

# 📖 Book2Dial: Generating Teacher Student Interactions from Textbooks for Cost-Effective Development of Educational Chatbots

Anonymous ACL submission

## Abstract

001 Educational chatbots are a promising tool for  
002 assisting student learning. However, the devel-  
003 opment of effective chatbots in education has  
004 been challenging, as high-quality data is sel-  
005 dom available in this domain. In this paper, we  
006 propose a framework for generating synthetic  
007 teacher-student interactions grounded in a set  
008 of textbooks. Our approaches capture a key as-  
009 pect of learning interactions where curious stu-  
010 dents with partial knowledge interactively ask  
011 teachers questions about the material in the text-  
012 book. We highlight various quality criteria that  
013 such dialogues must fulfill and compare sev-  
014 eral approaches relying on either prompting or  
015 finetuning large language models according to  
016 these criteria. We use the synthetic dialogues to  
017 train educational chatbots and show the benefits  
018 of further fine-tuning in educational domains.  
019 However, careful human evaluation shows that  
020 our best data synthesis method still suffers from  
021 hallucinations and tends to reiterate informa-  
022 tion from previous conversations. Our findings  
023 offer insights for future efforts in synthesizing  
024 conversational data that strikes a balance be-  
025 tween size and quality. We will open-source  
026 our data and code.

## 027 1 Introduction

028 Educational chatbots are a scalable way to improve  
029 learning outcomes among students (Kuhail et al.,  
030 2023). However, building educational chatbots  
031 has been challenging as high-quality data involv-  
032 ing teachers and students is difficult to obtain due  
033 to various practical reasons such as privacy con-  
034 cerns (Macina et al., 2023). In response to this,  
035 we study the task of generating synthetic teacher-  
036 student interactions from textbooks. We create a  
037 novel dataset of textbooks drawn from an open  
038 publisher of student textbooks and present a frame-  
039 work (📖 Book2Dial) to generate synthetic teacher-  
040 student interactions from these textbooks.

041 Our teacher-student interactions take the form

of conversational question-answering (QA) interac- 042  
tions (Choi et al., 2018; Reddy et al., 2019) where 043  
curious students ask teachers questions from the 044  
textbook and teachers answer these questions from 045  
the textbook. However, the task of generating high- 046  
quality synthetic data in the space of education is 047  
difficult (Kim et al., 2022a; Dai et al., 2022). Thus, 048  
it is important to have quality controls on such data, 049  
because students might otherwise receive wrong 050  
feedback, which could be detrimental to learning. 051

Thus, in this work, we also sketch various qual- 052  
ity requirements that measure the quality of educa- 053  
tional dialogues. For example, it is crucial that the 054  
chatbot does not provide students with incorrect 055  
information and stays grounded in the textbook, 056  
ensuring factual consistency with the knowledge 057  
taught. This is particularly important given that 058  
large language models (LLMs) are prone to 'hallu- 059  
cinations' or generating plausible but incorrect or 060  
unverified information (Rawte et al., 2023). While 061  
a simple teacher strategy would be to just answer 062  
with extracted passages from the textbook, this 063  
might hurt the coherence of the dialogue which is 064  
present in interactive educational situations (Baker 065  
et al., 2021). The teacher's response should both be 066  
relevant to the student's question (Ginzburg, 2010), 067  
as well as, informative as this ensures that key infor- 068  
mation from the textbook is covered in the dialogue 069  
(Tan et al., 2023). We formalize these requirements 070  
into 7 criteria, shown in Figure 1. 071

Our framework, 📖 Book2Dial, comprises of 072  
three approaches: multi-turn QG-QA (Kim et al., 073  
2022a), Dialogue Inpainting (Dai et al., 2022) and 074  
using role-playing abilities of LLMs to simulate 075  
teacher and student. We use the formatting infor- 076  
mation in the textbook, such as titles, key concepts, 077  
bold terms, etc to initialize student models with im- 078  
perfect information. In contrast, the teacher models 079  
have perfect information and are expected to gener- 080  
ate grounded responses based on the textbook. We 081  
fine-tune and prompt various open-source language 082

Formatting (C)		Textbook source text (S)						
<b>Subsection Title:</b> Planet <b>Key Concepts:</b> Sun, Earth, Mars <b>Learning Objectives:</b> Learn about Planets <b>Summary:</b> The Sun is the center of the solar system, Earth is ...		The Sun as the center of the solar system. <b>Earth</b> , the third planet from the Sun, <b>with one moon</b> . Mars, known for its red color, having two moons, Phobos and Deimos.						
	Answer Relevance	Coherence	Informativeness	Groundedness	Answerability	Factual Consistency	Specificity	
<b>Student:</b> What is the color of Mars? <b>Teacher:</b> Mars has moons.	✗	NA	✓	✓	✓	✗	✓	
<b>Student:</b> How many moons does it have? <b>Teacher:</b> I don't know how many moons Mars has.	✓	✓	✗	✗	✓	✓	✓	
<b>Student:</b> What is interesting about this passage? <b>Teacher:</b> Sun is the center of solar system.	✗	✗	✓	✓	✓	✓	✗	
<b>Student:</b> How many moons does Earth have? <b>Teacher:</b> Earth has moons, it has two moons.	✓	✗	✓	✓	✓	✗	✓	
<b>Student:</b> Mars is red. <b>Teacher:</b> Mars is red.	✓	✗	✓	✓	✗	✗	✗	

Figure 1: Example of a synthetic teacher-student interaction based on a textbook, along with various criteria for evaluating the quality of the interaction. The criteria include Answer Relevance of the answer with respect to the question, Coherence of the question-answer interaction with respect to the history, Informativeness of the overall interaction, Groundedness to the textbook, Answerability of the question from the textbook, Factual Consistency of the answer with respect to the question, and Specificity of the question. More details in Section 3.2.

models to generate teacher-student interactions.

We evaluated Book2Dial on the proposed quality criteria and also used human evaluations to support our findings. Our results reveal that data generated by role-playing LLMs scores highest in most criteria, as shown in Section 5.1.1 and 5.1.2, demonstrating reasonable efficacy in creating educational dialogues. While the generated dialogues consistently contain information grounded in the textbooks, they still fail to mirror natural educational conversations. The dialogues suffer from issues like hallucination and a tendency to reiterate information in the previous conversations, as shown in Section 5.3. Yet, despite these limitations, we were able to show that the generated synthetic data can be used to pre-train and finetune educational chatbots with some benefit in various educational domains, as shown in Section 5.4.

## 2 Related Work

### 2.1 Synthetic Data for Conversational QA

Prior work in educational research has focused on generating individual questions (Kurdi et al., 2020) under two common settings: answer-aware and answer-unaware generation. The former approach starts by identifying an answer and then generates a question accordingly, whereas the latter generates a question without pre-determining the answer. These approaches have also been extended to generating multiple questions (Rathod et al., 2022), causal question generation (Stasaski et al., 2021), prediction of question types to ask (Do et al., 2023), or decomposing problems into Socratic subquestions (Shridhar et al., 2022). However, most works do not address conversational settings.

Datasets like QuAC (Choi et al., 2018; Qu et al., 2020) and CoQA (Reddy et al., 2019) focus on con-

versational question answering in non-educational settings. Previous work has also explored strategies for creating such data with humans or automatically by using models. For example, Qi et al. (2020) withholds the context required for answers from the questioner, leading to information-seeking questions. SimSeek (Kim et al., 2022a) synthesizes datasets for conversational question answering from unlabeled documents. However, it fails to demonstrate significantly improved performance in downstream tasks. A recent work, Dialogizer (Hwang et al., 2023), proposes a framework for generating context-aware conversational QA dialogues. However, these methods do not take into account the needs and considerations of the educational domain.

### 2.2 Educational Dialogue Datasets

Development of educational chatbots is highly reliant on quality data. Yet such data is hard to obtain. Therefore, previous works such as MathDial (Macina et al., 2023) collect conversational data by pairing real teachers with an LLM that simulates students. Other datasets are commonly created by roleplaying both teacher and student, such as CIMA (Stasaski et al., 2020) or by transcribing classrooms (Suresh et al., 2022; Demszky and Hill, 2022) or recording online conversations (Caines et al., 2020). However, all of these methods are challenging to scale, and using non-experts often leads to data quality issues (Macina et al., 2023).

Thus, in this work, we explore data synthesis as a scalable way of creating such data. Data augmentation and synthetic data generation have gained attention as effective techniques to overcome the challenges associated with manual data annotation. Synthetic data generation, particularly

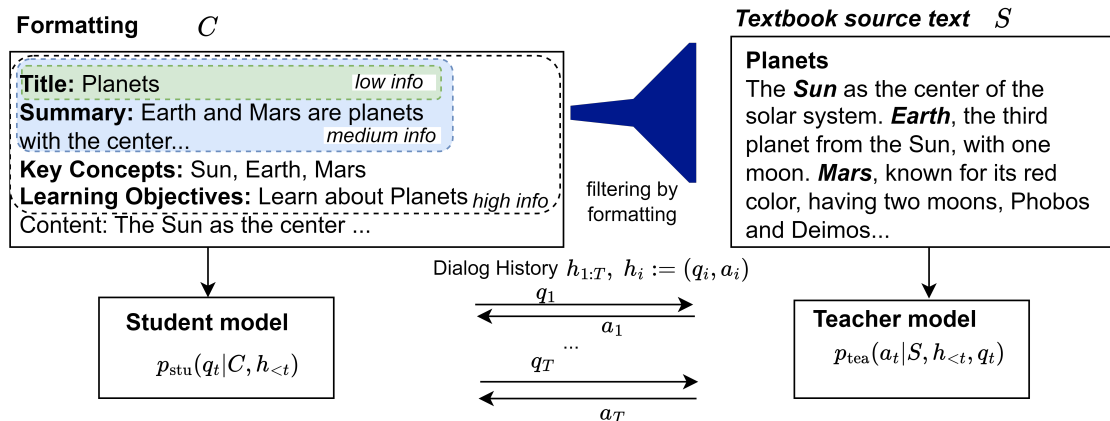


Figure 2: Book2Dial Framework for Generating Dialogues from Textbooks: Our approach uses two models – a Student model and a Teacher model. The Student model plays the role of a student, formulating questions from a limited context (document formatting). In contrast, the Teacher model assumes the role of a teacher, providing answers and guidance by referencing the (sub-)section in the textbook. This framework can be adapted to various instantiations of the two roles with varying formatting information, such as multi-turn QA-QG models (Kim et al., 2022a), Dialogue Inpainting (Dai et al., 2022), and a new approach utilizing role-playing LLMs.

utilizing LLMs, has been shown to be a promising approach. For instance, Kim et al. (2022b) demonstrated the potential of sourcing dialogue data from common sense knowledge. However, ensuring the objectivity of this generated data remains a concern. Similarly, Zhang et al. (2018) introduced innovative methods for task-oriented dialogue synthesis. However, its dependency on predefined schemas limits its scalability.

### 3 Educational Conversation Generation

We first introduce a framework for dialogue synthesis from textbooks in Section 3.1, and then discuss the quality criteria that the generated dialogues should fulfill in Section 3.2.

#### 3.1 Book2Dial Framework

We set out to create meaningful teacher-student interactions from educational textbooks in the form of conversational QA pairs between the teacher and the student. In order to generate these interactions, we assume that the “teacher” is familiar with the textbook content, and the “student” only knows limited information from the textbook. Thus, we intuitively provide the **teacher model** all the textbook information but withhold some information from a **student model**. For this, we can use the structuring and formatting elements found in textbooks, including 1) **Titles**: headings of sections and subsections; 2) **Summary**: summaries of chapters; 3) **Other Metadata**: key concepts, learning objectives, bold terms, and the introductory paragraph of each section; and assume that the student

model only has access to this information.

During the conversation, the “student” asks inquisitive questions about the textbook while the “teacher” guides them by answering these questions and including additional information in their response. Formally, a dialogue  $d$  comprises of a sequence of  $T$  question-answer interactions:  $d = \{(q_1, a_1), \dots, (q_T, a_T)\}$ . The formalization of the task is depicted in Figure 2. The student model  $p_{\text{stu}}(q_t | C, h_{<t})$  generates a question  $q_t$  given the dialog history  $h_{<t} = \{(q_i, a_i)\}_{i=1}^{t-1}$  and the partial context (formatting information)  $C$ . The teacher model  $p_{\text{tea}}(a_t | S, h_{<t}, q_t)$  generates the answer response  $a_t$  given the question, the dialogue history and the full textbook source  $S$ . We call this framework Book2Dial.

#### 3.2 Evaluation of Educational Conversations

To build a high-quality conversation, we want the student to ask questions that are **specific** enough to drive the conversation forward, and also **answerable** given the context. The teacher must then respond with an answer that is **relevant** to the question, **factually consistent** with the context, and **informative** to the student. Finally, the overall conversation should be **coherent** and **grounded** to the entire context, not just parts of it. We use this as our guiding principle and define 7 criteria to evaluate the quality of a good educational interaction. We detail these criteria in the rest of this subsection.

##### 3.2.1 Answer Relevance

Answer Relevance measures how directly related the answer is to the question in each QA pair in the

218	dialogue. This criterion is important in education	267
219	as it ensures students get relevant responses to their	268
220	questions. In order to compute Answer Relevance,	269
221	we assess the Answer Relevance of individual QA	270
222	pairs and then combine these assessments to de-	271
223	termine the dialogue’s overall Answer Relevance.	272
224	We use <b>BF1</b> ( $q_t, a_t$ ) and <b>QuestEval</b> as automatic	
225	metrics for Answer Relevance, the BF1 computing	
226	BERTScore F1 (Zhang et al., 2019) for seman-	
227	tic alignment between question and answer using	
228	BERT’s embeddings, while QuestEval (Scialom	
229	et al., 2021) generates questions from both the ques-	
230	tion and answer, then generates answers for these	
231	questions, comparing them to measure Answer Rel-	
232	evance. More details in Table 6.	
233	<b>3.2.2 Coherence of the Dialog</b>	
234	Coherence measures whether QA pairs in the dia-	
235	logue form a logical and smooth whole, rather than	
236	independent QA pairs. Coherence is an important	
237	aspect of good dialogue (Dziri et al., 2019), it is	
238	important in education because it helps students	
239	connect new information to what is already taught.	
240	We adapt two metrics, <b>BF1</b> ( $q_t, a_{<t}$ ) and <b>BF1</b> ( $q_t,$	
241	$a_{(t-1)}$ ), to measure coherence. The first metric cal-	
242	culates the BERTScore F1 considering the current	
243	question as the predicted sentence and each of the	
244	previous answers as reference sentences, while the	
245	second uses BERTScore F1 to compare the current	
246	question only against the immediately preceding	
247	answer. More details in Table 6.	
248	<b>3.2.3 Informativeness</b>	
249	Informativeness evaluates the amount of new infor-	
250	mation introduced by each student-teacher interac-	
251	tion in the dialogue. This criterion is important in	
252	education because the dialogue should teach the	
253	student about diverse information discussed in the	
254	textbook, rather than repeating or paraphrasing pre-	
255	viously stated facts. We use the <b>1 - Overlap</b> ( $a_t,$	
256	$a_{<t}$ ) metric for evaluating Informativeness. For	
257	each QA pair, this metric is calculated as one minus	
258	the ratio of the intersection over the union of to-	
259	kens in the current and all previous answers. More	
260	details in Table 6.	
261	<b>3.2.4 Groundedness to the Textbook</b>	
262	This criterion assesses the amount of information	
263	from the textbook incorporated into the dialogue.	
264	This metric is crucial in education as it ensures	
265	that the dialogue is grounded on information from	
266	the textbook and does not hallucinate information.	
	Two metrics are used for assessment: <b>Density</b> , eval-	267
	uating the average length of text spans extracted	268
	from textbook content $S$ and included in the di-	269
	alogues; <b>Coverage</b> , measuring the proportion of	270
	dialogue words originating from the textbook. The	271
	formulas of these metrics are shown in Table 6.	272
	<b>3.2.5 Answerability of the Questions</b>	273
	Answerability measures whether the student’s ques-	274
	tion is answerable given the textbook content. This	275
	criterion is important in education as unanswer-	276
	able questions lead to unproductive dialogues. We	277
	use the “distilbert-base-cased-distilled-squad” QA	278
	model <sup>1</sup> to judge whether each question is answer-	279
	able given the textbook content, and refer to this	280
	metric as <b>Answerability</b> . More details in Table 6.	281
	<b>3.2.6 Factual Consistency of the Answer</b>	282
	Factual Consistency measures whether the answer	283
	correctly responds to the student’s question, a key	284
	criterion for ensuring students receive accurate in-	285
	formation. This evaluation only applies if the ques-	286
	tion is answerable from the textbook. Existing	287
	metrics like $Q^2$ (Honovich et al., 2021) use a QA	288
	model to assess answer correctness, while RQUGE	289
	(Mohammadshahi et al., 2022) uses a QA model	290
	to evaluate the quality of the candidate question.	291
	In our scenario, we need to measure whether the	292
	answer contains correct information and accurately	293
	answers the question. To meet this new require-	294
	ment, we build on the idea of $Q^2$ and introduce a	295
	new metric, which we will refer to as <b>QFactScore</b> :	296
	$\alpha \cdot \text{sim}(\text{QA}(q_t, S), a_t) + \beta \cdot \text{sim}(q_t, a_t) \quad (1)$	297
	For each QA pair, it computes the cosine similar-	298
	ity between the embeddings of the QA model’s	299
	predicted answer and the original answer. Then,	300
	it assesses the similarity between the embeddings	301
	of the question and answer. The final score is the	302
	weighted sum of two similarity scores. More de-	303
	tails on this metric can be found in Table 6, and	304
	Appendix A.4.	305
	<b>3.2.7 Specificity of the Question</b>	306
	Specificity assesses whether the question is spe-	307
	cific, and could be posed in any setting, regard-	308
	less of the educational context. An example of a	309
	generic question is ‘What is interesting about this	310
	passage?’. The Specificity criterion is crucial in	311
	education as generic questions, which carry limited	312

<sup>1</sup><https://huggingface.co/distilbert-base-cased-distilled-squad>

educational value, are less preferred. We assess specificity through human evaluation, as there is no existing metric that captures specificity.

## 4 From Textbooks to Dialogues

In this section, we describe different methods used for generating dialogues from educational textbooks in `Book2Dialog`, namely:

1. **Multi-turn QG-QA models:** In this setting, we use fine-tuned QG and QA models interacting with each other.
2. **Dialogue Inpainting** Dai et al. (2021) uses a span extraction model over the textbook as a teacher model, where the response is copied from the textbook and the question is generated by a QG model acting as the student.
3. **Persona-based Generation.** This approach uses LLMs like GPT-3.5, and leverages prompting to interactively simulate the student and the teacher and generate dialogues.

In the following, we describe how each of these methods are implemented. Further details can be found in Appendix A.3.

### 4.1 Multi-turn QG-QA models

This scenario utilizes separate QG and QA models to interact in a multi-turn scenario. As a representation of this approach from related work, we consider the SimSeek-asym model (Kim et al., 2022a). The approach consists of two components:

1. A **Question Generation** (QG) model for generating conversational questions relying solely on prior information (i.e., formatting information relevant to the topic). The model generates question based on the dialog history and filtered Information  $C$ :  $p(q_t|C, h_{<t})$ .
2. A **Conversational Answer Finder** (CAF) to comprehend the generated question and provide an acceptable answer to the question from the evidence passage:  $p(a_t|S, h_{<t}, q_t)$ .

### 4.2 Dialogue Inpainting

Dialogue Inpainting (Dai et al., 2022) is an approach for dialogue generation characterized by its information-symmetric setting. In this framework, both the student and teacher model are provided with the complete textbook text  $S$ . The teacher

model is a simple model iterating over each sentence in  $S$  and copying it as an answer. The student model is a QG model. We use data from the OR-QuAC (Qu et al., 2020), QReCC (Anantha et al., 2020), and Taskmaster-2 (Byrne et al., 2020) datasets to train the student model. For the student model, a dialogue reconstruction task is employed. At training time rather than distinguishing questions and answers, the dialog reconstruction task treats a conversation as a sequence of utterances  $\{u_i\}_{i=1}^{2T}$ . To train it, a randomly chosen utterance  $u_i$  is masked to create a partial dialogue  $d_{m(i)} = u_1, \dots, u_{i-1}, \langle \text{mask} \rangle, u_{i+1}, \dots, u_{2T}$ . The model then predicts  $u_i$  and is trained by minimizing the loss:

$$\mathcal{L}(\theta) = - \sum_{d \in D} \mathbb{E}_{u_i \sim d} [\log p_\theta(u_i | d_{m(i)})] \quad (2)$$

During inference, the model uses each sentence in the textbook as a teacher’s utterance and only predicts student utterances accordingly,  $\{u_{2k-1}\}_{k=1}^T$  corresponding to  $\{q_i\}_{i=1}^T$  in our notation. We basing our model (eq 2) on FLAN-T5-XL (Chung et al., 2022). Further details are elaborated in Appendix A.3.2.

### 4.3 Persona-based Generation

Inspired by (Markel et al., 2023)’s idea of using LLMs to simulate student personas, we propose a method to simulate student and teacher personas using LLMs for dialogue generation. We use one GPT-3.5 model to play the student and another to play the teacher.<sup>2</sup> The teacher model is provided with all the information from the textbook, including content and all the formatting information. The information provided to the student model is varied. We consider four variants for generating dialogue in each subsection based on the amount of information provided to the student model: 1) Persona (Low Info) provides the student model with only the Title information, 2) Persona (Medium Info) provides both the Title and Summary information, 3) Persona (High Info) offers all formatting information, and 4) Persona (Single Instance) uses a single prompt to generate the entire dialogue, it provides one model with formatting and textbook content information. The detail for each variant is introduced in Table 8.

<sup>2</sup>We used the GPT-3.5-turbo API between 25th September and 4th October, 2023.

402 Considering that GPT-3.5 is proprietary and not  
403 open-source, we adopted prompting techniques  
404 to steer the models in dialogue generation. The  
405 prompt for Persona (High Info) and Persona (Single  
406 Instance) is detailed in Appendix A.3.3.

## 407 5 Results and Analyses

408 In this section, we aim to address the following  
409 research questions:

- 410 1. How does the choice of generation framework  
411 influence the quality of the generated data?
- 412 2. What is the optimal amount of information  
413 that should be incorporated into the student  
414 model to produce natural dialogues?
- 415 3. Does pretraining on our synthetically gen-  
416 erated data improve the downstream perfor-  
417 mance of models that are finetuned on existing  
418 educational datasets?

419 To address these questions, we generate dia-  
420 logues from textbooks across various domains and  
421 analyze the generated dataset.

422 **Textbook data:** We collected 35 textbooks avail-  
423 able on OpenStax<sup>3</sup>, spanning domains of math,  
424 business, science, and social science. From these,  
425 we select four textbooks to create our dialogue  
426 datasets. Table 7 provides statistics of the four  
427 textbooks. The first and second research questions  
428 are addressed in Sections 5.1 and 5.2, respectively,  
429 while the third question is answered in Section 5.4.

### 430 5.1 Automatic Evaluation

431 In this section, we discuss statistics and metrics  
432 for the generated datasets. Here, we present the  
433 average results of the datasets generated from four  
434 domains of textbooks in Tables 1 and 2. The results  
435 for each specific domain in the dataset can be found  
436 in Tables 11 and 12. To compensate for the effects  
437 of the different number of turns in the generated  
438 dialogues, we set the maximum number of turns  $T$   
439 to 12 for each model, similar to (Kim et al., 2022a).

#### 440 5.1.1 Statistical Analysis of the Datasets

441 In dialogue, different types of questions emphasize  
442 various aspects. We hypothesize that "what" and  
443 "which" questions focus on factual information.  
444 In contrast, other question types, such as "why"  
445 and "how," tend to reflect more complex inquiries,  
446 which are also important in educational contexts.

<sup>3</sup><https://openstax.org/>

447 In Table 1, we present the percentages of student  
448 questions beginning with words *what*, *which*, *why*,  
449 and *how*. Furthermore, the average token count for  
450 questions and answers across each dataset is also  
451 shown. The key findings are as follows:

452 **Less factual questions in the Persona (Single  
453 Instance) dataset** The Persona (Single Instance)  
454 model generates the fewest "what" or "which" ques-  
455 tions, suggesting more diverse questioning than  
456 other methods.

457 **More "how" questions in SimSeek and Dialogue  
458 Inpainting** The SimSeek and Dialogue Inpaint-  
459 ing generated datasets have a high ratio of "how"  
460 questions compared to the Persona-based model.

461 **High token counts in Persona datasets** Datasets  
462 from Persona models have the highest average to-  
463 ken counts in questions and answers, which means  
464 these dialogues are more verbose and contain more  
465 information.

#### 466 5.1.2 Data Quality Metrics

467 We report the various data quality metrics in Table  
468 2. Our key findings are as follows:

469 **Persona datasets excel in most of the criteria**  
470 The datasets generated by Persona models outper-  
471 form others in terms of metrics for various criteria:  
472 Answer Relevance, Coherence, Answerability, and  
473 Factual Consistency suggesting that persona-based  
474 generation models are more suitable for generating  
475 dialogues from textbooks.

476 **High Informativeness and Groundedness of Di-  
477 alogue Inpainting dataset:** Dialogue Inpainting  
478 models achieve the highest score across all mod-  
479 els in Informativeness and Groundedness. This is  
480 expected as this model uses sentences in the text-  
481 books as teachers' answers. In summary, Dialogue  
482 Inpainting generates dialogues that, despite not be-  
483 ing high quality, cover most textbook information.

484 **Students with more information access perform  
485 better in automatic metrics.** Datasets from Per-  
486 sona (High Info) generally perform as well as or  
487 better than those with less formatting information  
488 in terms of Answer Relevance, Informativeness,  
489 Groundedness, Coherence, and Answerability. Per-  
490 sona (Medium Info) has the highest score in Factual  
491 Consistency. This suggests that more information  
492 to the student may enhance key criteria. However,  
493 the impact differences among formatting levels are

	Question Type			Num. Tokens	
	% what/which	% why	% how	Tokens in Questions	Tokens in Answers
SimSeek	51.25	1.50	16.50	10.90	14.33
Dialogue Inpainting	55.75	<b>3.00</b>	<b>17.25</b>	6.85	19.63
Persona (Single Instance)	11.75	1.70	9.25	14.97	34.58
Persona (Low Info)	55.50	0.03	0.78	17.56	84.75
Persona (Medium Info)	52.00	0.01	0.40	17.69	85.19
Persona (High Info)	<b>59.25</b>	0.09	0.35	19.01	84.70

Table 1: Key statistics of the synthesized educational dialogue dataset.

	Answer Relevance		Informativeness	Groundedness		Coherence		Answerability	Factual Consistency
	BF1 ( $q_t, a_t$ )	QuestEval	1 - Overlap ( $a_t, a_{<t}$ )	Density	Coverage	BF1 ( $q_t, a_{<t}$ )	BF1 ( $q_t, a_{t-1}$ )	Answerable	QFactScore
SimSeek	0.53	0.25	0.71	11.66	0.82	0.51	0.55	0.84	0.32
Dialogue Inpainting	0.52	0.28	<b>0.91</b>	<b>22.62</b>	<b>0.90</b>	0.45	0.46	0.75	0.24
Persona (Sing. Inst.)	0.58	0.35	0.86	3.94	0.75	0.49	0.52	0.92	0.54
Persona (Low Info)	0.61	<b>0.44</b>	0.59	2.39	0.70	0.52	<b>0.59</b>	0.98	0.75
Persona (Med. Info)	0.61	<b>0.44</b>	0.59	2.43	0.71	0.52	<b>0.59</b>	<b>0.99</b>	<b>0.76</b>
Persona (High Info)	<b>0.62</b>	<b>0.44</b>	0.60	2.50	0.71	<b>0.53</b>	<b>0.59</b>	<b>0.99</b>	0.75

Table 2: Quality metrics computed for the synthesized dialogue data. Higher values mean better data quality.

not markedly significant, indicating a need for further research on this question.

## 5.2 Human Evaluation

To compensate for the limitations of automatic metrics, we evaluated dialogues from SimSeek, Dialogue Inpainting, and Persona (High Info) using human evaluation based on seven criteria - Answer Relevance (AnsRel), Informativeness (Info), Groundedness (Gro), Coherence (Coh), Factual Consistency (Fact), Answerability (Ans), and Specificity (Spe). The questions for judging each criterion are in Table 13. We recruited 3 expert annotators to evaluate 12 dialogues each. We report an average Cohen’s Kappa of  $\kappa = 0.66$ , indicating substantial agreement. Evaluation details are in Appendix A.8, and results in Table 3 and 4.

	AnsRel	Info	Gro	Coh	Fact
SimSeek	0.36	0.58	<b>1.00</b>	0.63	0.25
Dial. Inpaint.	0.72	<b>1.00</b>	<b>1.00</b>	0.83	0.68
Persona (High Info)	<b>0.93</b>	0.64	0.92	<b>0.90</b>	<b>0.76</b>

Table 3: Human Evaluation Result 1: Persona (High Info) generated dialogues score highest in Answer Relevance, Coherence and Factual Consistency, while Dialogue Inpainting generated dialogues excel in Informativeness and Groundedness

Persona (High Info) excels among the three models, leading in Answer Relevance, Coherence, Factual Consistency, Answerability, and Specificity, rendering it the most suitable choice for our dialogue generation objectives. This result aligns with the results of automatic metrics presented in Table 2. However, the dialogues generated by

the Persona-based method exhibit only an average score in Informativeness, with a score of 0.64 indicating that approximately 36% of QA pairs fail to contribute new information. The Persona-based model, while leading in Factual Consistency among the three models, scores only 0.76, which indicates that approximately 24% of the QA pairs **lack Factual Consistency**. For educational dialogues, it’s imperative to aim for high Factual Consistency to ensure the reliability of the knowledge imparted. The primary reason for this issue is the hallucination in LLMs, where LLMs respond to questions using fabricated or false information not grounded in the textbook. This poses a significant challenge and calls for further research into ways to better ground LLMs to text documents in the future.

	Answerability	Specificity
SimSeek	0.65	0.87
Dialogue Inpainting	0.89	0.64
Persona (High Info)	<b>0.90</b>	<b>1.00</b>

Table 4: Human Evaluation Result 2: Persona (High Info) generated dialogues score highest in Answerability and Specificity

## 5.3 Qualitative human analysis

We further analyzed the dialogues generated by each model. We find:

**Repeating answers in SimSeek and Persona** In the SimSeek and Persona datasets, we find that teacher answers often reiterate information from previous interactions. SimSeek often generates questions related to the same textbook sentence,

541 while Persona often provides summaries of text-  
542 book content in each answer.

### 543 **Insufficient follow-up ability of Persona models**

544 Dialogues generated by Persona models are unlike  
545 natural conversations and resemble a series of QA  
546 pairs about textbooks. The dialogue does not have  
547 enough follow-up questions and does not go into  
548 depth about a certain aspect.

### 549 **Insufficient Specificity of Dialogue Inpainting**

550 In alignment with the results of human evaluation,  
551 we find that the Dialogue Inpainting model tends  
552 to generate “general” questions, such as “What is  
553 interesting about this passage?” These types of  
554 questions, which are not specific to the textbook  
555 content, are less desirable in educational dialogue.

## 556 **5.4 Pretraining for Educational Chatbots**

557 We verify the suitability of our synthesized data  
558 for training educational chatbots in this section.  
559 We use the synthetic data to pre-train simple ed-  
560 ucational chatbot models, and evaluate them on  
561 downstream educational conversation tasks.

562 Specifically, we use text generation models  
563 based on language models to generate teacher re-  
564 sponses  $a_t$  given the dialogue history  $h_{<t}$ , text-  
565 book grounding information  $S$  and the question  
566  $q_t$ . We compare two scenarios: (1) a model pre-  
567 trained on our synthetic datasets, then fine-tuned  
568 and tested on various educational or information-  
569 seeking dialogue datasets; and (2) a model trained  
570 and tested solely on these dialogue datasets without  
571 pretraining. We used FLAN-T5-LARGE (Chung  
572 et al., 2022) as our base language model. For our  
573 test sets, we use the MCTest and CNN splits of  
574 the CoQA dataset (Reddy et al., 2019), as well as  
575 the NCTE dataset (Demszky and Hill, 2022). The  
576 MCTest split contains dialogues about children’s  
577 stories; the CNN split contains conversations about  
578 the news; the NCTE dataset contains transcripts of  
579 elementary math classrooms.

580 We pretrained the base model on four textbook-  
581 based synthetic datasets, each from a different  
582 subject: math, business, science, and social sci-  
583 ence. The datasets and training details are shown  
584 in Appendix A.9. The results are shown in Ta-  
585 ble 5. We report the BLEU score<sup>4</sup> of the sce-  
586 nario where we pretrained the base model on our  
587 textbook-generated dialogue dataset and the dif-  
588 ference between this pretrain version against the

589 version without this pretraining (in bracket).

590 We found that the model that was first pretrained  
591 on the social science textbook data achieved the  
592 highest score when tested on MCTest and CNN  
593 splits of the CoQA dataset, with improvements of  
594 4.16 and 1.99. Meanwhile, the model pretrained  
595 on the business textbook data achieved the highest  
596 score when tested on the NCTE dataset. The model  
597 pretrained on the math textbook data also shows  
598 improvements. As the social textbook dataset con-  
599 tain the least math expressions, it improves most in  
600 non-math domains but does worst in the math do-  
601 main. We conclude that **synthetic datasets created  
602 using our method are usually more effective for  
603 pretraining if they align with the target domain.**

604 Upon a more qualitative human examination of  
605 the generated results, we found that the pretrained  
606 models have a better understanding of the input  
607 context and generate more correct answers than  
608 the corresponding non-pretrained models. Some  
609 example generations are shown in Appendix A.10.

	CoQA (MCTest)	CoQA (CNN)	NCTE
Math	26.10 (+4.03)	13.95 (+0.82)	8.79 (+0.39)
Business	18.91 (-3.22)	13.29 (+0.16)	<b>8.99 (+0.59)</b>
Science	22.36 (+0.22)	14.96 (+1.83)	8.73 (+0.33)
Social	<b>26.30 (+4.16)</b>	<b>15.11 (+1.99)</b>	8.37 (-0.03)
All	23.05 (+0.92)	14.31 (+1.19)	8.41 (+0.01)

Table 5: Downstream Task Results. We use dialogues generated from one textbook from each domain for pre-training and evaluate on downstream benchmarks. Each cell displays BLEU score and the (difference from the baseline), where the baseline is derived from the same model without pre-training.

## 610 **6 Conclusion**

611 We introduced a new task of generating educational  
612 dialogues from textbooks to help pretrain educa-  
613 tional chatbots and detailed various approaches  
614 to simulate student-teacher interactions and create  
615 such data. We evaluated the generated dialogues,  
616 focusing on various measures of goodness, such  
617 as Answer Relevance, Informativeness, Coherence,  
618 and Factual Consistency. Our results indicate that  
619 the approach with LLMs role-playing as teachers  
620 and students for data synthesis excels in most met-  
621 rics. However, upon closer inspection, we also  
622 observed several issues with the synthesized data  
623 such as the problem of hallucinations and repeat-  
624 ing information. Despite these issues, we showed  
625 that the generated dialogues could be used to pre-  
626 train educational chatbots and achieve performance  
627 improvements in various educational settings.

<sup>4</sup><https://pypi.org/project/sacrebleu/>



## 7 Limitations

**Focus on a specific teaching scenario and limitations in educational contexts** In this work, we focus on a specific educational scenario where a curious student asks questions to a knowledgeable teacher. It has been shown that the quality of the student’s questions (with deep reasoning ones) is correlated with their learning (Graesser and Person, 1994; Person et al., 1994). We did not model any of these aspects in our approach. Furthermore, recent approaches of teachers asking Socratic questions or providing indirect scaffolds and hints instead of providing students directly with answers have also been shown to lead to better learning outcomes (Freeman et al., 2014). In our formulation, teachers directly provide students with answers. Future work could explore building a conversational dataset based on more nuances of student questioning patterns (Shridhar et al., 2022) or common scaffolding patterns by teachers.

**Achieving the highest Informativeness is not the overall goal for human learning** : While a dialogue rich in information suggests a potential for a greater extent of learning by a student, there exists a trade-off, as excessive information can increase the student’s cognitive load and become overwhelming (Kaylor, 2014). Therefore, finding the optimal amount of information that the dialogue should contain needs careful consideration in future work.

**Aspects of evaluation framework**: Although we tried to include various aspects of the evaluation in this work, it was not feasible to focus on all important educational aspects. We specifically focused on one setting, where students ask curious questions and the teacher provides answers. Additional aspects, such as the quality of teacher scaffolding, need to be considered to provide a more comprehensive assessment.

## 8 Ethics and Broader Impact Statement

We acknowledge the ethical implications and broader impacts of our work as follows:

### 8.1 Ethical Considerations

**Data Privacy and Anonymity** Our use of open-source textbooks from OpenStax ensures that the data is publicly available and free from privacy concerns. Additionally, in our human evaluation process, we rigorously removed all annotator information to maintain privacy and confidentiality.

**Content Accuracy and Misinformation** We recognize that our best data synthesis method has the problem of hallucinations, which may lead to misinformation. Continuous efforts to improve data accuracy and reduce misinformation are crucial.

### 8.2 Broader Impacts

**Accessibility and Inclusivity** By open-sourcing our data and code, we aim to enable a wider community to benefit from and contribute to this work.

**Potential Misuse** As with any AI-driven dataset, there is a potential for misuse. Our datasets and the accompanying code are intended to serve as supplementary resources in educational settings. It’s important to emphasize that they should not replace human interactions and traditional teaching methods.

### 8.3 Compliance with Ethical Standards

Our research adheres to the ethical code set out in the ACL Code of Ethics. We have taken care to ensure that our methodologies and applications align with these standards, especially regarding data privacy, accuracy, and the responsible use of AI.

## References

- Raviteja Anantha, Svitlana Vakulenko, Zhucheng Tu, Shayne Longpre, Stephen Pulman, and Srinivas Chappidi. 2020. Open-domain question answering goes conversational via question rewriting. *arXiv preprint arXiv:2010.04898*.
- Michael J Baker, Baruch B Schwarz, and Sten R Ludvigsen. 2021. Educational dialogues and computer supported collaborative learning: critical analysis and research perspectives. *International Journal of Computer-Supported Collaborative Learning*, pages 1–22.
- Bill Byrne, Karthik Krishnamoorthi, Saravanan Ganesh, Amit Dubey, Andy Cedilnik, and Kyu-Young Kim. 2020. Taskmaster-2. <https://github.com/google-research-datasets/Taskmaster/tree/master/TM-2-2020>. Second dataset in series of three.
- Andrew Caines, Helen Yannakoudakis, Helena Edmondson, Helen Allen, Pascual Pérez-Paredes, Bill Byrne, and Paula Buttery. 2020. The teacher-student chat-room corpus. *arXiv preprint arXiv:2011.07109*.
- Eunsol Choi, He He, Mohit Iyyer, Mark Yatskar, Wentaoh 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*, pages

726	2174–2184, Brussels, Belgium. Association for Computational Linguistics.	Or Honovich, Leshem Choshen, Roei Aharoni, Ella Neeman, Idan Szpektor, and Omri Abend. 2021. $\mathcal{Q}2$ : Evaluating factual consistency in knowledge-grounded dialogues via question generation and question answering. <i>arXiv preprint arXiv:2104.08202</i> .	779
727			780
728	Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Eric Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. 2022. Scaling instruction-finetuned language models. <i>arXiv preprint arXiv:2210.11416</i> .		781
729			782
730			783
731		Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models. <i>arXiv preprint arXiv:2106.09685</i> .	784
732			785
733	Nico Daheim, Nouha Dziri, Mrinmaya Sachan, Iryna Gurevych, and Edoardo M Ponti. 2023. Elastic weight removal for faithful and abstractive dialogue generation. <i>arXiv preprint arXiv:2303.17574</i> .		786
734			787
735			788
736		Yerin Hwang, Yongil Kim, Hyunkyung Bae, Jeessoo Bang, Hwanhee Lee, and Kyomin Jung. 2023. Dialogizer: Context-aware conversational-qa dataset generation from textual sources. <i>arXiv preprint arXiv:2311.07589</i> .	789
737	Shuyang Dai, Guoyin Wang, Sunghyun Park, and Sungjin Lee. 2021. Dialogue response generation via contrastive latent representation learning. In <i>Proceedings of the 3rd Workshop on Natural Language Processing for Conversational AI</i> , pages 189–197. Online. Association for Computational Linguistics.		790
738			791
739			792
740			793
741		Sara K Kaylor. 2014. Preventing information overload: Cognitive load theory as an instructional framework for teaching pharmacology. <i>Journal of Nursing Education</i> , 53(2):108–111.	794
742			795
743	Zhuyun Dai, Arun Tejasvi Chaganty, Vincent Y Zhao, Aida Amini, Qazi Mamunur Rashid, Mike Green, and Kelvin Guu. 2022. Dialog inpainting: Turning documents into dialogs. In <i>International Conference on Machine Learning</i> , pages 4558–4586. PMLR.		796
744			797
745		Gangwoo Kim, Sungdong Kim, Kang Min Yoo, and Jaewoo Kang. 2022a. Generating information-seeking conversations from unlabeled documents. In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 2362–2378.	798
746			799
747			800
748	Dorottya Demszky and Heather Hill. 2022. The ncte transcripts: A dataset of elementary math classroom transcripts. <i>arXiv preprint arXiv:2211.11772</i> .		801
749			802
750			803
751	Xuan Long Do, Bowei Zou, Shafiq Joty, Tran Tai, Liangming Pan, Nancy Chen, and Ai Ti Aw. 2023. Modeling what-to-ask and how-to-ask for answer-unaware conversational question generation. In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 10785–10803, Toronto, Canada. Association for Computational Linguistics.	Hyunwoo Kim, Jack Hessel, Liwei Jiang, Ximing Lu, Youngjae Yu, Pei Zhou, Ronan Le Bras, Malihe Alikhani, Gunhee Kim, Maarten Sap, et al. 2022b. Soda: Million-scale dialogue distillation with social commonsense contextualization. <i>arXiv preprint arXiv:2212.10465</i> .	804
752			805
753			806
754			807
755			808
756			809
757		Mohammad Amin Kuhail, Nazik Alturki, Salwa Alramlawi, and Kholood Alhejori. 2023. Interacting with educational chatbots: A systematic review. <i>Education and Information Technologies</i> , 28(1):973–1018.	810
758			811
759	Nouha Dziri, Ehsan Kamaloo, Kory W Mathewson, and Osmar Zaiane. 2019. Evaluating coherence in dialogue systems using entailment. <i>arXiv preprint arXiv:1904.03371</i> .		812
760			813
761		Ghader Kurdi, Jared Leo, Bijan Parsia, Uli Sattler, and Salam Al-Emari. 2020. A systematic review of automatic question generation for educational purposes. <i>International Journal of Artificial Intelligence in Education</i> , 30:121–204.	814
762			815
763	Scott Freeman, Sarah L Eddy, Miles McDonough, Michelle K Smith, Nnadozie Okoroafor, Hannah Jordt, and Mary Pat Wenderoth. 2014. Active learning increases student performance in science, engineering, and mathematics. <i>Proceedings of the national academy of sciences</i> , 111(23):8410–8415.		816
764			817
765			818
766		I Loshchilov and F Hutter. 2019. "decoupled weight decay regularization", 7th international conference on learning representations, iclr. <i>New Orleans, LA, USA, May</i> , (6-9):2019.	819
767			820
768			821
769	Jonathan Ginzburg. 2010. Relevance for dialogue. In <i>SemDial: Workshop on the Semantics and Pragmatics of Dialogue (PozDial)</i> , pages 121–129.		822
770		Jakub Macina, Nico Daheim, Sankalan Pal Chowdhury, Tanmay Sinha, Manu Kapur, Iryna Gurevych, and Mrinmaya Sachan. 2023. Mathdial: A dialogue tutoring dataset with rich pedagogical properties grounded in math reasoning problems. <i>arXiv preprint arXiv:2305.14536</i> .	823
771			824
772	Arthur C Graesser and Natalie K Person. 1994. Question asking during tutoring. <i>American educational research journal</i> , 31(1):104–137.		825
773			826
774			827
775	Max Grusky, Mor Naaman, and Yoav Artzi. 2018. Newsroom: A dataset of 1.3 million summaries with diverse extractive strategies. <i>arXiv preprint arXiv:1804.11283</i> .		828
776			829
777		Julia M Markel, Steven G Opferman, James A Landay, and Chris Piech. 2023. Gpteach: Interactive training with gpt based students.	830
778			831

832	Alireza Mohammadshahi, Thomas Scialom, Majid Yazdani, Pouya Yanki, Angela Fan, James Henderson, and Marzieh Saeidi. 2022. Rquge: Reference-free metric for evaluating question generation by answering the question. <i>arXiv preprint arXiv:2211.01482</i> .	Katherine Stasaski, Kimberly Kao, and Marti A Hearst. 2020. Cima: A large open access dialogue dataset for tutoring. In <i>Proceedings of the Fifteenth Workshop on Innovative Use of NLP for Building Educational Applications</i> , pages 52–64.	887
833			888
834			889
835			890
836			891
837	Natalie K Person, Arthur C Graesser, Joseph P Magliano, and Roger J Kreuz. 1994. Inferring what the student knows in one-to-one tutoring: The role of student questions and answers. <i>Learning and individual differences</i> , 6(2):205–229.	Katherine Stasaski, Manav Rathod, Tony Tu, Yunfang Xiao, and Marti A Hearst. 2021. Automatically generating cause-and-effect questions from passages. In <i>Proceedings of the 16th Workshop on Innovative Use of NLP for Building Educational Applications</i> , pages 158–170.	892
838			893
839			894
840			895
841			896
842	Matt Post. 2018. A call for clarity in reporting bleu scores. <i>arXiv preprint arXiv:1804.08771</i> .	Abhijit Suresh, Jennifer Jacobs, Margaret Perkoff, James H. Martin, and Tamara Sumner. 2022. <b>Fine-tuning transformers with additional context to classify discursive moves in mathematics classrooms</b> . In <i>Proceedings of the 17th Workshop on Innovative Use of NLP for Building Educational Applications (BEA 2022)</i> , pages 71–81, Seattle, Washington. Association for Computational Linguistics.	898
843			899
844	Peng Qi, Yuhao Zhang, and Christopher D Manning. 2020. Stay hungry, stay focused: Generating informative and specific questions in information-seeking conversations. <i>arXiv preprint arXiv:2004.14530</i> .		900
845			901
846			902
847			903
848	Chen Qu, Liu Yang, Cen Chen, Minghui Qiu, W Bruce Croft, and Mohit Iyyer. 2020. Open-retrieval conversational question answering. In <i>Proceedings of the 43rd International ACM SIGIR conference on research and development in Information Retrieval</i> , pages 539–548.		904
849			905
850			906
851			907
852			908
853			909
854	Manav Rathod, Tony Tu, and Katherine Stasaski. 2022. Educational multi-question generation for reading comprehension. In <i>Proceedings of the 17th Workshop on Innovative Use of NLP for Building Educational Applications (BEA 2022)</i> , pages 216–223.	Wei Tan, Jionghao Lin, David Lang, Guanliang Chen, Dragan Gašević, Lan Du, and Wray Buntine. 2023. Does informativeness matter? active learning for educational dialogue act classification. In <i>International Conference on Artificial Intelligence in Education</i> , pages 176–188. Springer.	910
855			911
856			912
857			913
858			914
859	Vipula Rawte, Amit Sheth, and Amitava Das. 2023. A survey of hallucination in large foundation models. <i>arXiv preprint arXiv:2309.05922</i> .	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, et al. 2020. Transformers: State-of-the-art natural language processing. In <i>Proceedings of the 2020 conference on empirical methods in natural language processing: system demonstrations</i> , pages 38–45.	915
860			916
861			917
862	Siva Reddy, Danqi Chen, and Christopher D Manning. 2019. Coqa: A conversational question answering challenge. <i>Transactions of the Association for Computational Linguistics</i> , 7:249–266.	Lichao Zhang, Abel Gonzalez-Garcia, Joost Van De Weijer, Martin Danelljan, and Fahad Shahbaz Khan. 2018. Synthetic data generation for end-to-end thermal infrared tracking. <i>IEEE Transactions on Image Processing</i> , 28(4):1837–1850.	918
863			919
864			920
865			921
866	Nils Reimers and Iryna Gurevych. 2019. <b>Sentence-BERT: Sentence embeddings using Siamese BERT-networks</b> . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 3982–3992, Hong Kong, China. Association for Computational Linguistics.	Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. 2019. Bertscore: Evaluating text generation with bert. <i>arXiv preprint arXiv:1904.09675</i> .	922
867			923
868			924
869			925
870			926
871			927
872			928
873			929
874	Thomas Scialom, Paul-Alexis Dray, Patrick Gallinari, Sylvain Lamprier, Benjamin Piwowarski, Jacopo Staiano, and Alex Wang. 2021. Questeval: Summarization asks for fact-based evaluation. <i>arXiv preprint arXiv:2103.12693</i> .	<b>A Appendix</b>	930
875			931
876			932
877			933
878			934
879	Kumar Shridhar, Jakub Macina, Mennatallah El-Assady, Tanmay Sinha, Manu Kapur, and Mrinmaya Sachan. 2022. <b>Automatic generation of socratic subquestions for teaching math word problems</b> . In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 4136–4149, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	<b>A.1 Metrics Formulas</b>	935
880			936
881			937
882			938
883			939
884			940
885			941
886			942

Criterion	Metric	Definition	Explanation
Answer Relevance	BF1 ( $q_t, a_t$ )	BERTScoreF1( $q_t, a_t$ )	For each QA pair, we compute the BERTScore F1 (Zhang et al., 2019), treating the question as the predicted sentence and the answer as the reference sentence. It evaluates the semantic correspondence between the question and answer using BERT’s contextual embeddings.
	QuestEval	QuestEval( $q_t, a_t$ )	For each QA pair, we compute the QuestEval score (Scialom et al., 2021), treating the question as the predicted sentence and the answer as the reference sentence. QuestEval generates questions from both the original question and the answer, then generates answers for these questions, comparing their consistency and completeness to evaluate Answer Relevance.
Coherence	BF1 ( $q_t, a_{<t}$ )	BERTScoreF1( $q_t, a_{<t}$ )	It computes the BERTScore F1 for each dialogue question, treating it as the predicted sentence against all preceding answers as references. Aggregated scores reflect the dialogue’s coherence.
	BF1 ( $q_t, a_{(t-1)}$ )	BERTScoreF1( $q_t, a_{(t-1)}$ )	It computes the BERTScore F1 for each dialogue question against the immediately preceding answer as the reference. Aggregated scores provide a measure of overall coherence.
Informativeness	1-Overlap ( $a_t, a_{<t}$ )	$1 - \frac{ a_t \cap a_{<t} }{ a_t \cup a_{<t} }$	For each answer in a dialogue, the proportion of its intersection with previous answers to their union is computed using word-level tokens. This value is then subtracted from 1.
Content Match	Density	$\frac{1}{ h_{1:T} } \sum_{f \in \mathcal{F}(S, h_{1:T})}  f ^2$ ., $\mathcal{F}(S, h_{1:T})$ : the set of extractive phrases in dialogue $h_{1:T}$ extracted from textbook content $S$ .	Density refer to Extractive Fragment Density (Grusky et al., 2018), as the average length of text spans that are directly extracted from textbook content $S$ and included in the dialogues.
	Coverage	$\frac{1}{ h_{1:T} } \sum_{f \in \mathcal{F}(S, h_{1:T})}  f $	Coverage refer to Extractive Fragment Coverage (Grusky et al., 2018), as the percentage of words in a dialogue that originated from the textbook content.
Answerability	Answerable	Valid(QA( $q_t, S$ ))	We use the “distilbert-base-cased-distilled-squad” QA model to determine if a question is answerable from the textbook content. If it generates an empty string or an invalid answer such as “CANNOTANSWER”, the question is deemed unanswerable. We report the ratio of answerable questions as 1 minus the ratio of unanswerable questions.
Factual Consistency	QFactScore	$\alpha sim(QA(q_t, S), a_t) + \beta sim(q_t, a_t)$	For each QA pair, it computes the cosine similarity between the embeddings of the QA model’s predicted answer and the original answer. Then, it assesses the similarity between the embeddings of the question and answer. The final score is the weighted sum of two similarity scores.
Specificity	NA	NA	We lack automatic metrics for evaluating this criterion.

Table 6: Criteria with Formulas and Explanations

Domain	Name	Chapters	Paragraphs	Pages	Words
Math	Introductory Statistics	13	1,412	65	35,182
Business	Business Ethics	11	795	42	85,626
Science	Physics	23	1,918	89	106,712
Social science	Psychology 2e	16	1,710	88	191,273

Table 7: Summary of the textbook statistics.

## A.3 From Textbooks to Dialogues Details

### A.3.1 Information-seeking scenario

In the SimSeek-ASYM setup, the CQG model ingests the title and summary information, each separated by special tokens. We use T5-Large as the student’s model and Longformer-Large as the teacher’s model.

The SimSeek-ASYM code<sup>5</sup> can be executed with minor modifications. We use the same CQG and CAF models as in (Kim et al., 2022a), which utilize T5 as the student’s model and Longformer as the teacher’s model.

### A.3.2 Dialogue Inpainting

We adopt a training regimen that integrates data from the OR-QuAC (Qu et al., 2020), QReCC (Anantha et al., 2020), and the movie and restaurant datasets from Taskmaster-2 (Byrne et al., 2020), employing the technique as described in (Dai et al., 2022). We randomly selected 80% of the data as the training set, while the remaining 20% as the test set. We implement Dialogue Inpainting using the

<sup>5</sup><https://github.com/naver-ai/simseek>

962	code framework of (Daheim et al., 2023), basing	ask questions about information you already	1009
963	our model (eq 2) on FLAN-T5-XL (Chung et al.,	ready have. Only ask one question at a	1010
964	2022), and train it with LoRA (Hu et al., 2021)	time.	1011
965	to reduce computational load. We used one V100	Expected Output: Please phrase your	1012
966	GPU to train the model, the FLAN-T5-XL model	question as a string.	1013
967	has 3 Billion parameters and took 12 hours to train.		
968	The model, while fundamentally designed to	<b>Prompt for simulating teacher</b>	1014
969	predict single utterances, is used autoregressively.	Task: You are a teacher preparing to an-	1015
970	It begins with the input $s_{prompt}, < mask >, s_1$	swer a student’s question about a subsec-	1016
971	and sequentially generates questions using top-p	tion of a textbook. The student’s ques-	1017
972	sampling. This autoregressive process continues	tion is: {question}. Provide a concise,	1018
973	until the dialogue is wholly formed.	specific response, ensuring it’s not a sum-	1019
		mary and distinct from any previous an-	1020
974	<b>A.3.3 Persona-based Generation</b>	swers you’ve given.	1021
975	<b>Prompt for Persona (High Info)</b> The design of	Information Provided:	1022
976	our prompts was chiefly driven by the requisites	1. Section Title: ...	1023
977	of context-awareness, speaker identification, and	2. Subsection Title: ...	1024
978	specificity. We incorporated guidelines and an-	3. Subsection Content: ...	1025
979	notations to ensure GPT yields concise responses	4. Section Summary: ...	1026
980	and minimizes redundant information. To distin-	5. Bold Terms in Section: ...	1027
981	guish between speakers, we prefixed dialogues with	6. Learning Objectives: ...	1028
982	labels:“Teacher:” or “Student:”. The prompt is	7. Concepts in Section: ...	1029
983	shown below.	8. Section Introduction: ...	1030
984	<b>Prompt for simulating student</b>	Previous Conversation:	1031
985	Task: You are a student preparing to ask	Student:...	1032
986	questions about a textbook subsection	Teacher:...	1033
987	to a teacher. Your goal is to uncover	*Note:* When crafting your response,	1034
988	the key information from this subsection.	consider all the information above. Be	1035
989	Based on the teacher’s responses, you’ll	sure your answer directly addresses the	1036
990	further inquire to get a comprehensive	student’s question and is not a repetition	1037
991	understanding. Make sure to ask specific	of prior information.	1038
992	questions about the subsection’s content	Expected Output: Please phrase your an-	1039
993	and avoid repeating queries from prior	swer as a string.	1040
994	discussions.		
995	Information Provided:	<b>Prompt for Persona (Single Instance)</b> The	1041
996	1. Section Title: ...	prompt for the Persona (Single Instance) method is	1042
997	2. Subsection Title: ...	shown below. It uses one prompt to generate one	1043
998	3. Section Summary: ...	dialogue.	1044
999	4. Bold Terms in Section: ...	Task: generate a conversation between	1045
1000	5. Learning Objectives: ...	a student and a teacher using the given	1046
1001	6. Concepts in Section: ...	section.	1047
1002	7. Section Introduction: ...	Introduction:	1048
1003	Previous Conversation:	1. The conversation should contain 6	1049
1004	Student:...	question-answer pairs.	1050
1005	Teacher:...	2. The output conversation should be	1051
1006	*Note:* Frame your questions consid-	in this format: student: ... teacher:	1052
1007	ering the information above and ensure	... student: ...	1053
1008	they’re relevant to the content. Do not	3. The given section: ...	1054

#### A.4 QFactScore Implementation

For computing the embeddings of questions and answers, we use the “msmarco-distilbert-cos-v5” model from (Reimers and Gurevych, 2019). This model is suitable for computing cosine similarity and performs well in our task.

It is important to ensure that the QA model used in QFactScore is different from the QA model used for generating dialogue datasets. This is because if the same QA model is used, the predicted answer is likely to be similar to the original answer in the dialogue. In QFactScore, we use the ‘distilbert-base-cased-distilled-squad’ model, which differs from the GPT-3.5, T5, and Flan-T5 models that we used for generating the dataset.

QFactScore computes as the below equation. For each QA pair, it computes the cosine similarity between the embeddings of the QA model’s predicted answer and the original answer. Then, it assesses the similarity between the embeddings of the question and answer. The final score is the weighted sum of two similarity scores. The weight can be adjusted according to different applications, in our study we use  $\alpha = 1$  and  $\beta = 1$ .

$$\alpha \cdot \text{sim}(\text{QA}(q_t, S), a_t) + \beta \cdot \text{sim}(q_t, a_t) \quad (3)$$

We further evaluate the correlation between QFactScore and human evaluation of Factual Consistency in Appendix A.6. We also provide correlation between 1 - Overlap( $a_t, a_{<t}$ ) and human evaluation of Informativeness in Appendix A.6.

#### A.5 Model Comparison

The details of the different models are listed in Table 8. The term “Formatting” refers to formatting information, which contains a title, summary, introduction, learning objectives, bold terms, and key concepts from textbooks, which is introduced in Section-3.1. The “COPY” in the teacher’s model of Dialogue Inpainting indicates this method just copying a sentence from the textbook as the answer.

#### A.6 Metric Evaluation

To validate the effectiveness of the metrics introduced in this study, we calculated both Pearson and Spearman correlations between the metrics’ outcomes and the corresponding results from human evaluations. The results are shown in Table 9 and Table 10. The “1 - Overlap( $a_t, a_{<t}$ )” exhibits

a Pearson correlation of 0.84 and a Spearman correlation of 0.82 with the Informativeness score in human evaluation, both with p-values below 0.005, suggesting that this F1 score could effectively represent Informativeness in evaluations.

On the other hand, QFactScore exhibits a Pearson correlation of 0.36 and a Spearman correlation of 0.36 with Factual Consistency in human evaluation, both with p-values below 0.0005. We interpret this as indicative of a moderate correlation, suggesting that this metric can approximate factual consistency to a certain extent. When comparing the correlation results with existing methods, including the use of GPT-3.5 scores derived from prompts, QuestEval, and QrelScore, the findings indicate that QFactScore’s correlation score surpasses others. However, Factual Consistency is a nuanced criterion that necessitates an assessment of whether the answer accurately addresses the question within the given context. Existing metrics struggle with this task, highlighting the need for more comprehensive evaluations in the future.

#### A.7 Metrics Results Details

We provide the complete results of different metrics for datasets in four domains in this section. The results are shown in Table 11 and Table 12.

#### A.8 Human Evaluation Details

##### A.8.1 Experiment Details

We have adopted a human evaluation approach to assess the performance of dialogues generated by various methods. We recruited three expert annotators to undertake the task; all annotators have at least a master’s degree in computer science. The annotators have educational backgrounds in Europe and age between 20-25. We recruit them by advertising on social media and bonus with some gifts for each annotator. As all annotators are satisfied with this payment, we consider this as adequate. To alleviate the burden on participants, we selected the 3 models from each method category for evaluation. To ensure the consistency of results across different domains, we chose datasets from four textbooks, each covering a different subject area: mathematics, business, science, and social sciences. From each textbook, we randomly selected a subsection. For each subsection, we generated one dialogue using a different method, preparing each dialogue separately for evaluation. We use only the first 12 turns (6 QA pairs) of each dialogue for evaluation, similar to what is described in Section 5.1.2. During the

Models	Student’s Model	Teacher’s Model	Input to Student	Input to Teacher
SimSeek Dialog Inpainting	T5 FLAN-T5	Longformer COPY	Title + Summary Contents + Formatting	Contents+Formatting
Persona (Low Info) Persona (Medium Info) Persona (High Info) GPT (Single Instance)	GPT-3.5	GPT-3.5	Title Title + Summary Formatting Contents + Formatting	

Table 8: Model Comparison

		Correlation	P-Value
1 - Overlap( $a_t, a_{<t}$ )	vs Informativeness	<b>0.84</b>	<b>0.0006</b>
1 - BF1( $a_t, a_{<t}$ )	vs Informativeness	0.75	0.006
QFactScore	vs Factual Consistency	<b>0.36</b>	<b>0.0002</b>
GPT-3.5	vs Factual Consistency	0.30	0.002
QuestEval	vs Factual Consistency	0.21	0.04
QrelScore	vs Factual Consistency	-0.02	0.9

Table 9: Pearson correlation of metrics and human evaluation

		Correlation	P-Value
1 - Overlap( $a_t, a_{<t}$ )	vs Informativeness	<b>0.82</b>	<b>0.001</b>
1 - BF1( $a_t, a_{<t}$ )	vs Informativeness	0.77	0.003
QFactScore	vs Factual Consistency	<b>0.36</b>	<b>0.0002</b>
GPT-3.5	vs Factual Consistency	0.27	0.007
QuestEval	vs Factual Consistency	0.21	0.03
QrelScore	vs Factual Consistency	0.06	0.6

Table 10: Spearmans correlation of metrics and human evaluation

evaluation, each of the three participants received 12 dialogues, with every dialogue corresponding to a related textbook subsection. Evaluators rated each question-answer (QA) pair within a dialogue based on eight criteria. The overall evaluation score for a dialogue was determined by averaging the scores of all its QA pairs. The specific evaluation criterion and corresponding questions are detailed in Table 13. Participants responded to each question with “yes” or “no”. The “yes” is recorded as a score of 1, while the “no” is recorded as a score of 0.

We provide the specific question the participants will be asked during human evaluation as shown in Table 13. The task is straight forward, that we provide QA pairs for evaluation in an excel file and

the annotators just read the QA pair and give score based on their judgement of each question.

We further show the Cohen’s Kappa score between each participant in Table 14, which proves that each pair of participants has substantial agreement.

### A.8.2 Disclaimer for Annotators

Thank you for participating in our evaluation process. Please read the following important points before you begin:

- **Voluntary Participation:** Your participation is completely voluntary. You have the freedom to withdraw from the task at any time without any consequences.
- **Confidentiality:** All data you will be working with is anonymized and does not contain any personal information. Your responses and scores will also be kept confidential.
- **Risk Disclaimer:** This task does not involve any significant risks. It primarily consists of reading and scoring QA pairs.
- **Queries:** If you have any questions or concerns during the task, please feel free to reach out to us.

### A.8.3 Instructions for Experiments

Thank you for participating in our evaluation experiment. The data collected through this process will be used to assess the quality of our methods.

Follow these steps to score each QA pair:

1. **Accessing the Data:** Open the provided Excel file, which contains the QA pairs for evaluation.
2. **Scoring Each QA Pair:** For each pair, read the question and the corresponding answer carefully.

Domain	Models	Question Type			Number of Tokens	
		% “what” or “which”	% “why”	% “how”	Avg Tokens in Questions	Avg Tokens in Answers
Math	SimSeek	49	2	23	11.24	11.66
	Dialogue Inpainting	48	3	20	7.55	15.41
	Persona (Low Info)	47	0.1	0.5	18.19	80.16
	Persona (Medium Info)	55	0	0.5	18.28	81.83
	Persona (High Info)	69	0.1	0.2	19.95	77.96
	GPT (Single Instance)	9	0.8	7	15.57	29.96
Business	SimSeek	49	2	14	11.19	16.17
	Dialogue Inpainting	62	3	17	6.75	23.79
	Persona (Low Info)	65	0	1	17.74	99.03
	Persona (Medium Info)	47	0	0.4	18.94	99.36
	Persona (High Info)	60	0	0	19.52	98.86
	GPT (Single Instance)	12	3	10	16.28	40.84
Science	SimSeek	54	1	16	10.73	14.73
	Dialogue Inpainting	53	3	17	6.55	17.91
	Persona (Low Info)	57	0	1	17.50	83.12
	Persona (Medium Info)	55	0	0.4	16.71	83.70
	Persona (High Info)	56	0.2	0.2	18.40	84.05
	GPT (Single Instance)	14	2	12	13.43	31.31
Social Science	SimSeek	53	1	13	10.42	14.74
	Dialogue Inpainting	60	3	15	6.53	21.42
	Persona (Low Info)	53	0	0.6	16.82	76.68
	Persona (Medium Info)	51	0.04	0.3	16.82	75.87
	Persona (High Info)	52	0.04	1	18.15	77.91
	GPT (Single Instance)	12	1	8	14.58	36.20

Table 11: Dataset statistics in more detail

Domain	Models	Answer Relevance		Informativeness	Groundedness		Coherence		Answerability	Factual Consistency
		BFI ( $q_t, a_t$ )	QuestEval	1 - Overlap ( $a_t, a_{<t}$ )	Density	Coverage	BFI ( $q_t, a_{<t}$ )	BFI ( $q_t, a_{t-1}$ )	Answerable	QFactScore
Math	SimSeek	0.51	0.24	0.61	9.5	0.71	0.49	0.53	0.74	0.27
	Dialogue Inpainting	0.57	0.30	<b>0.88</b>	<b>19.37</b>	<b>0.88</b>	0.46	0.47	0.52	0.19
	Persona (Single Instance)	0.58	0.32	0.85	2.94	0.62	0.50	0.52	0.87	0.53
	Persona (Low Info)	<b>0.62</b>	<b>0.43</b>	0.54	1.94	0.59	0.51	0.59	0.99	0.80
	Persona (Medium Info)	<b>0.62</b>	<b>0.43</b>	0.55	2.09	0.60	0.51	0.59	<b>1.00</b>	<b>0.81</b>
	Persona (High Info)	<b>0.62</b>	<b>0.43</b>	0.56	2.07	0.60	<b>0.52</b>	<b>0.60</b>	0.99	<b>0.81</b>
Business	SimSeek	0.54	0.25	0.77	13.16	0.88	0.52	0.56	0.89	0.32
	Dialogue Inpainting	0.49	0.26	<b>0.94</b>	<b>26.44</b>	<b>0.92</b>	0.43	0.45	0.88	0.23
	Persona (Single Instance)	0.58	0.36	0.88	4.07	0.82	0.50	0.53	0.95	0.52
	Persona (Low Info)	0.62	<b>0.46</b>	0.61	2.38	0.76	0.52	<b>0.60</b>	0.99	0.73
	Persona (Medium Info)	0.62	<b>0.46</b>	0.61	2.31	0.77	0.53	<b>0.60</b>	<b>1.00</b>	0.73
	Persona (High Info)	<b>0.63</b>	<b>0.46</b>	0.62	2.44	0.77	<b>0.54</b>	<b>0.60</b>	<b>1.00</b>	<b>0.74</b>
Science	SimSeek	0.52	0.25	0.71	11.78	0.83	0.50	0.54	0.89	0.34
	Dialogue Inpainting	0.51	0.27	<b>0.92</b>	<b>20.43</b>	<b>0.90</b>	0.44	0.44	0.72	0.24
	Persona (Single Instance)	0.58	0.35	0.85	4.65	0.79	0.48	0.51	0.94	0.61
	Persona (Low Info)	<b>0.59</b>	<b>0.43</b>	0.57	2.55	0.73	<b>0.51</b>	<b>0.57</b>	0.98	0.79
	Persona (Medium Info)	<b>0.59</b>	<b>0.43</b>	0.57	2.63	0.73	0.50	<b>0.57</b>	<b>0.99</b>	<b>0.80</b>
	Persona (High Info)	<b>0.59</b>	<b>0.43</b>	0.58	2.68	0.74	<b>0.51</b>	<b>0.57</b>	<b>0.99</b>	0.76
Social Science	SimSeek	0.53	0.27	0.74	12.21	0.84	0.51	0.55	0.89	0.34
	Dialogue Inpainting	0.51	0.28	<b>0.91</b>	<b>24.22</b>	<b>0.91</b>	0.45	0.48	0.86	0.29
	Persona (Single Instance)	0.57	0.36	0.87	4.09	0.77	0.49	0.52	0.92	0.50
	Persona (Low Info)	<b>0.62</b>	<b>0.45</b>	0.63	2.67	0.73	0.52	0.59	0.98	0.69
	Persona (Medium Info)	<b>0.62</b>	<b>0.45</b>	0.64	2.69	0.73	0.52	<b>0.60</b>	<b>1.00</b>	<b>0.71</b>
	Persona (High Info)	<b>0.62</b>	<b>0.45</b>	0.63	2.79	0.74	<b>0.53</b>	0.59	0.99	0.69

Table 12: Metrics results of different datasets



Criterion	Questions for each QA pair
Answer Relevance	<b>Question:</b> Is the response directly addressing the posed question? (answer no if it is answering a different question)
Informativeness	<b>Question:</b> Does the current answer introduce new information that was not mentioned in previous answers within the same conversation?
Groundedness	<b>Question:</b> Does the answer contain specific details or data points mentioned in the contextual background or previous dialogue?
Coherence	<b>Question 1:</b> Does the current question directly follow up on the immediate previous answer? (Ignore the first QA pair)
Factual Consistency	<b>Question:</b> Does the answer correctly address the question, considering the context provided? (If 'answerability' is 'no,' then this criterion should also be 'no.')
Answerability	<b>Question:</b> Can the question be answered given the context?
Specificity	<b>Question:</b> Does this question exhibit generality, such that it could be relevant beyond the immediate context provided? (e.g. What is interesting about this passage?)

Table 13: Exact framing of questions asked during the human evaluation.

Participants Pairs	Cohen's Kappa
P1 vs. P2	0.67
P1 vs. P3	0.60
P2 vs. P3	0.71

Table 14: The Cohen's Kappa score between each pair of participants.

3. **Scoring Scale:** Answer each question with "yes" or "no".
4. **Entering Scores:** Enter your score for each QA pair in the designated column in the Excel sheet. Please stick to the scoring scale provided.
5. **Consistency:** Try to maintain consistency in your scoring. Refer to the example evaluations provided if you're unsure.
6. **Completion:** Once you have scored all the QA pairs, save the file and return it to us as instructed.

We appreciate your time and effort in this task.

#### A.8.4 Ethics Review

In our study, the data collection protocol was strictly devised in accordance with the ethical guidelines of our university. According to these

regulations, it did not need to be reviewed by the university's ethics review board, as this experiment does not involve any medical devices, human body effects, or diseases.

#### A.9 Pretraining for Educational Chatbots Details

We sourced four textbooks from the OpenStax website for our study. These include 'Introductory Statistics' for math, 'Business Ethics' for business studies, 'Physics' for science, and 'Psychology 2e' for social science. We use the entire textbook dialogue dataset for pretraining.

In line with the methodology described in (Macina et al., 2023), the models with pretrain were trained 10 epochs during pretrain and trained 10 epochs during finetune. The models without pretrain trained 10 epochs during training. For CoQA CNN and MCTest dialogue datasets for finetune or training, we use 60% of data for training, 20% for validation, and 20% for testing. For the NCTE dataset, we randomly select 10,000 dialogues for training, 2,000 dialogues for validation, and 2,000 dialogues for testing. We set an initial learning rate of  $6.25e-5$  and employed linear learning rate decay without warmup. For model optimization, we utilized checkpoints from the transformers library (Wolf et al., 2020). The negative log-likelihood of the ground-truth response was minimized using the AdamW optimizer, as detailed in (Loshchilov and Hutter, 2019). Model performance was assessed using the sacrebleu implementation of the BLEU metric, following (Post, 2018). We used one V100 GPU to train the model, the FLAN-T5-LARGE model has 0.8 Billion parameters and took 7 hours to train. The result in Table 5 is the average of 3 runs.

#### A.10 Pretraining for Educational Chatbots Generation Results Examples

The ground truth example of CoQA CNN split is shown in Table 15, the generation results of the same dialogue using model without pretraining is shown in Table 16, the generation results of the same dialogue using model with pretrain on social textbook dialogue datasets is shown in Table 17.

The ground truth example of CoQA MCTest split is shown in Table 18, and the generation results of the same dialogue using model without pretraining is shown in Table 19, the generation results of the same dialogue using model with pretrain on social textbook dialogue datasets is shown in Table 20.

Q: What is the name of the Nigerian performer?  
A: Thomas Chinedu Ehiem.  
Q: What did he arrange?  
A: gay liaisons.  
Q: For who?  
A: an Italian government official.  
Q: Who caught him doing so?  
A: Italian authorities  
Q: How?  
A: wiretaps.  
Q: When did this happen?  
A: 2008, to 2010.  
Q: What other role did Balducci serve?  
A: a papal usher.

Table 15: CoQA CNN split dialogue example (ground truth).

Q: What is the name of the Nigerian performer?  
A: **Thomas Ayala.**  
Q: What did he arrange?  
A: gay liaisons.  
Q: For who?  
A: an Italian government official.  
Q: Who caught him doing so?  
A: Italian authorities  
Q: How?  
A: wiretaps.  
Q: When did this happen?  
A: **as part of an investigation into how public-works contracts were awarded.**  
Q: What other role did Balducci serve?  
A: **oversaw the Italian government’s awarding of construction contracts.**

Table 16: CoQA CNN split dialogue generation example (without pretraining).

Q: What is the name of the Nigerian performer?  
A: Thomas.  
Q: What did he arrange?  
A: gay liaisons.  
Q: For who?  
A: an Italian government official.  
Q: Who caught him doing so?  
A: Italian authorities  
Q: How?  
A: wiretaps.  
Q: When did this happen?  
A: **2013.**  
Q: What other role did Balducci serve?  
A: a papal usher.

Table 17: CoQA CNN split dialogue generation example (with pretraining on social textbook dataset).

Q: What else did he like?  
A: Looking at the pictures.  
Q: Who called him?  
A: His mother.  
Q: What did he do first?  
A: Closed the book.  
Q: What was watching the cat watching?  
A: a bird  
Q: Where was it?  
A: in a tree  
Q: what color was it?  
A: blue.

Table 18: CoQA MCTest split dialogue example (ground truth).

Q: What else did he like?  
A: Looking at the pictures.  
Q: Who called him?  
A: His mother.  
Q: What did he do first?  
A: **Joe closed the book, and went downstairs.**  
Q: What was watching the cat watching?  
A: a bird  
Q: Where was it?  
A: in a tree  
Q: what color was it?  
A: blue.

Table 19: CoQA MCTest split dialogue generation example (without pretraining).

Q: What else did he like?  
A: Looking at the pictures.  
Q: Who called him?  
A: His mother.  
Q: What did he do first?  
A: Closed the book.  
Q: What was watching the cat watching?  
A: a bird  
Q: Where was it?  
A: in a tree  
Q: what color was it?  
A: blue.


Table 20: CoQA MCTest split dialogue generation example (with pretraining on social textbook dataset).

terms:	1276
<b>Prohibited Use</b> The code and datasets shall not be used for commercial purposes without prior written consent from the authors.	1277 1278 1279
<b>Attribution</b> When using or referencing the code and datasets, users must provide proper attribution to the original authors.	1280 1281 1282
<b>No Warranty</b> This project is provided "as is" without any warranties of any kind, either expressed or implied, including but not limited to fitness for a particular purpose. The authors are not responsible for any damage or loss resulting from the use of this project.	1283 1284 1285 1286 1287 1288

## A.11 Datasets Overview


We provide the overview of our generated dataset in Table 21.


## A.12 Terms of Use

This section outlines the terms and conditions for the use of  Book2Dial. By using the code and datasets in this project, users agree to the following

Domain	Generation Method	Dialogues	Dialogic Pairs	Bigram Entropy	Avg. words per utterance
Math	Persona (High Info)	142	852	6.08	48.95
	Dialog Inpainting	142	1444	4.07	11.05
Business	Persona (High Info)	123	738	6.61	59.01
	Dialog Inpainting	123	3575	4.46	14.39
Science	Persona (High Info)	228	1368	6.22	48.03
	Dialog Inpainting	228	5898	4.56	13.99
Social	Persona (High Info)	396	2376	6.2	51.04
	Dialog Inpainting	396	7503	4.34	11.69
Total		1778	23754	5.3175	19.48875

Table 21: Detailed Overview of the Synthetic dataset

**Liability** The authors shall not be held liable for any direct, indirect, incidental, special, exemplary, or consequential damages arising in any way out of the use of the  Book2Dial project.

**Updates and Changes** The authors reserve the right to make changes to the terms of this license or the  Book2Dial itself at any time.

### A.13 Compliance with Artifact Usage and Intended Use Specifications

#### A.13.1 Compliance with Existing Artifact Usage

In our study, we utilized a range of existing artifacts, such as open-source textbooks from OpenStax, to develop our research datasets. We rigorously ensured that our usage of these materials was in strict accordance with their intended purposes, aligning with OpenStax’s vision of freely accessible educational content. Additionally, we employed various computational tools within their prescribed licensing terms, thus adhering to ethical and legal standards.

#### A.13.2 Specification of Intended Use for Created Artifacts

Our research led to the development of two significant artifacts:

**Framework for Generating Dialogues from Textbooks** **Intended Use:** This framework is designed for academic research and educational technology development. It facilitates the generation of synthetic dialogues, aiming to enhance AI-driven educational tools. **Restrictions:** The framework should be used within the bounds of educational and research settings. Any commercial or high-stakes educational application is advised against without further validation and ethical review. **Ethical Considerations:** We emphasize the responsible use of this framework, particularly in maintaining the integrity and context of the source textbooks.

### Dataset of Generated Dialogues **Intended Use:**

The dataset is primarily intended for research in educational chatbots and conversational AI. It offers a resource for developing and testing dialogue systems in educational contexts. **Restrictions:** This dataset is not recommended for direct application in live educational settings without substantial vetting, as it may contain synthetic inaccuracies. **Data Ethics:** As the dataset is derived from open-source textbooks, it respects the principles of open access. We encourage users to keep the dataset within academic and research domains, in line with the ethos of the source material.

### A.14 Data Collection and Anonymization Procedures

In our research, rigorous steps were taken to ensure that the data collected and used did not contain any personally identifiable information or offensive content. The data, primarily sourced from open-access textbooks, inherently lacked individual personal data. For the components involving human interaction, such as feedback or evaluation, all identifying information was carefully removed to maintain anonymity. Additionally, we implemented a thorough review process to screen for and exclude any potentially offensive or sensitive material from our dataset. These measures were taken to uphold the highest standards of privacy, ethical data usage, and respect for individual confidentiality.

### A.15 Artifact Documentation

#### A.15.1 Dialogue Generation Framework

**Domain Coverage** The framework is designed to generate dialogues across a range of academic subjects, as exemplified by the textbooks used (math, business, science, social science).

**Linguistic Phenomena** It captures various linguistic phenomena, including question-answering patterns and dialogue quality regarding different

criteria.

### A.15.2 Dataset of Generated Dialogues

**Language and Style** The dialogues are primarily in English, reflecting the language of the source textbooks. The style is educational and academic, suited for educational purposes.

**Content Diversity** The dataset spans multiple academic disciplines, offering a rich variety of topics and themes.

**Demographic Representation** While the dataset itself does not directly represent demographic groups (as it is synthesized from textbooks), the diversity in the source material reflects a broad spectrum of cultural and societal contexts.

### A.16 Use of AI Assistants in Research

In our study, AI assistants were used sparingly and in accordance with ACL’s ethical guidelines. GPT-3.5 was employed for data generation tasks, integral to our research objectives. Additionally, we utilized ChatGPT and Grammarly for basic paraphrasing and grammar checks, respectively. These tools were applied minimally to ensure the authenticity of our work and to adhere strictly to the regulatory standards set by ACL. Our use of these AI tools was focused, responsible, and aimed at supplementing rather than replacing human input and expertise in our research process.

### A.17 Experimental Details

We implement Dialogue Inpainting using the code framework of Daheim et al. (2023), basing our model (eq 2) on FLAN-T5-XL (Chung et al., 2022), and train it with LoRA (Hu et al., 2021) to reduce computational load. We set an initial learning rate of  $6.25e-5$  and employed linear learning rate decay without warmup. For model optimization, we utilized checkpoints from the transformers library (Wolf et al., 2020). The negative log-likelihood of the ground-truth response was minimized using the AdamW optimizer, as detailed in (Loshchilov and Hutter, 2019). Model performance was assessed using the sacrebleu implementation of the BLEU metric, following (Post, 2018). We used one V100 GPU to train the model. The FLAN-T5-XL model has 3 billion parameters and took 12 hours to train.

For each educational chatbot, we used similar settings: we used the code framework of Daheim et al. (2023) with an initial learning rate of  $6.25e-5$  and linear learning rate decay without warmup.

For model optimization, we utilized checkpoints from the transformers library (Wolf et al., 2020). The negative log-likelihood of the ground-truth response was minimized using the AdamW optimizer, as detailed in (Loshchilov and Hutter, 2019). Model performance was assessed using the sacrebleu implementation of the BLEU metric, following (Post, 2018). We used one V100 GPU to train the model. The FLAN-T5-LARGE model has 0.8 billion parameters and took 7 hours to train.

### A.18 Example of Generated Dialogues

We show examples of generated dialogues from each model in this section. Example dialogue generated by SimSeek is shown in Table 22. Example dialogue generated by Dialogue Inpainting is shown in Table 23. Example dialogue generated by Persona (Low Info) is shown in Table 24. Example dialogue generated by Persona (Medium Info) is shown in Table 25. Example dialogue generated by Persona (High Info) is shown in Table 26. Example dialogue generated by Persona (Single Instance) is shown in Table 27.

**Student:** What is a major fact about physics in the past and present?

**Teacher:** The word physics is thought to come from the Greek word physis, meaning nature.

**Student:** What is the name of a person or thing that has a significance regarding physics in the past and present?

**Teacher:** Albert Einstein

**Student:** How is Einstein regarded as a pioneer of physics in the past and present?

**Teacher:** Einstein devised radical new ideas about time and space.

**Student:** How does Einstein relate to the theory of relativity in the past and present?

**Teacher:** The theory of relativity was developed by Albert Einstein in 1905.

**Student:** What are some other major aspects of physics in the past and present?

**Teacher:** Quantum mechanics is the second major theory of modern physics.

**Student:** What year was Einstein regarded as one of the pioneers of quantum mechanics in the past and present?

**Teacher:** 1905.

Table 22: Example of SimSeek Generated Dialogue

**Student:** What is depth perception?  
**Teacher:** Our ability to perceive spatial relationships in three-dimensional (3-D) space is known as depth perception.  
**Student:** What does that mean?  
**Teacher:** With depth perception, we can describe things as being in front, behind, above, below, or to the side of other things.  
**Student:** Why is this important?  
**Teacher:** Our world is three-dimensional, so it makes sense that our mental representation of the world has three-dimensional properties.  
**Student:** How do we perceive depth?  
**Teacher:** We use a variety of cues in a visual scene to establish our sense of depth.  
**Student:** What kind of cues?  
**Teacher:** Some of these are binocular cues, which means that they rely on the use of both eyes.  
**Student:** What is one example?  
**Teacher:** One example of a binocular depth cue is binocular disparity, the slightly different view of the world that each of our eyes receives.

Table 23: Example of Dialogue Inpainting Generated Dialogue

**Student:** What is the difference between elastic and inelastic collisions?

**Teacher:** In an elastic collision, objects separate after impact and do not lose any kinetic energy. The law of conservation of momentum is applicable in elastic collisions. On the other hand, in an inelastic collision, objects stick together after impact and kinetic energy is not conserved. Some of the kinetic energy gets converted to other forms of energy, like thermal energy.

**Student:** How does the law of conservation of momentum apply in elastic collisions?

**Teacher:** The law of conservation of momentum applies in elastic collisions by stating that the total momentum before the collision is equal to the total momentum after the collision, as long as there is no net external force acting on the system. In an elastic collision, the objects separate after impact and do not lose any of their kinetic energy. This means that the momentum of the system is conserved, even though the objects may change their velocities.

**Student:** Can you explain what happens to the kinetic energy in an inelastic collision?

**Teacher:** In an inelastic collision, the kinetic energy is not conserved. Instead, the kinetic energy of the colliding objects may be converted to other forms of energy, such as thermal energy. This means that some of the energy of motion gets transferred into heat. Even though the objects stick together after the collision, their final velocity may not be the same as their initial velocity.

**Student:** Can you explain why the objects in an elastic collision do not lose any kinetic energy?

**Teacher:** In an elastic collision, the objects do not lose any kinetic energy because they separate after impact. The law of conservation of momentum applies in elastic collisions, stating that the total momentum before the collision is equal to the total momentum after the collision, as long as there is no net external force acting on the system. This conservation of momentum ensures that the kinetic energy of the system remains constant, as the objects do not transfer any energy to other forms. Therefore, the objects in an elastic collision retain all of their initial kinetic energy, resulting in no loss of energy.

**Student:** How does the final velocity of objects in an inelastic collision compare to their initial velocity?

**Teacher:** The final velocity of objects in an inelastic collision is not necessarily the same as their initial velocity. In an inelastic collision, the objects stick together after impact and kinetic energy is not conserved. Some of the kinetic energy gets converted to other forms of energy, such as thermal energy. Therefore, the objects may have a different final velocity compared to their initial velocity.

**Student:** How does the law of conservation of momentum apply in inelastic collisions?

**Teacher:** The law of conservation of momentum applies in inelastic collisions by stating that the total momentum before the collision is equal to the total momentum after the collision, as long as there is no net external force acting on the system. In an inelastic collision, objects stick together after impact and kinetic energy is not conserved. Some of the kinetic energy gets converted to other forms of energy, such as thermal energy. Therefore, the objects may have a different final velocity compared to their initial velocity, but the total momentum of the system remains constant.

Table 24: Example of Persona (Low Info) Generated Dialogue

**Student:** What is the difference between elastic and inelastic collisions?

**Teacher:** An elastic collision is one where the objects separate after impact and don't lose any kinetic energy. Momentum is conserved in an elastic collision. On the other hand, an inelastic collision is one where objects stick together after impact and kinetic energy is not conserved. The forces between colliding objects in an inelastic collision may convert kinetic energy to other forms such as potential energy or heat.

**Student:** Can you give an example of an inelastic collision where objects stick together after impact?

**Teacher:** An example of an inelastic collision where objects stick together after impact is when two objects of equal mass collide head-on and then stick together, coming to rest. This type of collision conserves momentum but not kinetic energy, as some of the energy of motion is converted to thermal energy, or heat.

**Student:** How does an inelastic collision differ from an elastic collision in terms of kinetic energy conservation?

**Teacher:** In an inelastic collision, objects stick together after impact and kinetic energy is not conserved. The forces between colliding objects may convert kinetic energy to other forms such as potential energy or thermal energy. Consequently, some of the energy of motion gets converted into heat. On the other hand, in an elastic collision, the objects separate after impact and no kinetic energy is lost.

**Student:** What happens to the kinetic energy in an inelastic collision?

**Teacher:** In an inelastic collision, the kinetic energy is not conserved. The forces between the colliding objects may convert kinetic energy into other forms such as potential energy or thermal energy. As a result, some of the energy of motion gets converted into heat.

**Student:** How does an elastic collision differ from an inelastic collision in terms of conservation of momentum?

**Teacher:** An inelastic collision differs from an elastic collision in terms of conservation of momentum in that, in an inelastic collision, the objects stick together after impact and kinetic energy is not conserved. In an inelastic collision, the forces between colliding objects may convert kinetic energy to other forms of energy, such as potential energy or thermal energy. This means that some of the energy of motion gets converted into heat.

**Student:** Can you provide an example of an elastic collision where the objects separate after impact?

**Teacher:** An example of an elastic collision where the objects separate after impact is when two steel blocks collide on ice. In this scenario, the objects bounce off each other and retain all of their kinetic energy.

Table 25: Example of Persona (Medium Info) Generated Dialogue

**Student:** What is the difference between elastic and inelastic collisions?

**Teacher:** In an elastic collision, the objects separate after impact and don't lose any of their kinetic energy. On the other hand, in an inelastic collision, the objects stick together after impact and kinetic energy is not conserved. Inelastic collisions involve a conversion of kinetic energy to other forms, such as potential energy or thermal energy.

**Student:** What happens to the objects after an elastic collision?

**Teacher:** In an elastic collision, the objects separate after impact and do not lose any of their kinetic energy. This means that the objects bounce off each other and continue to move independently. However, in an inelastic collision, the objects stick together after impact and kinetic energy is not conserved. This means that the objects combine and move together as one, converting some of their kinetic energy into other forms like potential energy or thermal energy.

**Student:** In an inelastic collision, what forms can kinetic energy be converted into?

**Teacher:** In an inelastic collision, kinetic energy can be converted into other forms of energy, such as potential energy or thermal energy. This means that the forces between colliding objects may cause a loss of kinetic energy, which gets converted into heat.

**Student:** What happens to the objects after an inelastic collision?

**Teacher:** In an inelastic collision, the objects stick together after impact. This means that they combine and move together as one, converting some of their kinetic energy into other forms such as potential energy or thermal energy.

**Student:** What are the different forms that kinetic energy can be converted into during an inelastic collision?

**Teacher:** In an inelastic collision, kinetic energy can be converted into other forms of energy, such as potential energy or thermal energy. This means that the forces between colliding objects may cause a loss of kinetic energy, which gets converted into heat.

**Student:** What is the main difference between elastic and inelastic collisions in terms of how the objects behave after impact?

**Teacher:** In an elastic collision, the objects separate after impact and don't lose any of their kinetic energy. This means that the objects bounce off each other and continue to move independently. On the other hand, in an inelastic collision, the objects stick together after impact and kinetic energy is not conserved. This means that the objects combine and move together as one, converting some of their kinetic energy into other forms like potential energy or thermal energy.

Table 26: Example of Persona (High Info) Generated Dialogue

**Student:** Hi, Professor! I have a question about collisions. When objects collide, can they either stick together or bounce off each other?

**Teacher:** Yes, that's correct. When objects collide, they can either stick together or bounce off one another, remaining separate.

**Student:** I see. So, when they stick together, it's called an inelastic collision, right?

**Teacher:** Exactly. An inelastic collision is one in which objects stick together after impact, and kinetic energy is not conserved.

**Student:** And what about when they bounce off each other?

**Teacher:** When objects bounce off each other, it's called an elastic collision. In an elastic collision, the objects separate after impact and don't lose any of their kinetic energy.

**Student:** Oh, I see. So, in an elastic collision, kinetic energy is conserved?

**Teacher:** Yes, that's correct. In an elastic collision, kinetic energy is conserved. However, in an inelastic collision, kinetic energy is not conserved as it may be converted to other forms of energy.

**Student:** Got it. So, are there any real-life examples of perfectly elastic collisions?

**Teacher:** Perfectly elastic collisions can only happen with subatomic particles. Everyday observable examples of perfectly elastic collisions don't exist as some kinetic energy is always lost due to heat transfer caused by friction.

**Student:** I understand. So, when solving problems involving one-dimensional elastic collisions, we can use the equation for conservation of momentum, right?

**Teacher:** Yes, that's correct. The equation for conservation of momentum for two objects in a one-dimensional elastic collision is  $m_1v_1 + m_2v_2 = m_1v'_1 + m_2v'_2$ , assuming the mass of each object does not change during the collision.

**Student:** Great. And for inelastic collisions, where the objects stick together, is there a different equation?

**Teacher:** Yes, for inelastic collisions, where the objects stick together, we can simplify the conservation of momentum equation to  $m_1v_1 + m_2v_2 = (m_1 + m_2)v'$ , where  $v'$  is the final velocity for both objects as they are stuck together.

Table 27: Example of Persona (Single Instance)  
Generated Dialogue