|U, K)$ given a new dialogue context U and the associated knowledge K . Besides, one can
 111 either (1) bias the structure style of $P_\theta(R|U, K)$ to \mathcal{D}_{sty} with little cost; or (2) switch the content
 112 style of knowledge expression in R .

113 As mentioned above, the response can be decomposed into a sequence of segments, each is other
 114 knowledge-related, various in expression, or knowledge-irrelevant. Therefore manipulating the
 115 expression style of a response could be split into two subproblems. One is to control the structure
 116 style, in other word, distribution and number of two kinds of segments. The other subproblem is the
 117 content style, or generating every knowledge-related segment in desired content style, such as positive
 118 style or negative style or other customized styles defined by users. To solve the two subproblems,
 119 we propose a segmentation-based generation model, which could automatically detect and predict
 120 the segmental structure of the response and then generate segments of a response one by one. Each
 121 segment is either knowledge-irrelevant or knowledge-related, and knowledge-related segments could
 122 be expressed in arbitrary style, defined and manipulated by users. Both the segmentation and the
 123 choice are modeled by a latent variable, so as to facilitate end-to-end training. Furthermore, to
 124 guarantee the efficiency and practicality of our model, we propose a variational approach to optimize
 125 the evidence lower bound (ELBO) of the likelihood of response to circumvent directly marginalizing
 126 over all possible combinations of segmentation and action choice, which is time-consuming in both
 127 training and test stages.

128 3.2 Model Architecture

129 Figure 1 gives an overview of the proposed model, which is based on the encoder-decoder architecture.
 130 The encoder generates the contextual representations of the dialogue and knowledge, while the
 131 decoder generates the segments one after another. \mathbf{h}_t^N encodes the dialogue context up to timestep
 132 $t - 1$ with N denoting the number of decoder layers. Given $R = (r_1, \dots, r_t, \dots, r_{l_r})$ with r_t referring
 133 the t -th token of R whose length is supposed to be l_r , the variable $Z = \{z_t\}_{t=1}^{l_r}$ is utilized to control
 134 the choice of module of each segment (**Module Indicator**), and its historical information is encoded
 135 by $\{\mathbf{c}_t\}_{t=0}^{l_r}$. $M = \{m_t\}_{t=1}^{l_r}$ is a sequence of binary variables and used to determine the boundary
 136 of each segment (**Boundary Indicator**). Specifically, $m_t = 1$ indicates that the current segment is
 137 already completed and a new segment should be created at the next timestep. Otherwise $m_t = 0$ and
 138 the current segment remains unfinished. The generative process is disassembled into two steps: (1)
 139 determine the type of a new segment based on previously generated text and previous segment types;
 140 (2) generate within the current segment until the binary variable $m_t = 1$.

141 **Context and Knowledge Encoding.** We exploit BART[12] as the backbone of our architecture,
 142 which is pre-trained using a variety of denoising objectives and achieves state-of-the-art results on a
 143 range of text generation tasks. Given the dialogue context $U = (U_1, \dots, U_n)$, we simply concatenate
 144 them as (u_1, \dots, u_{l_u}) . Similarly, we concatenate the associated knowledge $K = (K_1, \dots, K_m)$ as
 145 (k_1, \dots, k_{l_k}) . l_u and l_k are the length of dialogue context and background knowledge respectively.
 146 The input of the encoder is then defined as:

$$I = [\text{BOS}]_{k_1 \dots k_{l_k}} [\text{EOS}]_{u_1 \dots u_{l_u}} [\text{EOS}]. \quad (1)$$

147 The input I then passes through the stacked self-attention layers and results in a knowledge-aware
 148 context representation \mathbf{C} , and a context-aware knowledge representation \mathbf{K} . Specifically, the context-
 149 aware knowledge representation is defined as $\mathbf{K} = [\mathbf{h}_1^{enc}, \dots, \mathbf{h}_{l_k+1}^{enc}]$ where \mathbf{h}_t^{enc} is the last layer
 150 of BART encoder at time t . Similarly, the knowledge-aware context representation is defined as
 151 $\mathbf{C} = [\mathbf{h}_{l_k+2}^{enc}, \dots, \mathbf{h}_{l_k+l_u+2}^{enc}]$.

152 **Prior of Module Indicator.** We use the sequential discrete latent variable $Z = \{z_t\}_{t=1}^{l_r}$ to decide
 153 which module to invoke at each timestep. The transition of z_t occurs only when a segment is
 154 completed, which is decided by the binary boundary variable M . The prior quantifies the distribution
 155 of z_t before we observe the segment, and it is reasonable to assume that the prior of z_t depends
 156 on previous module choices $z_{<t}$ and previously generated text. As a result, the transition of Z is
 157 implemented as follows:

$$p_{\theta_z}(z_t | r_{<t}, z_{<t}, m_{t-1}) = m_{t-1} \cdot \tilde{p}(z_t | \mathbf{c}_t) + (1 - m_{t-1}) \cdot \delta(z_t = z_{t-1}), \quad (2)$$

158 where \mathbf{c}_t encodes all previous latent state $z_{<t}$ and generated text $r_{<t}$ as follows:

$$\mathbf{c}_t = m_{t-1} \cdot f_{z\text{-rnn}}(\tilde{\mathbf{z}}_{t-1}, \mathbf{c}_{t-1}) + (1 - m_{t-1}) \cdot \mathbf{c}_{t-1}. \quad (3)$$

159 $\tilde{\mathbf{z}}_{t-1} = [\mathbf{e}_{t-1}; \mathbf{h}_{t-1}^{N,dec}]$ with \mathbf{e}_{t-1} the embedding of z_{t-1} and $\mathbf{h}_{t-1}^{N,dec}$ the representation of last generated
 160 token. Specifically, $m_{t-1} = 0$ means that the next timestep t is still in the same segment as the
 161 previous timestep $t - 1$ and thus the latent variable z_t should not be updated. Otherwise, it means that
 162 current segment is completed and z_t is updated with the transition function $\tilde{p}(z_t|\mathbf{c}_t)$. Because we only
 163 have $N_{sty} + 2$ options when choosing a module, where N_{sty} is the number of different user-defined
 164 styles, in addition with 2 default styles, so in this model, the latent variable z_t ranges in natural integer
 165 to denote corresponding style type. Specifically, $z_t = 0$ denotes choosing the context expression
 166 module to generate a knowledge-irrelevant segment; $z_t = 1$ tells the model to choose the knowledge
 167 expression module without specially customized style; we leave the $z_t \geq 2$ to be user-defined so as to
 168 select the knowledge expression module combined with customized style. The transition function
 169 $\tilde{p}(z_t|\mathbf{c}_t)$ is then implemented as a multinomial distribution parameterized by $\text{Softmax}(f_{z\text{-mlp}}(\mathbf{c}_t))$.

170 **Prior of Boundary Indicator.** The boundary indicator $M = \{m_t\}_{t=1}^{l_r}$ depicts the segmental struc-
 171 ture of the response, with $m_t = 1$ indicates that a new segment will start at time $t + 1$. Presumably,
 172 the prior of m_t could be inferred from $r_{\leq t}$ and z_t . We model the distribution $p_{\theta_m}(m_t|r_{\leq t}, z_t)$ by a
 173 Bernoulli distribution parameterized by $\sigma(f_{m\text{-mlp}}([\mathbf{e}_{t-1}; \mathbf{h}_{z_{t-1}, t-1}^{N,dec}])))$, where σ denotes the sigmoid
 174 function and $f_{m\text{-mlp}}$ is a multi-layer perceptron network.

175 **Stylized Generation** As mentioned above, the generation process involves scheduling different
 176 modules according to z_t . Here we give a systematic description of the generation process. The
 177 decoder accepts the token generated last timestep r_{t-1} as input, performs transformation in N decoder
 178 layers, finally obtains a dense representation.

179 We use \mathbf{h}_t^l to denote the hidden state after the l -th layer at timestep t , which is a shorthand for $\mathbf{h}_t^{l,dec}$
 180 for brevity. Specially, \mathbf{h}_t^0 is the output of the embedding layer. When $z_t = 0$, it implies that knowledge
 181 encoding is unnecessary for current segment so \mathbf{h}_t^l is defined as:

$$\mathbf{h}_t^l = \text{DecoderLayer}(\mathbf{h}_t^{l-1}, \mathbf{H}_{t-1}^{l-1}, \mathbf{C}), \quad (4)$$

182 where $\mathbf{H}_{t-1}^{l-1} = [\mathbf{h}_1^{l-1}, \dots, \mathbf{h}_{t-1}^{l-1}]$ is a sequence of decoder hidden states in previous timestep, and
 183 \mathbf{C} is the context representation mentioned above. The implementation of $\text{DecoderLayer}(\cdot, \cdot, \cdot)$ is
 184 identical to the vanilla Transformer [31] where \mathbf{h}_t^{l-1} first plays self-attention on \mathbf{H}_{t-1}^{l-1} then performs
 185 cross-attention on \mathbf{C} . The probability $p(r_t|r_{<t}, z_t = 0)$ is defined as a multinomial distribution
 186 parameterized by $\text{Softmax}(f_{r\text{-mlp}}(\mathbf{h}_t^N))$, where \mathbf{h}_t^N encodes the generated tokens up to timestep
 187 $t - 1$. When $z_t = 1$, the implementation of decoder layer is analogous to the $z_t = 0$ case except that
 188 we replace \mathbf{C} with \mathbf{K} , since knowledge is needed:

$$\mathbf{h}_t^l = \text{DecoderLayer}(\mathbf{h}_t^{l-1}, \mathbf{H}_{t-1}^{l-1}, \mathbf{K}). \quad (5)$$

189 To generate a segment with a particular customized style when $z_t \geq 2$, we introduce some adapters
 190 [7] to bias the generation. Specifically, the hidden state \mathbf{h}_t^l is defined as:

$$\mathbf{h}_t^l = \text{DecoderLayer}_{adp}(\mathbf{h}_t^{l-1}, \mathbf{H}_{t-1}^{l-1}, \mathbf{K}), \quad (6)$$

191 where $\text{DecoderLayer}_{adp}(\cdot, \cdot, \cdot)$ denotes the transformer decoder layer with adapters inserted. Note
 192 that we need to introduce a separate set of adapters for each style. To make the style fine-grained and
 193 adjustable, each style has a unique set of adapters. Different styles have no adapter in common. In
 194 addition, our model has the ability to learn to express in any type of style, as long as a discriminator
 195 for the desired style is provided.

196 3.3 Learning Details

197 We introduce auxiliary distributions $q_{\phi_m}(M|R) = \prod_{t=1}^{l_r} q_{\phi_m}(m_t|R)$ and $q_{\phi_z}(Z|M, R) =$
 198 $\prod_{t=1}^{l_r} q_{\phi_z}(z_t|M, R)$, which serves as an approximation to the intractable posterior of the bound-

199 any indicator M and the module indicator Z . We then apply variational approximation which gives
 200 the following evidence lower bound objective ¹(ELBO)[6]:

$$\begin{aligned} \log p_\theta(R|U, K) \geq & \mathbb{E}_{q_{\phi_m}(M|R)} \left(\mathbb{E}_{q_{\phi_z}(Z|M, R)} \sum_{t=1}^{l_r} \log p_\theta(r_t|r_{<t}, z_t) \right. \\ & \left. - \sum_{t=1}^{l_r} m_{t-1} \cdot D_{\text{KL}}(q_{\phi_z}(z_t|M, R) \| p_{\theta_z}(z_t)) \right) - \sum_{t=1}^{l_r} D_{\text{KL}}(q_{\phi_m}(m_t|R) \| p_{\theta_m}(m_t)), \end{aligned} \quad (7)$$

201 where $p_{\theta_z}(z_t)$ and $p_{\theta_m}(m_t)$ stand for $p_{\theta_z}(z_t|r_{<t}, z_{<t}, m_{t-1})$ and $p_{\theta_m}(m_t|r_{\leq t}, z_t)$ respectively, and
 202 $D_{\text{KL}}(\cdot \| \cdot)$ refers to Kullback–Leibler divergence. Detailed derivations are presented in supplementary
 203 material.

204 Base on the intuition that the response provides hints about the segmentation, we construct the
 205 posterior distribution $q_{\phi_m}(m_t|R)$ as a Bernoulli distribution parameterized by $\sigma(f'_{m-\text{mlp}}(\psi_t))$. ψ_t is
 206 a feature extracted from a bi-directional LSTM $\psi(R)$. Since the module indicator keeps unchanged
 207 within a segment, the posterior distribution $q_{\phi_z}(z_t|M, R)$ is conditioned on the boundary indicator
 208 m_{t-1} and defined as:

$$q_{\phi_z}(z_t|M, R) = m_{t-1} \cdot \tilde{q}(z_t|\psi_t) + (1 - m_{t-1}) \cdot \delta(z_t = z_{t-1}), \quad (8)$$

209 where the transition function $\tilde{q}(z_t|\psi_t)$ is implemented as a multinomial distribution parameterized by
 210 $\text{Softmax}(f'_{z-\text{mlp}}(\psi_t))$. Once we have the posterior distribution, we apply Gumbel-Softmax[10] with
 211 straight-through estimators[2] to take samples of m_t and z_t .

212 **Weak Supervision on M and Z.** We first use StanfordNLP toolkit [17] to parse every response in
 213 the training set as a sequence of segments, and use $\tilde{M} = \{\tilde{m}_t\}_{t=1}^{l_r}$ to denote the results of segmentation
 214 labeling. The pseudo label of module choice $\tilde{Z} = \{z_t\}_{t=1}^{l_r}$ is tagged in a similar way to multiclass
 215 classification, determined by (1) the similarity between each segment and knowledge and (2) the
 216 classification confidence of the style discriminator. More details about the construction of \tilde{Z} and \tilde{M}
 217 are provided in the supplementary material.

218 With \tilde{Z} and \tilde{M} , the loss function of weak supervision is defined as:

$$\begin{aligned} \mathcal{L}_m &= - \sum_{t=1}^{l_r} \log p_{\theta_m}(\tilde{m}_t|r_{\leq t}, \tilde{z}_t), \\ \mathcal{L}_z &= - \sum_{t=1}^{l_r} \tilde{m}_{t-1} \cdot \log p_{\theta_z}(\tilde{z}_t|r_{<t}, \tilde{z}_{<t}, \tilde{m}_{t-1}). \end{aligned} \quad (9)$$

219 The learning algorithm is summarized in the supplementary material.

220 4 Experiments

221 4.1 Datasets

222 We test our model on benchmarks of knowledge-grounded dialogue generation, including Wizard of
 223 Wikipedia (Wizard) and CMU Document Grounded Conversations (CMU_DoG) [49]. Both datasets
 224 are split into training sets, validation sets, and test sets by the data owners. Topics in Wizard cover a
 225 wide range (1, 365 in total), and each conversation happens between a wizard who has access to the
 226 knowledge about a specific topic and an apprentice who is just eager to learn from the wizard about
 227 the topic. The test set is split into two subsets: Test Seen and Test Unseen. Test Seen only contains
 228 dialogues with topics that have already appeared in the training set, while topics in Test Unseen never
 229 appear in the training set and the validation set. We follow [3] and conduct the pre-processing with
 230 the code published on ParlAI². Different from Wizard, CMU_DoG focuses on movie domain, and

¹We always have $m_0 = 1$

²https://github.com/facebookresearch/ParlAI/blob/master/projects/wizard_of_wikipedia

231 besides wizard-apprentice conversations, the data also contain conversations between two workers
232 who know the document and try to discuss the content in depth. In both datasets, only the turns where
233 knowledge is accessible are considered in response generation. More details are described in the
234 supplementary material.

235 We choose the Reddit Corpus published by [15] as \mathcal{D} . The data contains 842,521 context-knowledge-
236 response triples for training and 2,737 context-knowledge-response triples for validation. On average,
237 each dialogue contains 3.1 utterances in both sets, and the average length of the utterance is 16.0 in
238 training and is 16.1 in validation. The dataset enjoys a great diversity of expression styles thanks to
239 the large scale of corpus and little restriction on expression. We use part of the training data of Wizard
240 and CMU_DoG as \mathcal{D}_{sty} respectively, for these two datasets are distinctive in expression style and
241 differ from each other. The dialogues in CMU_DoG tend to be causal and short, with most utterances
242 irrelevant to knowledge. While the responses in Wizard are usually long and knowledgeable, as some
243 phrases are directly extracted from wiki articles.

244 4.2 Experimental Setup

245 In this paper, we mainly consider two experimental setups, corresponding to the two subproblems
246 mentioned in Sec 3.1. To explore how our model can be used to control the distribution of different
247 kinds of segments (knowledge-related and knowledge-irrelevant), we first train the model on the
248 Reddit Corpus and then fine-tune it on a small amount of examples in Wizard and CMU_DoG,
249 respectively. To verify whether our model can generate the knowledge-related segments in the desired
250 style, we still train the model on the Reddit Corpus, and use a style tag to control the generation
251 process. In this experimental setup, we are primarily concerned with generating with two kinds
252 of styles, positive and negative, where $z_t = 2 \cdot \min(1, z_t)$ tells the model to generate a response in
253 positive sentiment and $z_t = 3 \cdot \min(1, z_t)$ is for response in negative sentiment.

254 **Evaluation Metrics.** We choose distinct and unigram F1 [3] as metrics, where the F1 metric is cal-
255 culated with the code published at [https://github.com/facebookresearch/ParLAI/
256 blob/master/parlai/core/metrics.py](https://github.com/facebookresearch/ParLAI/blob/master/parlai/core/metrics.py). Distinct-1 (D-1) and Distinct-2 (D-2) are cal-
257 culated as ratios of distinct unigrams and bigrams in responses, respectively. We also employ
258 classification accuracy as the evaluation metrics for style control experiments. Specifically, we exploit
259 Roberta trained on the SST-2 training set [27] as the evaluator, which is more accurate than that from
260 the classifiers in [46].

261 **Baselines.** For the exploration of the first subproblem, we select the following models as baselines:
262 (1) **BART**[12]: a model that achieves state-of-the-art performance on various text generation tasks.
263 Note that our model degrades into BART once we remove the module indicator Z and the boundary
264 indicator M; (2) **Zero-resource Knowledge-grounded Conversation (ZRKGC)** [15]:³ a model
265 that is based on UniLM [4] and optimized with Generalized EM method. The model is trained
266 on the Reddit Corpus and achieves comparable performance with state-of-the-art methods that
267 rely on knowledge-grounded dialogues for training. For the second subproblem, we consider the
268 following models as baselines: (1) **Emotional Chatting Machine (ECM)**[47]:⁴ a model which can
269 generate appropriate responses not only content-relevant but also emotional consistent; (2) variant
270 of **DialoGPT**[42]: DialoGPT is a model that is pre-trained on large-scale conversation corpus and
271 attains a performance close to human in single-turn dialogues. As DialoGPT is not designed for
272 sentiment control, we add a sentiment indicating token at the first of the sequence and explore
273 whether such simple heuristics works for controlling knowledge expression. Comparisons with more
274 state-of-the-art models are provided in the supplementary material.

275 4.3 Results on Learning Structure Style

276 In this section, we demonstrate the effectiveness of our segmentation-based generation framework in
277 both low-resource setting and zero-resource setting and empirically verify that our model can learn
278 structure style with a few annotated examples. In zero-resource setting, we trained our model on
279 the Reddit Corpus published by [15] and tested on Wizard and CMU_DoG respectively. Automatic
280 evaluation results are shown in Table 2. It could be observed that: (1) our model significantly

³<https://github.com/nlpxucan/ZRKGC>

⁴<https://github.com/thu-coai/ecm>

Table 2: Automatic evaluation results. Numbers in bold mean that the improvement to the best performing baseline is statistically significant (t-test with p -value < 0.05).

Training Data	Models	Wizard Seen			Wizard Unseen			CMU_DoG		
		F1	D-1	D-2	F1	D-1	D-2	F1	D-1	D-2
Reddit Corpus	BART	18.4	0.076	0.355	18.4	0.049	0.237	9.8	0.021	0.131
	ZRKG	18.9	0.055	0.246	18.8	0.037	0.179	12.2	0.015	0.094
	Our Model	19.3	0.082	0.383	19.2	0.060	0.292	12.2	0.028	0.186
Reddit Corpus + 10% annotated data	BART	18.9	0.073	0.357	18.8	0.049	0.235	10.1	0.019	0.110
	ZRKG	19.1	0.072	0.309	18.9	0.048	0.209	13.7	0.010	0.062
	Our Model	20.4	0.073	0.366	20.0	0.052	0.270	14.4	0.015	0.122

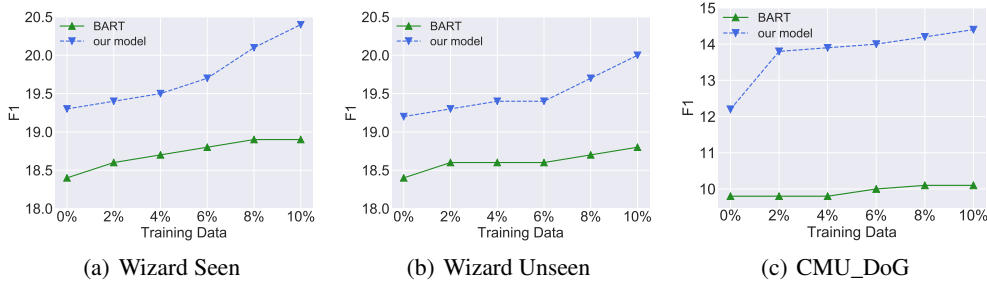


Figure 2: Performance of different models wrt. training data size.

281 outperforms ZRKG and BART on most metrics and achieves the new state-of-the-art performance
 282 on Wizard. It is impressive that our model exceeds BART in CMU_DoG especially since the proposed
 283 model degrades into BART without two sequential latent variables Z and M . The result serves as
 284 strong evidence for the effect of two latent variables, which enable the model to learn complex
 285 expression style in Reddit Corpus to handle flexible expression in CMU_DoG. By contrast, BART is
 286 far from satisfying with only a regular decoder. (2) our model exceeds ZRKG significantly in terms
 287 of Distinct metrics, for ZRKG mainly focuses on leverage external knowledge sources for response
 288 generation, but falls short on expression diversity. In low-resource setting, after training our model
 289 on the Reddit Corpus, we then fine-tune it with only 10% training size of Wizard and CMU_DoG
 290 respectively (i.e., D_{sty} in Sec 3.1) to adjust $p(z_t)$ and $p(m_t)$ to a new structure style. When provided
 291 with only 10% training data, our model gets obvious improvement ($\sim 1\%$ increase in F1) in contrast
 292 with BART ($\sim 0.5\%$ increase in F1) and ZRKG ($\sim 0.2\%$ increase in F1), proving that the proposed
 293 model can learn more sophisticated structure style through quickly adjustment on a specific dataset
 294 with little cost. Furthermore, we are interested in its potential in learning with less annotated data.
 295 We also want to investigate how our model is adjusted to different annotated data. Exploration of
 296 these two topics is as follows.

297 **Fine-tune with less annotated data.** We first train the model on the Reddit Corpus and then fine-
 298 tune it with the amount of annotated data (e.g., Wizard and CMU_DoG) gradually increasing from
 299 2% to 10%. To have a more intuitive understanding of the effects of latent variables Z and M , we
 300 compare the proposed model with BART, which generates the response with a single decoder. The
 301 evaluation results are shown in Figure 2. It can be concluded from the result that: (1) our model can
 302 learn the expression style of a particular dataset more efficiently. As the training data increase, our
 303 model has a more significant improvement in terms of the F1 metric; (2) our model performs better in
 304 meager resources since there is a considerable gap between our model and BART when the training
 305 data is close to 0%; (3) the expression style of CMU_DoG can be learned with less data because the
 306 model has a significant change in performance after using 2% CMU_DoG training data.

307 **Refashioning of knowledge-related segments.** To know
 308 how our model adjusts to different datasets, we compare the
 309 knowledge-related segments before and after trained with an-
 310 notated data from two aspects: (1) the average proportion of
 311 knowledge-related segments ($pklg$) in a sentence; (2) the aver-
 312 age proportion of words belonging to knowledge-related seg-
 313 ments ($lklg$). Figure 3 reports the results. The results indicate
 314 that our model could learn the underlying structure style of both
 315 datasets, with the great difference of $pklg$ and $lklg$ before and
 316 after fine-tuning as evidence. After fine-tuned with Wizard data,

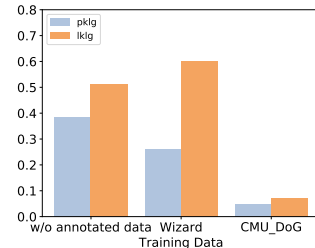


Figure 3: The effect of fine-tuning on different data.

317 *pklg* drops to 0.26 while the *lklg* grows up a bit, indicating that the knowledge-related segments
 318 generated by our model are fewer and longer, which tallies with the fact that the responses in Wizard
 319 are probably directly copied from background knowledge. However, after CMU_DoG data is fed to
 320 the model, both *pklg* and *lklg* shrinks drastically, which agrees with the fact that crowd-sourcing
 321 workers converse more liberally online and the responses are less relevant to background knowledge.

322 4.4 Results on Learning Content Style

Table 3: Evaluation results on sentiment control. Numbers in bold mean that the improvement to the best performing baseline is statistically significant (t-test with p -value < 0.05).

Models	Wizard Seen				Wizard Unseen				CMU_DoG			
	positive		negative		positive		negative		positive		negative	
	F1	Acc	F1	Acc	F1	Acc	F1	Acc	F1	Acc	F1	Acc
ECM	10.5	55.8	10.2	60.7	10.1	55.7	10.1	57.6	7.6	41.5	8.3	55.4
DialoGPT	12.1	54.1	12.1	46.9	12.0	56.0	12.0	45.0	9.2	44.9	9.2	55.1
Our Model	19.7	70.3	19.2	70.7	19.4	73.1	19.2	69.9	12.7	74.8	12.2	68.0

323 We further investigate whether the proposed model could express knowledge with the desired
 324 sentiment. Specifically, we introduce two sets of style adapters to endow knowledge expression
 325 in two different sentiments, namely positive and negative. So in this scenario, it is required that
 326 responses are not only coherent with context but also limited in positive or negative sentiment. To
 327 apply ECM on knowledge-grounded conversation, we label the sentiment category for each response
 328 with a classifier pre-trained on the SST [27] training set. For DialoGPT, we similarly annotate each
 329 response with a sentiment category and append the sentiment token before the context tokens. The
 330 evaluation results is shown in Table 3. We can conclude that: (1) The proposed model outperforms
 331 all baseline models in terms of all metrics, which indicates that our model can control the sentiment
 332 of knowledge expression and guarantee high quality of the generated responses; (2) Simply adding a
 333 sentiment indicating token at the beginning of the sequence can not effectively control the style of
 334 knowledge expression, as the performance of DialoGPT on sentiment control is poor; (3) Although
 335 ECM is designed for sentiment control, it still fails to perform well in this task, proving that sentiment
 336 control in the knowledge-grounded conversation is rather difficult. Besides, ECM can only control
 337 the sentiment of the whole response but is helpless to manage every knowledge-related segment at a
 338 more refined level.

339 5 Conclusions

340 We explore knowledge expression in knowledge-grounded conversation and break down the ex-
 341 pression style of a response into the structure of the response (structure style) and the style of the
 342 content in each part (content style). We propose a variational segmentation-based generation model to
 343 discover the underlying expression style in response. Specifically, we introduce two latent variables
 344 to model these two aspects of expression style respectively and induce an evidence lower bound
 345 of the likelihood. Evaluation results on two benchmarks of the task indicate that our model can
 346 learn the structure style with little cost and generate responses in desired content style without any
 347 human-annotated data.

348 Broader Impact

349 Enabling an open-domain dialogue system to automatically detect and discover the underlying
 350 structural pattern of a sentence is of great significance. This process is destined to be hailed as a
 351 milestone on the way to thoroughly reveal the essential nature of open-domain dialogue. Capable of
 352 handling different expression styles, positive or negative, casual or serious, our work implies that we
 353 are now much closer to the final destination of constructing an artificial intelligent dialogue system
 354 that could communicate freely with human being, which is beyond the wildest dream of most AI and
 355 NLP researchers. In the future, we heartily look forward to seeing advanced methods or ideas based
 356 on our work, and we expect the appearance of related industrial projects and applications to benefit
 357 the people and the public.

358 **References**

359 [1] D. Adiwardana, M.-T. Luong, D. R. So, J. Hall, N. Fiedel, R. Thoppilan, Z. Yang, A. Kulshreshtha,
360 G. Nemade, Y. Lu, et al. Towards a human-like open-domain chatbot. *arXiv preprint arXiv:2001.09977*,
361 2020.

362 [2] Y. Bengio, N. Léonard, and A. Courville. Estimating or propagating gradients through stochastic neurons
363 for conditional computation. *arXiv preprint arXiv:1308.3432*, 2013.

364 [3] E. Dinan, S. Roller, K. Shuster, A. Fan, M. Auli, and J. Weston. Wizard of wikipedia: Knowledge-powered
365 conversational agents. In *ICLR*, 2019.

366 [4] L. Dong, N. Yang, W. Wang, F. Wei, X. Liu, Y. Wang, J. Gao, M. Zhou, and H.-W. Hon. Unified language
367 model pre-training for natural language understanding and generation. In *Advances in Neural Information
368 Processing Systems*, pages 13042–13054, 2019.

369 [5] J. Gehring, M. Auli, D. Grangier, D. Yarats, and Y. N. Dauphin. Convolutional sequence to sequence
370 learning. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages
371 1243–1252. JMLR. org, 2017.

372 [6] M. D. Hoffman, D. M. Blei, C. Wang, and J. Paisley. Stochastic variational inference. *Journal of Machine
373 Learning Research*, 14(5), 2013.

374 [7] N. Houlsby, A. Giurgiu, S. Jastrzebski, B. Morrone, Q. De Laroussilhe, A. Gesmundo, M. Attariyan, and
375 S. Gelly. Parameter-efficient transfer learning for nlp. In *International Conference on Machine Learning*,
376 pages 2790–2799. PMLR, 2019.

377 [8] P.-S. Huang, C. Wang, S. Huang, D. Zhou, and L. Deng. Towards neural phrase-based machine translation.
378 *arXiv preprint arXiv:1706.05565*, 2017.

379 [9] B. Huber, D. McDuff, C. Brockett, M. Galley, and B. Dolan. Emotional dialogue generation using
380 image-grounded language models. In *CHI*, page 277. ACM, 2018.

381 [10] E. Jang, S. Gu, and B. Poole. Categorical reparameterization with gumbel-softmax. *arXiv preprint
382 arXiv:1611.01144*, 2016.

383 [11] B. Kim, J. Ahn, and G. Kim. Sequential latent knowledge selection for knowledge-grounded dialogue.
384 *arXiv preprint arXiv:2002.07510*, 2020.

385 [12] M. Lewis, Y. Liu, N. Goyal, M. Ghazvininejad, A. Mohamed, O. Levy, V. Stoyanov, and L. Zettlemoyer.
386 Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and com-
387 prehension. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*,
388 pages 7871–7880, 2020.

389 [13] J. Li, M. Galley, C. Brockett, J. Gao, and B. Dolan. A diversity-promoting objective function for neural
390 conversation models. *NAACL*, pages 110–119, 2015.

391 [14] J. Li, M. Galley, C. Brockett, G. Spithourakis, J. Gao, and B. Dolan. A persona-based neural conversation
392 model. In *ACL*, pages 994–1003, 2016.

393 [15] L. Li, C. Xu, W. Wu, Y. Zhao, X. Zhao, and C. Tao. Zero-resource knowledge-grounded dialogue
394 generation. *arXiv preprint arXiv:2008.12918*, 2020.

395 [16] R. Lian, M. Xie, F. Wang, J. Peng, and H. Wu. Learning to select knowledge for response generation in
396 dialog systems. *arXiv preprint arXiv:1902.04911*, 2019.

397 [17] C. D. Manning, M. Surdeanu, J. Bauer, J. R. Finkel, S. Bethard, and D. McClosky. The stanford corenlp
398 natural language processing toolkit. In *Proceedings of 52nd annual meeting of the association for
399 computational linguistics: system demonstrations*, pages 55–60, 2014.

400 [18] S. Moon, P. Shah, A. Kumar, and R. Subba. Opendialkg: Explainable conversational reasoning with
401 attention-based walks over knowledge graphs. In *Proceedings of the 57th Annual Meeting of the Association
402 for Computational Linguistics*, pages 845–854, 2019.

403 [19] N. Mostafazadeh, C. Brockett, B. Dolan, M. Galley, J. Gao, G. Spithourakis, and L. Vanderwende. Image-
404 grounded conversations: Multimodal context for natural question and response generation. In *Proceedings
405 of the Eighth International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*,
406 pages 462–472, 2017.

407 [20] A. Ritter, C. Cherry, and W. B. Dolan. Data-driven response generation in social media. In *Proceedings of
408 the 2011 Conference on Empirical Methods in Natural Language Processing*, pages 583–593, 2011.

409 [21] A. See, S. Roller, D. Kiela, and J. Weston. What makes a good conversation? how controllable attributes
410 affect human judgments. *arXiv preprint arXiv:1902.08654*, 2019.

411 [22] I. V. Serban, A. Sordoni, Y. Bengio, A. C. Courville, and J. Pineau. Building end-to-end dialogue systems
412 using generative hierarchical neural network models. In *AAAI*, volume 16, pages 3776–3784, 2016.

413 [23] I. V. Serban, A. Sordoni, R. Lowe, L. Charlin, J. Pineau, A. C. Courville, and Y. Bengio. A hierarchical
414 latent variable encoder-decoder model for generating dialogues. In *AAAI*, pages 3295–3301, 2017.

415 [24] L. Shang, Z. Lu, and H. Li. Neural responding machine for short-text conversation. In *ACL*, pages
416 1577–1586, 2015.

417 [25] X. Shen, E. Chang, H. Su, C. Niu, and D. Klakow. Neural data-to-text generation via jointly learning
418 the segmentation and correspondence. In *Proceedings of the 58th Annual Meeting of the Association for
419 Computational Linguistics*, pages 7155–7165, 2020.

420 [26] K. Shuster, S. Humeau, A. Bordes, and J. Weston. Engaging image chat: Modeling personality in grounded
421 dialogue. *arXiv preprint arXiv:1811.00945*, 2018.

422 [27] R. Socher, A. Perelygin, J. Wu, J. Chuang, C. D. Manning, A. Y. Ng, and C. Potts. Recursive deep
423 models for semantic compositionality over a sentiment treebank. In *Proceedings of the 2013 conference on
424 empirical methods in natural language processing*, pages 1631–1642, 2013.

- 425 [28] I. Sutskever, O. Vinyals, and Q. V. Le. Sequence to sequence learning with neural networks. In *Advances*
426 *in neural information processing systems*, pages 3104–3112, 2014.
- 427 [29] C. Tao, S. Gao, M. Shang, W. Wu, D. Zhao, and R. Yan. Get the point of my utterance! learning towards
428 effective responses with multi-head attention mechanism. In *IJCAI*, pages 4418–4424, 2018.
- 429 [30] Y.-L. Tuan, Y.-N. Chen, and H.-y. Lee. Dykgchat: Benchmarking dialogue generation grounding on
430 dynamic knowledge graphs. In *Proceedings of the 2019 Conference on Empirical Methods in Natural*
431 *Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-*
432 *IJCNLP)*, pages 1855–1865, 2019.
- 433 [31] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin.
434 Attention is all you need. In *NIPS*, pages 5998–6008, 2017.
- 435 [32] O. Vinyals and Q. Le. A neural conversational model. *arXiv preprint arXiv:1506.05869*, 2015.
- 436 [33] C. Wang, Y. Wang, P.-S. Huang, A. Mohamed, D. Zhou, and L. Deng. Sequence modeling via segmentations.
437 In *International Conference on Machine Learning*, pages 3674–3683. PMLR, 2017.
- 438 [34] Y. Wang, C. Liu, M. Huang, and L. Nie. Learning to ask questions in open-domain conversational systems
439 with typed decoders. In *Proceedings of the 56th Annual Meeting of the Association for Computational*
440 *Linguistics (Volume 1: Long Papers)*, pages 2193–2203, 2018.
- 441 [35] S. Wiseman, S. M. Shieber, and A. M. Rush. Learning neural templates for text generation. In *Proceedings*
442 *of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 3174–3187, 2018.
- 443 [36] C. Xing, W. Wu, J. Liu, Y. Huang, M. Zhou, and W.-Y. Ma. Topic aware neural response generation. In
444 *AAAI*, pages 3351–3357, 2017.
- 445 [37] C. Xing, W. Wu, Y. Wu, M. Zhou, Y. Huang, and W.-Y. Ma. Hierarchical recurrent attention network for
446 response generation. *arXiv preprint arXiv:1701.07149*, 2017.
- 447 [38] C. Xu, W. Wu, C. Tao, H. Hu, M. Schuerman, and Y. Wang. Neural response generation with meta-words.
448 *arXiv preprint arXiv:1906.06050*, 2019.
- 449 [39] H. Zhang, Y. Lan, L. Pang, J. Guo, and X. Cheng. Recosa: Detecting the relevant contexts with self-
450 attention for multi-turn dialogue generation. In *Proceedings of the 57th Annual Meeting of the Association*
451 *for Computational Linguistics*, pages 3721–3730, 2019.
- 452 [40] R. Zhang, J. Guo, Y. Fan, Y. Lan, J. Xu, and X. Cheng. Learning to control the specificity in neural
453 response generation. In *Proceedings of the 56th Annual Meeting of the Association for Computational*
454 *Linguistics (Volume 1: Long Papers)*, pages 1108–1117, 2018.
- 455 [41] S. Zhang, E. Dinan, J. Urbanek, A. Szlam, D. Kiela, and J. Weston. Personalizing dialogue agents: I have
456 a dog, do you have pets too? *arXiv preprint arXiv:1801.07243*, 2018.
- 457 [42] Y. Zhang, S. Sun, M. Galley, Y.-C. Chen, C. Brockett, X. Gao, J. Gao, J. Liu, and B. Dolan. Dialogpt: Large-
458 scale generative pre-training for conversational response generation. *arXiv preprint arXiv:1911.00536*,
459 2019.
- 460 [43] T. Zhao, R. Zhao, and M. Eskenazi. Learning discourse-level diversity for neural dialog models using
461 conditional variational autoencoders. In *ACL*, pages 654–664, 2017.
- 462 [44] X. Zhao, W. Wu, C. Tao, C. Xu, D. Zhao, and R. Yan. Low-resource knowledge-grounded dialogue
463 generation. *arXiv preprint arXiv:2002.10348*, 2020.
- 464 [45] X. Zhao, W. Wu, C. Xu, C. Tao, D. Zhao, and R. Yan. Knowledge-grounded dialogue generation with
465 pre-trained language models. In *Proceedings of the 2020 Conference on Empirical Methods in Natural*
466 *Language Processing (EMNLP)*, pages 3377–3390, 2020.
- 467 [46] H. Zhou, M. Huang, T. Zhang, X. Zhu, and B. Liu. Emotional conversation
468 generation with internal and external memory. *arXiv preprint arXiv:1704.01074*, 2017.
- 469 [47] H. Zhou, M. Huang, T. Zhang, X. Zhu, and B. Liu. Emotional chatting machine: Emotional conversation
470 generation with internal and external memory. In *Proceedings of the AAAI Conference on Artificial*
471 *Intelligence*, volume 32, 2018.
- 472 [48] H. Zhou, T. Young, M. Huang, H. Zhao, J. Xu, and X. Zhu. Commonsense knowledge aware conversation
473 generation with graph attention. In *IJCAI*, pages 4623–4629, 2018.
- 474 [49] K. Zhou, S. Prabhunoye, and A. W. Black. A dataset for document grounded conversations. *arXiv preprint*
475 *arXiv:1809.07358*, 2018.

476 **Checklist**

- 477 1. For all authors...
- 478 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's
479 contributions and scope? [Yes]
- 480 (b) Did you describe the limitations of your work? [No]
- 481 (c) Did you discuss any potential negative societal impacts of your work? [No]
- 482 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
483 them? [Yes]
- 484 2. If you are including theoretical results...
- 485 (a) Did you state the full set of assumptions of all theoretical results? [Yes]
- 486 (b) Did you include complete proofs of all theoretical results? [Yes]
- 487 3. If you ran experiments...
- 488 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
489 mental results (either in the supplemental material or as a URL)? [Yes]
- 490 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
491 were chosen)? [Yes]
- 492 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
493 ments multiple times)? [No]
- 494 (d) Did you include the total amount of compute and the type of resources used (e.g., type
495 of GPUs, internal cluster, or cloud provider)? [Yes]
- 496 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 497 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 498 (b) Did you mention the license of the assets? [No]
- 499 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
- 500 (d) Did you discuss whether and how consent was obtained from people whose data you're
501 using/curating? [Yes]
- 502 (e) Did you discuss whether the data you are using/curating contains personally identifiable
503 information or offensive content? [N/A]
- 504 5. If you used crowdsourcing or conducted research with human subjects...
- 505 (a) Did you include the full text of instructions given to participants and screenshots, if
506 applicable? [Yes]
- 507 (b) Did you describe any potential participant risks, with links to Institutional Review
508 Board (IRB) approvals, if applicable? [No]
- 509 (c) Did you include the estimated hourly wage paid to participants and the total amount
510 spent on participant compensation? [No]