

DG²: Data Augmentation Through Document Grounded Dialogue Generation

Anonymous ACL submission

Abstract

Collecting data for training dialog systems can be extremely expensive due to the involvement of human participants and the need for extensive annotation. Especially in document-grounded dialog systems, human experts need to carefully read the unstructured documents to answer the users' questions. As a result, existing document-grounded dialog datasets are relatively small-scale and obstruct the effective training of dialogue systems. In this paper, we propose an automatic data augmentation technique grounded on documents through a generative dialogue model. The dialogue model consists of a user bot and agent bot that can synthesize diverse dialogues given an input document, which are then used to train a downstream model. When supplementing the original dataset, our method achieves significant improvement over traditional data augmentation methods. We also achieve great performance in the low-resource setting.

1 Introduction

Most of human knowledge is stored in the form of documents, ranging from answering factoid questions (Reddy et al., 2019) to providing how-tos on millions of tasks (Zhang et al., 2020a). How to comprehend and retrieve relevant knowledge from documents given a user query is a challenging research problem. Inspired by real-world applications, there have been more works (Rajpurkar et al., 2016a, 2018; Kwiatkowski et al., 2019; Yang et al., 2015) that aims to tackle this challenge. In this work, we focus on the task of conversational information seeking based on the associated documents, which are often referred to as document-grounded dialogue systems (Ma et al., 2020).

Recent works have introduced various datasets for building document-grounded conversational question answering and dialogue systems. Some work such as QuAC (Choi et al., 2018) and CoQA (Reddy et al., 2019) first explored the direction of

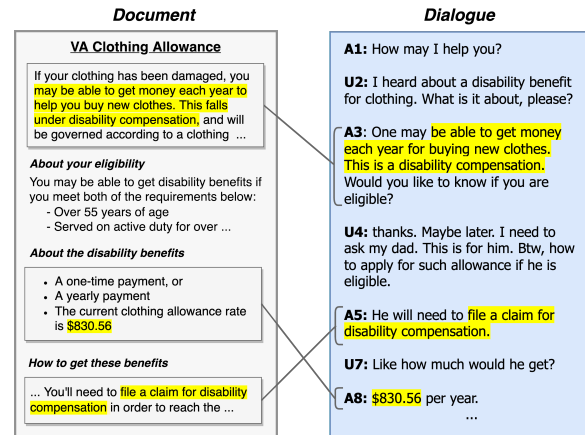


Figure 1: An example from Doc2Dial of dialogue conversation produced from grounding to an associated document. The agent must select the correct spans and engage in a fluent manner to generate a proper response.

conversational question answering. Then, ShARC (Saeidi et al., 2018) added follow-up questions by agents. Later, Doc2Dial (Feng et al., 2020a) further expands the dialogue actions and domains, which aims to simulate more kinds of real-life scenarios. However, such dataset is typically hard to scale up and to new domains, as it requires carefully crafted dialogue flows and expensive human annotations.

As a consequence, one main obstacle for developing scalable and effective document grounded dialog systems is the lack of sufficient data. In chat scenarios, recent works such as DialoGPT (Zhang et al., 2020b), Meena (Adiwardana et al., 2020), and Blender (Roller et al., 2021) have achieved human-like performance by taking the advantage of training on a large-scale corpus. Similarly, task-oriented dialog systems such as ARDM (Wu et al., 2021) and SimpleTOD (Hosseini-Asl et al., 2020) have also utilized large-scale corpora or pre-trained models to achieve good performance. The aforementioned models were trained with millions of samples, while the current document-grounded dialogue datasets like Doc2Dial (Feng

065 et al., 2020a) only contains thousands of conver- 116
066 sations. Training on such a small-scale dataset 117
067 constrains the performance of neural network mod- 118
068 els. Therefore, augmenting existing datasets can 119
069 help build a more effective document-grounded 120
070 dialogue system. 121

071 One popular approach to augmenting datasets is 122
072 to paraphrase existing seed data. The most straight- 123
073 forward form of paraphrasing is to directly use a 124
074 model trained to generate paraphrase pairs (Gao 125
075 et al., 2020). Back-translation serves as another 126
076 type of paraphrasing, which first translates a sen- 127
077 tence into another language and then back again 128
078 (Chadha and Sood, 2019; Bornea et al., 2021).
079 Back-translation ensures the quality and correct-
080 ness of the augmented data and often shows im-
081 provement in downstream models. Both methods
082 aim to provide variety to the training data without
083 greatly altering the semantics of the original sen-
084 tences. However, these methods only operate on
085 the existing dialogue data and fail to take advantage
086 of the available document for augmentation.

087 Another direction for data augmentation is to 129
088 generate examples from scratch by grounding to 130
089 auxiliary documentation. Lewis et al. (2021) gener- 131
090 ate question-answer pairs with a model pre-trained 132
091 on available training data. This often requires ad- 133
092 ditional filtering or denoising measures to ensure 134
093 correctness of generated data. Also, these models 135
094 are built for the purposes of single-turn question 136
095 answering, rather than multi-turn dialogues. 137

096 Inspired by Alberti et al. (2019), we propose 138
097 an automatic document-grounded dialogue gener- 139
098 ation (DG^2) method that augments the amount of 140
099 data available for training a dialogue system. The 141
100 model consists of a user bot and an agent bot that 142
101 alternately generates utterances to complete a con- 143
102 versation. The user bot includes a span extraction 144
103 model that can first select a passage and then pre- 145
104 dict the rationale start and end positions inside a 146
105 passage. The agent bot has a denoising mechanism 147
106 to filter out generated rationales irrelevant to the 148
107 conversation. The user bot begins by selecting a 149
108 passage from the document that is most relevant to 150
109 the current context. It then selects a rationale span 151
110 from this passage and generates the user utterance. 152
111 The agent bot takes the selected span from the user 153
112 bot, and then checks if it can find the correct ratio- 154
113 nale span, and finally generates the agent response. 155
114 This process repeats until an entire dialogue is gen- 156
115 erated. 157

116 We evaluate our model on a representative 117
118 document-grounded dialog dataset Doc2Dial (Feng 119
120 et al., 2020a). We test and generate additional 121
122 dialogs with both the seen documents and un- 123
124 seen documents. We augment the original dataset 124
125 and train it on a downstream model. The results 125
126 show that our method improves the performance 126
127 of the downstream model after augmentation. We 127
128 also test scenarios of low-resource settings. We 128
129 train and evaluate the generative models with only 129
130 25%, 50%, 75% data. Experimental results show 130
131 that our method perform well even when training 131
132 data is scarce. 132

2 Related Work 129

2.1 Document Grounded Dialogue Systems 130

133 Document Grounded Dialogue System (DGDS) 134
134 is the type of dialogue systems that the dialogues 135
135 are grounded on the given documents. It helps 136
136 humans to better retrieve information they want as 137
137 most of human knowledge is stored in the form of 138
138 documents. The study of DGDS can greatly impact 139
139 the future way of interacting with knowledge. 140

141 Recently, there are many document grounded 141
142 dialogue datasets proposed. Doc2Dial (Feng et al., 142
143 2020b) is a representative document grounded di- 143
144 alogue dataset which involved human-to-human 144
145 conversations and focused on real scenarios under 145
146 social welfare domains. Previous datasets such as 146
147 CoQA (Reddy et al., 2019) and QuAC (Choi et al., 147
148 2018) focused on machine reading comprehensions. 148
149 SharC (Saeidi et al., 2018) is close to Doc2Dial. Its 149
150 conversations are grounded to short text snippets, 150
151 and contains follow-up questions. ABCD (Chen 151
152 et al., 2021) supports customer service interactions 152
153 by providing Agent Guidelines as additional docu- 153
154 mentation to aid in task-oriented conversations. 154

155 An example of DGDS from Doc2Dial is shown 155
156 in Figure 1. For each turn, the agent needs to look 156
157 at the specific paragraph inside the document to be 157
158 capable of answering the user’s questions. More- 158
159 over, the agent can also ask follow-up questions. 159
160 For A3, the agent asks “Would you like to know if 160
161 you are eligible?”. In this way, the agent guides the 161
162 user to center more on the details in the document. 162

2.2 Data Augmentation 162

163 Data augmentation for question answering and di- 163
164 alogue systems has been well-studied in the past. 164

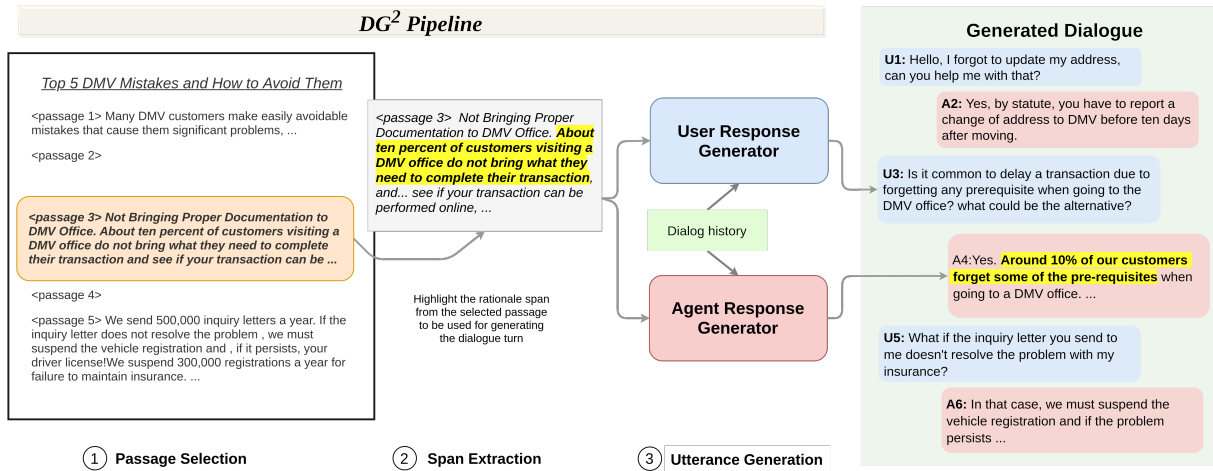


Figure 2: Overall pipeline of DG^2 . Given a document and the dialogue history, DG^2 iteratively performs (1) passage selection, (2) rationale extraction, and (3) utterance generation to produce a completed dialogue.

There are two major directions: paraphrasing existing QA pairs from seed data or generating new QA pairs from scratch.

Paraphrasing is a simple and effective technique to augment natural language datasets. It has been widely used in many NLP tasks including natural language understanding, question answering, and task-oriented dialog systems (Gao et al., 2020) to improve the downstream models’ performance. In question answering, paraphrasing with back-translation (Chadha and Sood, 2019; Bornea et al., 2021) is well-studied for datasets such as SQUAD (Rajpurkar et al., 2016b).

Another approach is generating new question-answer pairs. Early question-answer generation models used rule-based methods (Rajpurkar et al., 2016b). More recently, neural network-based question-answer pair generation models have been studied. PAQ (Lewis et al., 2021) generates 65 million question-answer pairs based on Wikipedia and trained a retrieval based that.

However, existing approaches have not explored applications for conversational question answering yet, especially for document grounded dialog systems. Compared to single-turn question answering datasets like SQUAD (Rajpurkar et al., 2016b), it involves additional complexity of modeling dialog flow and interconnection naturalness. Also, instead of only providing an answer span, datasets like Doc2Dial (Feng et al., 2020b) have free-form agent responses. The agent needs to produce natural utterances conditional to the selected rationale.

Also, existing conversational question generation models (Gu et al., 2021) only focused on

the quality of generations but did not address the improvement on downstream models. We design a specific dialog augmentation approach for document-grounded dialog systems. Our work can synthesize the entire conversation, and can be used to improve down-stream task’s performance.

3 Document-Grounded Dialogue Setup

A dialogue can be thought of as a series of turns between two interlocutors. Within goal-oriented dialogues, we refer to the first speaker as the user, and the second speaker as the agent, whom we model as $d = [(u_1, a_1), (u_2, a_2), \dots, (u_t, a_t)]$. In a document-grounded setting, the conversation revolves around the topics and entities mentioned in the associated document. A document is composed of a series of text passages, which are themselves broken down further into spans.

Dialogue success is determined by following the typical success metrics for any given task, where the only difference is that the outcome of the conversation is likely to depend on the ability to reason about the contents of the document. While sophisticated architectures are certainly capable of improving document-grounding, we take a data-centric approach instead by generating new dialogues from the documents to serve as additional training data for the downstream model.

4 Data Augmentation via DG^2

We propose Document-Grounded Dialogue Generation (DG^2) as a method of data augmentation. We aim to generate a complete and coherent dialogue given a document by building two bots

231 talking to each other.

232 Given a document C , we can model a dialog d
233 between the user and the agent with:

$$234 \quad p(d|C) = \prod_{i=1}^t p(u_i, a_i | c_i \in C) \quad (1)$$

235 where u_i is the user turn utterance, a_i is the agent
236 turn utterance, and c_i is the selected passage at i -th
237 turn.

238 We further decompose the model into three parts:
239 passage selection, rationale extraction, and utter-
240 ance generation. We also apply a filtering model to
241 ensure the quality of generated utterances.

242 4.1 Passage Selection

243 A document can often be very long, so it must
244 be divided into smaller passages first. Then, we
245 need to rank the passages, and select a relevant
246 passage given the dialogue context. We can maxi-
247 mize the passage probability for c_t with contrastive
248 loss where the positive passages are from ground
249 truth, and the negative passages are from the same
250 document.

$$251 \quad p(c_t | \{u_i, a_i\}_{i < t}, C) \quad (2)$$

252 During generation, we sample from the probability
253 distribution to select the passage. We choose to
254 sample rather than perform greedy selection since
255 this allows for choosing different passages given
256 the same dialogue context, thereby increasing the
257 diversity of the augmentation.

258 4.2 Rationale Extraction

259 Next, we further extract a rationale span from the
260 selected passage.

$$261 \quad p(r_t | \{u_i, a_i\}_{i < t}, c_t)$$

262 Span extraction systems typically model the start
263 and end position of a span independently as
264 $p(r_{\text{start}}|c) \times p(r_{\text{end}}|c)$. This settings works well when
265 the span is short, as is often the case for stan-
266 dard question answering tasks. However, the spans
267 encountered in some document-grounded dialog
268 datasets are much longer causing problems in tra-
269 ditional approaches. As an alternative, we propose
270 an autoregressive method that samples the start and
271 end position in sequentially with:

$$272 \quad p(r_t) = p(r_{\text{start}}|c) \times p(r_{\text{end}}|r_{\text{start}}, c) \quad (3)$$

To ensure that the autoregressive property holds, we
273 add the predicted start position’s hidden state H_{start}
274 and each position’s hidden state H_i , and then we
275 project the combined hidden state with a learnable
276 function f_r to get the final predicted end position.
277 Thus, the training objective becomes to maximize
278

$$279 \quad r_{\text{end}} = \arg \max_i f_r(H_{\text{start}} + H_i) \quad (4)$$

280 When extracting a rationale, we first sample a start
281 position from top-k options. Conditioned on this
282 start index, we then sample the end position. This
283 allows us to extract different rationales given the
284 same context, which greatly improves the diversity
285 of generated dialogues compared to using the same
286 rationale.

287 4.3 Utterance Generation

288 Given the selected passage and the extracted ratio-
289 nale, we can now start to generate the user utterance
290 and the agent utterance.

291 **User Utterance** As seen in Figure 2, user model
292 generates a user utterance conditioned on the di-
293 alog history and the extracted rationale. Instead
294 of only using the rationale to generate utterances,
295 we provide the context passage along with the ra-
296 tionale for better performance. To tell the model
297 where the rationale is in the passage, we highlight
298 the rationale span by wrapping its text in the in-
299 put with “[” and “]”. The new passage with the
300 rationale span information is defined as c'_t .

301 We then model the user utterance with a encoder-
302 decoder where the input is the dialogue history and
303 the passage c'_t , and the output is the user utterance.
304

$$305 \quad p(u_t) = p(u_t | \{u_i, a_i\}_{i < t}, c'_t) \quad (5)$$

306 **Agent Utterance** Similar to user utterance gener-
307 ation, we model the agent utterance with a encoder-
308 decoder.

$$309 \quad p(a_t) = p(a_t | \{a_i, u_i\}_{i < t}, c'_t) \quad (6)$$

310 The difference is that the dialogue history now
311 includes the previous generated user utterance. The
312 rationale position information in the passage is pro-
313 cessed similarly as in user utterance generation. We
314 can repeat the user utterance and agent utterance
315 generation process to generate the entire dialogue.

316 4.4 Filtering the Augmented Data

317 Roundtrip consistency checking (Alberti et al.,
318 2019; Zhong et al., 2020) has previously been used

to improve the correctness of generated augmentation data. It utilizes a model to double-check whether the answer span is the same as the span used to generate the question. Based on this insight, rather than tuning a sampling temperature to trade-off against noise and diversity, we instead greedily pick the rationale span and use consistency checking to filter for quality. For our purposes, we expect the extracted rationale to be aligned with the dialogue context as well as the user utterance.

We build a new passage selector and rationale extraction model such that:

$$p(\hat{c}_t | \{u_i, a_i\}_{i < t}, u_t, C) \quad (7)$$

$$p(\hat{r}_t | \{u_i, a_i\}_{i < t}, u_t, \hat{c}_t) \quad (8)$$

where \hat{c}_t is the predicted passage from the document C with the dialogue context and the generated user utterance, and \hat{r}_t is the prediction rationale within \hat{c}_t . When \hat{r}_t is not aligned the previous r_t , we remove this utterance u_t . Because rationale spans can be very long, filtering based on exact match will be too strict, so we filter based on f1 word overlap.

4.5 Document Positional Information

When a document is divided into passages, it loses positional information between different passages. As a dialogue progresses, we can expect to focus more on the later part of a document, which involves more details of a topic. Therefore, it is important to incorporate the turn information and the passage position information into the model.

We use a simple yet effective method to combine the dialogue turn positional information and passage positional information. For the speaker positions we use a prompt “user{num}:" or “agent{num}:", where “num” is replaced with the number of turns so far. This allows the model to track how many turns have passed, leading to a more coherent dialog structure. For the passage positions, we embed a passage index to indicate the location of the passage within the document. Combining the two flows together, the model is able to have conversations focused on the beginning of the document at the first, and naturally shift towards the end of document later.

5 Experiments

We first introduce the datasets evaluated with our method, then the baselines for comparisons, and in the end our method’s implementation details.

5.1 Datasets

	Dialogue Level				Document Level	
	#dial	#turns	#tok	span	#doc	#tok
train	3,474	11.8	15.0	26.5	415	834
valid	661	12.1	15.3	25.8	273	821
test	661	12.0	14.9	24.5	273	809
DG^2	3,474	12.0	14.2	42.2	415	834

Table 2: Doc2Dial dataset statistics. The following abbreviations are made: ‘dial’ is short for dialogue, ‘tok’ is short for tokens, and ‘doc’ is short for documents.

Doc2Dial consists of two subtasks around identifying relevant spans based on dialogue context and producing cohesive responses based on extracted rationales (Feng et al., 2020a). Formulated as a span selection task, user utterance understanding requires an agent to interpret user queries in the context of the dialogue history and then select the relevant span from the associated document. Predicted spans are graded based on Exact match (EM) and F1-score. Exact match is when the predicted span exactly lines up with the actual span. F1-score balances the recall and precision of the predicted uni-grams compared to the gold span.

The second subtask is agent response prediction, which requires an agent to generate a natural language response to the user query given the dialogue context and the document. Response quality is measured by SacreBLEU metric (Post, 2018) which aims to capture how closely the predicted response lines up with the gold response. Table 2 shows Doc2Dial’s dialogue-level statistics and document-level statistics.

5.2 Baselines

We compare against a number of baselines typically used to augment natural language data. In contrast to our technique, these methods all operate on the existing dialogues, whereas our method generates new dialogues from scratch from the associated document.

Easy Data Augmentation Wei and Zou (2019) propose to augment data through a series of surface form alterations. In particular, Easy Data Augmentation (EDA) consists of inserting new tokens, deleting random tokens, swapping pairs of tokens, or replacing tokens with their synonyms.

Back-translation Back-translation is another strong augmentation method which first translates

Model	Validation			Test			Span Coverage
	EM	F1	BLEU	EM	F1	BLEU	
Original data	58.13	72.61	37.08	58.34	73.25	36.89	48.27
+ EDA	60.40	74.30	37.72	59.71	73.62	37.63	48.27*
+ Back-translation	60.15	73.74	36.68	60.17	73.35	37.32	48.27*
+ Paraphrase	59.97	73.92	37.76	57.98	72.71	38.40	48.27*
+ DG^2	60.30	74.34	38.07	60.92	74.53	38.57	57.65

Table 1: Experimental results on the Doc2Dial dataset. EM stands for Exact Match. DG^2 outperforms all other data augmentation methods on almost every metric. *EDA, Back-translation, and Paraphrase do not modify span information and thus are unable to increase span coverage in relation to the original data.

some text into a separate language and then back-translates to the original language. We follow BERT-QA (Chadha and Sood, 2019), in translating all user utterances to French and then back to English to augment the original dialogues.

Paraphrase Paraphrasing can be achieved by training a sequence-to-sequence model on parallel paraphrase pairs corpora. In particular, we train a BART-base model (Lewis et al., 2020a) on the MRPC (Dolan and Brockett, 2005), QQP (Iyer et al., 2017) and PAWS (Zhang et al., 2019) datasets.

5.3 Coverage Metric

During inference, any section within the document is fair game for discussion. A model trained on dialogues that cover larger portions of the given documents should therefore perform better later on. Consequently, a strong data augmentation method should aim to generate dialogues that cover as much of the document as possible. We formalize this intuition with the span coverage metric, which we calculate as:

$$\text{Coverage} = \frac{\sum_{\text{span}} |\bigcup_{d \in \text{doc}_i} \bigcup_{s \in d} s|}{|\text{document}_i|}$$

where s refers to spans within a document and doc refers to the number of documents in the corpus.

5.4 Implementation Details

For passage ranker, and rationale extraction model, we fine-tuned RoBERTa-base (Liu et al., 2019) on the downstream training datasets. For utterance generators, we fine-tuned BART-base (Lewis et al., 2020b). We set total input length of 512-tokens which is 128 tokens for dialogue followed by 360 tokens for the document, with some room left over

for special tokens. The augmented data is generated with beam size 4, top-p 0.9, and temperature 0.9. When utilizing the augmented data, we pre-trained the downstream model on the augmented data for one epoch before fine-tuning (Alberti et al., 2019). The default f1 threshold is set to 0.9, which we determined by validating against the dev set. For fine-tuning, we train for five epochs, and use the same optimizer of AdamW (Loshchilov and Hutter, 2019) and learning rate of $3e - 5$ for all experiments.

6 Results and Analysis

This section shows the results for the full dataset and low-resource settings. We also conduct human evaluation on the generated dialogues. Afterwards, we discuss the results by analyzing generated examples.

6.1 Main Results

As shown in Table 1, DG^2 achieves the overall best performance compared to other baselines that only augment the original human-annotated data. Other baselines all show some improvements over the downstream model only trained using the original data. EDA has very high EM and F1 scores for the rationale extraction task, but suffers at producing coherent dialogues as measured by BLEU. Paraphrase has relatively lower EM and F1 scores, but it achieves better BLEU scores than EDA and Back-translation. We suspect that this is because Paraphrase contains more diverse utterances as the inputs than other baselines.

When evaluating the augmented dialogues with the original training set’s documents, we find that DG^2 achieves higher span coverage. Unlike the other methods, DG^2 is able to generate novel rationales to increase the diversity of the augmented data, which we believe plays a large factor in im-

Speaker	Utterance
User	I'm thinking about getting some retirement benefits.
Agent	Are you wondering what age to start your benefits? <i>Rationale: What Is The Best Age To Start Your Benefits?</i>
User	Yes, I want to know what age I should start my benefits.
Agent	Are you planning to apply for Medicare? <i>Rationale: Should I apply for Medicare?</i>
User	Yes, I am.
Agent	Then you should know that Medicare usually starts when you reach age 65. <i>Rationale: Remember, Medicare usually starts when you reach age 65. If you decide to delay starting your benefits, be sure to contact Social Security about 3 months before you turn age 65 to check about applying for Medicare. Even if you have health insurance through a ...</i>
User	What if I need to enroll but don't?
Agent	Your Medicare coverage may be delayed and cost more. <i>Rationale: your Medicare coverage may be delayed and cost more. Please read the general and special enrollment period information on page 12 of our Medicare booklet to find out what may happen if you delay.</i>
User	What age is the best age to begin benefits?
Agent	The answer is that there is no one "best age" for everyone and, ultimately, it is your choice. You should make an informed decision about when to apply for benefits based on your individual and family circumstances. <i>Rationale: The answer is that there is no one "best age" for everyone and, ultimately, it is your choice. You should make an informed decision about when to apply for benefits based on your individual and family circumstances.</i>

Table 4: Generated dialogue showing typical agent and user interaction. Unlike question answering datasets, the agent can also ask questions to the user to guide the direction of the conversation.

475 proving downstream metrics.

Filtering	#Spans	EM	F1
None	-	57.78	73.27
f1 < 0.5	top-1	57.73	73.01
f1 < 0.9	top-10	58.23	73.05
f1 < 0.9	top-1	60.80	74.38
f1 < 0.95	top-1	59.21	74.00
f1 < 0.98	top-1	59.26	73.84

Table 5: We test different quality thresholds to determine the optimal level of filtering. A higher F1-score means that more samples are filtered.

476 6.2 Low Resource Setting

477 To further illustrate the performance of DG^2 , we
478 train all the models with only 25%, 50%, 75% of
479 the original training data. We generate the dia-
480 logues based on the documents in the knowledge
481 base. In this limited data setting, our model gen-
482 erally outperformed Back-translation. However,
483 compared to EDA, there is still some performance
484 gap. We suspect that this is because when training
485 with less data, the generative models' performance

degenerates faster than the downstream model. We
486 hope to overcome these issues with further improve-
487 ments on data quality filtering. 488

489 6.3 Different Filtering Strategies

Prior works in data augmentation have shown that
490 filtering the synthetically generated examples can
491 provide a meaningful boost in the data quality
492 (Chen and Yu, 2021). As a result, we tune against
493 different F1-score thresholds and span counts on
494 the validation set. When the generated dialogue
495 produces a higher F1-score, then this example is
496 more likely to also produce better results during
497 testing. The span count determines how many ex-
498 amples we consider when calculating this score.
499 While raising the F1-score threshold increases the
500 potential quality of the data, it comes at the expense
501 of keeping fewer of the generated examples. Based
502 on Table 5, we observe a sweet spot at 0.9, where
503 a stricter filtering process would remove too many
504 examples while a looser filtering process would
505 lower the quality too much. 506

Model	25%			50%			75%		
	EM	F1	BLEU	EM	F1	BLEU	EM	F1	BLEU
Baseline	43.08	64.01	32.76	41.61	62.25	34.35	58.03	72.61	36.48
+ EDA	46.68	64.68	33.97	56.09	70.51	35.84	59.84	73.40	36.24
+ Back-translation	47.48	65.18	33.00	54.44	69.52	35.30	58.66	72.75	36.08
+ DG^2	46.48	65.58	32.90	54.51	71.40	35.74	58.89	73.38	37.01

Table 6: Experimental results on low-resource settings.

	Consistency	Fluency	Naturalness	Overall
Human	3.80	3.96	3.56	3.70
DG^2	3.60	4.18	2.98	3.38

Table 7: Human evaluation results on the generated dialogues.

6.4 Human Evaluation

We conduct human evaluation on the human dialogues and the generated dialogues. We randomly sample 50 dialogues from each class. We shuffled the sampled dialogues and ask annotators to rate the dialogues with a score 1-5 in four different aspects: consistency, fluency, naturalness, and overall quality.

From the evaluation results, the generated dialogues show better fluency than original human dialogues with $p < 0.05$. We observed that some human dialogues contain typos and grammar errors, while the generated dialogues are more grammatically correct, which explains the score difference. In terms of other human evaluation metrics, the generated dialogues are still worse than the original human dialogues.

6.5 Qualitative Analysis

We now compare and contrast two examples generated by our procedure. Table 4 shows a good example from the document-grounded dialogue dataset. In the first four turns, the agent guides the user’s focus by asking relevant questions. When the user wants to know more details, the agent then switches to provide the relevant knowledge retrieved from the rationale. This behavior is different from traditional question answering datasets where the agent simply reacts to user requests rather than exhibiting proactive behavior. On the flip side, one major problem of the current approach is repetition. The user continues to ask about forgetting to update their address despite attempts by the agent to answer their query. Although the surface form of the

user utterances are different, the semantic meaning remains the same. This repetition confuses the agent who then extracts irrelevant rationales, further exacerbating the situation.

7 Ethical Consideration

The models and approaches introduced in our work involve using synthetic data as an enhancement to existing datasets for modeling document-grounded dialogue. For the existing datasets, they are often dialogue simulation data generated by human workers based on their understanding of the associated document content and dialogue context. There are potential biases or toxic content introduced in the existing simulation during data collection. We can address such concerns by making efforts to improve the quality of the generated data that has shown its effectiveness in the downstream task. Therefore, our method can add an extra layer of safety and privacy if we only use generated data for training downstream models. Future work can explore how data augmentation can help to build a more private and safe dataset.

8 Conclusion

To address the problem of limited data in document-grounded dialogue systems, we propose DG^2 to perform data augmentation via dialogue generation. Our technique generates diverse utterances grounded on the given document while filtering the utterances to ensure quality and correctness when training on the downstream model. We demonstrated the effectiveness of our pipeline by showing the improvement over the previous data augmentation methods. We additionally show competitive results in the low-resource setting when a limited amount of human annotated data is available for training. Future work will explore more techniques of filtering to improve data quality. We hope this spurs further research into document-grounded augmentation techniques for dialogue systems.

References

Daniel Adiwardana, Minh-Thang Luong, David R. So, Jamie Hall, Noah Fiedel, Romal Thoppilan, Zi Yang, Apoorv Kulshreshtha, Gaurav Nemade, Yifeng Lu, and Quoc V. Le. 2020. [Towards a human-like open-domain chatbot](#). *CoRR*, abs/2001.09977.

Chris Alberti, Daniel Andor, Emily Pitler, Jacob Devlin, and Michael Collins. 2019. [Synthetic QA corpora generation with roundtrip consistency](#). In *Proceedings of the 57th Conference of the Association for Computational Linguistics, ACL 2019, Florence, Italy, July 28- August 2, 2019, Volume 1: Long Papers*, pages 6168–6173. Association for Computational Linguistics.

Mihaela A. Bornea, Lin Pan, Sara Rosenthal, Radu Florian, and Avirup Sil. 2021. [Multilingual transfer learning for QA using translation as data augmentation](#). In *Thirty-Fifth AAAI Conference on Artificial Intelligence, AAAI 2021, Thirty-Third Conference on Innovative Applications of Artificial Intelligence, IAAI 2021, The Eleventh Symposium on Educational Advances in Artificial Intelligence, EAAI 2021, Virtual Event, February 2-9, 2021*, pages 12583–12591. AAAI Press.

Ankit Chadha and Rewa Sood. 2019. [BERTQA - attention on steroids](#). *CoRR*, abs/1912.10435.

Derek Chen, Howard Chen, Yi Yang, Alexander Lin, and Zhou Yu. 2021. [Action-based conversations dataset: A corpus for building more in-depth task-oriented dialogue systems](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2021, Online, June 6-11, 2021*, pages 3002–3017. Association for Computational Linguistics.

Derek Chen and Zhou Yu. 2021. [GOLD: improving out-of-scope detection in dialogues using data augmentation](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic, 7-11 November, 2021*, pages 429–442. Association for Computational Linguistics.

Eunsol Choi, He He, Mohit Iyyer, Mark Yatskar, Wentau Yih, Yejin Choi, Percy Liang, and Luke Zettlemoyer. 2018. [Quac: Question answering in context](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Brussels, Belgium, October 31 - November 4, 2018*, pages 2174–2184. Association for Computational Linguistics.

William B. Dolan and Chris Brockett. 2005. [Automatically constructing a corpus of sentential paraphrases](#). In *Proceedings of the Third International Workshop on Paraphrasing, IWP@IJCNLP 2005, Jeju Island, Korea, October 2005, 2005*. Asian Federation of Natural Language Processing.

Song Feng, Kshitij P. Fadnis, Q. Vera Liao, and Luis A. Lastras. 2020a. [Doc2dial: A framework for dialogue composition grounded in documents](#). In *The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020*, pages 13604–13605. AAAI Press.

Song Feng, Hui Wan, R. Chulaka Gunasekara, Siva Sankalp Patel, Sachindra Joshi, and Luis A. Lastras. 2020b. [doc2dial: A goal-oriented document-grounded dialogue dataset](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing, EMNLP 2020, Online, November 16-20, 2020*, pages 8118–8128. Association for Computational Linguistics.

Silin Gao, Yichi Zhang, Zhijian Ou, and Zhou Yu. 2020. [Paraphrase augmented task-oriented dialog generation](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, ACL 2020, Online, July 5-10, 2020*, pages 639–649. Association for Computational Linguistics.

Jing Gu, Mostafa Mirshekari, Zhou Yu, and Aaron Sisto. 2021. [Chaincqq: Flow-aware conversational question generation](#). In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume, EACL 2021, Online, April 19 - 23, 2021*, pages 2061–2070. Association for Computational Linguistics.

Ehsan Hosseini-Asl, Bryan McCann, Chien-Sheng Wu, Semih Yavuz, and Richard Socher. 2020. [A simple language model for task-oriented dialogue](#). In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*.

Shankar Iyer, Nikhil Dandekar, and Kornel Csernai. 2017. [First quora dataset release: Question pairs](#). *Kaggle Competition*.

Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Matthew Kelcey, Jacob Devlin, Kenton Lee, Kristina N. Toutanova, Llion Jones, Ming-Wei Chang, Andrew Dai, Jakob Uszkoreit, Quoc Le, and Slav Petrov. 2019. [Natural questions: a benchmark for question answering research](#). *Transactions of the Association of Computational Linguistics*.

Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020a. [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, ACL 2020, Online, July 5-10, 2020*, pages 7871–7880. Association for Computational Linguistics.

809 Yuan Zhang, Jason Baldridge, and Luheng He. 2019.
810 PAWS: paraphrase adversaries from word scrambling.
811 In *Proceedings of the 2019 Conference of the North*
812 *American Chapter of the Association for Computa-*
813 *tional Linguistics: Human Language Technologies,*
814 *NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7,*
815 *2019, Volume 1 (Long and Short Papers)*, pages 1298–
816 1308. Association for Computational Linguistics.

817 Victor Zhong, Mike Lewis, Sida I. Wang, and Luke
818 Zettlemoyer. 2020. Grounded adaptation for zero-
819 shot executable semantic parsing. In *Proceedings of*
820 *the 2020 Conference on Empirical Methods in Nat-*
821 *ural Language Processing, EMNLP 2020, Online,*
822 *November 16-20, 2020*, pages 6869–6882. Associa-
823 tion for Computational Linguistics.

A Training Details

	Hyper-parameters
batch size	16
warm-up steps	1,000
learning rate	3e-5
dropout	0.1
weight decay	0.01
max gradient norm	1.0
training epochs	5

Table 8: Training Details

825 We trained our model on and A4000 and A6000
826 GPUs. The detailed hyper-parameters can be found
827 in the table above.

B Bad Generation Examples

828 We show a bad generation example here to illustrate
829 one of the generation problem.
830

Speaker	Utterance
User	I was planning to move and forgot to update my address. What do I do?
Agent	By statute, you must report a change of address to DMV within ten days of moving. That is the case for the address associated with your license, as well as all the addresses associated with each registered vehicle, which may differ <i>Rationale: By statute, you must report a change of address to DMV within ten days of moving. That is the case for the address associated with your license, as well as all the addresses associated with each registered vehicle, which may differ.</i>
User	Am I able to do anything if I forget to update my address?
Agent	It is not sufficient to only: write your new address on the back of your old license; tell the United States Postal Service; or inform the police officer writing you a ticket. <i>Rationale: It is not sufficient to only: write your new address on the back of your old license; tell the United States Postal Service; or inform the police officer writing you a ticket.</i>
User	What happens if I forget to update my address?
Agent	You will miss a suspension order and may be charged with operating an unregistered vehicle and/or aggravated unlicensed operation, both misdemeanors. This really happens, but the good news is this is a problem that is easily avoidable. <i>Rationale: you will miss a suspension order and may be charged with operating an unregistered vehicle and/or aggravated unlicensed operation, both misdemeanors. This really happens, but the good news is this is a problem that is easily avoidable. Learn more ...</i>

Table 9: A bad generated dialogue here shows repetition. Although the user utterances' surface form are different, their semantic meaning is the same.