Consistency Training by Synthetic Question Generation for Conversational Question Answering

Anonymous ACL submission

Abstract

Efficiently modeling historical information is a critical component in addressing user queries within a conversational question-answering (QA) context, as historical context plays a vital role in clarifying the user's questions. However, irrelevant history induces noise in the reasoning process, especially for those questions with a considerable historical context. In our novel model-agnostic approach, referred to as CoTaH (Consistency-Trained augmented History), we augment the historical informa-012 tion with synthetic questions and subsequently employ consistency training to train a model that utilizes both real and augmented historical data to implicitly make the reasoning robust to 016 irrelevant history. To the best of our knowledge, this is the first instance of research using data augmentation to model conversational QA settings. By citing a common modeling 020 error prevalent in previous research, we introduce a new baseline model and compare our model's performance against it, demonstrating 022 an improvement in results, particularly when 024 dealing with questions that include a substantial amount of historical context.

1 Introduction

017

021

037

041

Humans often seek data through an informationseeking process, in which users engage in multiple interactions with machines to acquire information about a particular concept. Prominent examples of this phenomenon include the introduction of Chat-GPT (OpenAI, 2023) and the adoption of industrial systems like Amazon Alexa. Conversational Question-Answering (CQA) systems address user questions within the context of information-seeking interactions. In CQA, unlike conventional question answering, questions are interconnected, relying on previous questions and their corresponding answers to be fully understood without ambiguities. Although many researchers have proposed solutions to model history in CQA, a common modeling mistake made in these studies is using the gold answers of the history instead of the predicted ones. Our work aligns with the framework of addressing irrelevant history, as introduced by Qiu et al. (2021). However, unlike Qiu et al. (2021), our method abstains from utilizing the gold answers of history. Moreover, unlike Qiu et al. (2021), we utilize only one transformer during prediction, resulting in reduced time and memory. Initially, we augment the history of questions in the training set with synthetic questions. Our underlying idea is to maintain the model's consistency in its reasoning, whether utilizing the original historical data or the augmented version. Bert-HAE (Qu et al., 2019a) and HAM (Qu et al., 2019b) have previously served as baselines for several prior methods, but Siblini et al. (2021) conducted a re-implementation of these models using predicted history answers, which resulted in a significant performance decrease. As a result, in this paper, we employ the base transformer of our method as the baseline, as its performance surpasses the re-implementation of the mentioned methods. Our method results in a 1.8% upgrade in overall F1 score, with causing a significant improvement in the scores of questions with a large historical context.

042

043

044

047

048

053

054

056

060

061

062

063

064

065

066

067

068

069

070

071

072

074

076

077

079

2 **Related Works**

The task of CQA has been introduced to extend question answering to a conversational setting. CoQA (Reddy et al., 2019) and QuAC (Choi et al., 2018) have been proposed as two extractive datasets in the CQA task. Bert-HAE (Qu et al., 2019a) employs a manually defined embedding layer to annotate tokens from previous answers within the document, and Qu et al. (2019b) extends this approach introducing an ordering to these annotations. FlowQA (Huang et al., 2019) utilizes multiple blocks of Flow and Context Integration to facilitate the transfer of information between the

170

171

172

173

174

175

176

context, the question, and the history. Qiu et al. (2021) introduces the idea of irrelevant history and its effect in degrading performance, proposing a policy network to select the relevant history before reasoning. However, Qu et al. (2019a,b); Huang et al. (2019); Qiu et al. (2021) employ the gold answers from history in their modeling. This approach deviates from real-world scenarios, where systems should rely on their previous predictions to answer current questions (Siblini et al., 2021). Siblini et al. (2021) re-implements Bert-HAE and HAM using the model's predictions, reporting a sharp decrease in performance. FlowQA experiences a performance drop from 64.4% to 59.0% on the development set when gold answers in history are not used (Huang et al., 2019).

3 Problem Definition

081

087

094

100

101

102

103

104

105

106

107

108

110

111

112

113

114

115

116

117

118

119

To model a CQA setting, at dialog turn k, a model receives a question (q_k) , a document containing the answer (D), and the history of the question (H_k) , which is represented as a set of tuples, such as $H_k = \{(q_0, a_0^{pred}), \dots, (q_{k-1}, a_{k-1}^{pred})\}$, where a_i^{pred} is the model's prediction for q_i . It's important to note that the model may utilize only some of this information. For instance, we only employ history questions while excluding history answers. The primary objective is to predict the answer a_k^{pred} for q_k .

$$a_k^{pred} = \underset{a_k}{\arg\max} P(a_k | q_k, H_k, D)$$
(1)

4 Methodology

We seek to make the reasoning robust to irrelevant history implicitly by augmenting the dataset. To this end, for question q_k , we augment its history by injecting some synthetic questions. Let H_k^* be the augmented history. The intuition is that irrespective of whether the reasoning is performed with H_k or with H_k^* , the result should be the same. In other words:

$$P(a_k|q_k, H_k, D) = P(a_k|q_k, H_k^{\star}, D) \quad (2)$$

To achieve this goal, we establish a two-stage pipeline. Our pipeline consists of a history augmentation module, whose goal is to augment the history and a question-answering module, whose objective is to consistently train a QA network so that the reasoning is consistent. The overall architecture of our model is depicted in Figure 1.

4.1 History Augmentation Module

This module includes a conversational question generator, denoted as CQG_{θ} , where θ represents the parameter set of the generator, and a question selector, denoted as QS, which is responsible for choosing a set of S synthetic questions generated to augment the history.

Training The first step involves training CQG_{θ} . While there has been research aimed at generating conversational questions (Gu et al., 2021; Pan et al., 2019), for the sake of simplifying the implementation, we employ a straightforward generative transformer for this task. To train this network, we input D, H_k , and a_k into the network, intending to generate q_k . We train this network using crossentropy loss in an auto-regressive manner. In 9.2, question generation result is described.

Question Generation After training CQG_{θ} , we aim to generate synthetic conversational questions for the training set. Suppose that we want to generate synthetic conversational questions between q_k and q_{k+1} . Suppose that a_k is located in the *i*-th sentence of the document. We extract noun phrases from sentences i - 1, i, and i + 1 as potential answers. We make this choice as we want these answers to be similar to the flow of conversation and if these answers are extracted from local regions, the likelihood increases. Let one of these answers be called a_k^{syn} . We feed D, H_k , a_k^{syn} to CQG_{θ} to obtain q_k^{syn} . We iterate this process for all $0 \le j \le k$, and generate synthetic questions. We refer to all generated synthetic questions and real questions of the history as pool of questions (P_k) for q_k .

Question Filtering & Injection We could set P_k as H_k^{\star} , however, P_k contains a multitude of synthetic questions which induces too much noise. Additionally, in the consistency training setting, the noise (perturbation) should be small. Thus, we only select S of synthetic questions from P_k , where S is a hyperparameter. Not all synthetic questions are helpful, necessitating the need to filter out degenerate ones. We want our selected synthetic questions to be similar to the trend of the conversation. To this end, we compute a score for each synthetic question and only keep the top Msynthetic questions with the highest score. To compute the score, each question (real or synthetic) is encoded with LaBSE (Feng et al., 2022). For each synthetic question q^{syn} which is located between



Figure 1: Architecture of the Model: For a given question q_k , the conversational question generator CQG_{θ} constructs a pool of questions denoted as P_k . questions in H_k are shown in blue, and synthetic questions are depicted in green and red. The synthetic questions, which are similar to H_k questions, are marked in red, while dissimilar ones are in green. The question selector QS discards red synthetic questions, selects M ones with the highest scores, and chooses S = 3 synthetic questions from the green questions according to uniform distribution, along with H_k questions, to create H_k^* . The QA network $QA_{\theta'}$ computes its output using both H_k and H_k^* as input. The QA network is trained by minimizing L_{CE} and L_{Cons} .

history turns q_i and q_{i+1} , the score is computed 177 as $Sim(h(q_i), h(q^{syn})) + Sim(h(q_{i+1}), h(q^{syn})),$ where Sim is the cosine similarity function and 179 $h(\mathbf{x})$ is the LaBSE's encoding of the sentence \mathbf{x} . 180 Additionally, Sometimes, we generate questions 181 that are too similar to previous or future questions, which are invaluable. Thus, we compare the similarity of generated question q^{syn} with questions in $\{q_k\} \bigcup H_k$ and if the similarity is above γ , q^{syn} 185 is discarded. This situation is depicted in Figure 186 1, where P_k contains real history questions, depicted in blue, and synthetic questions, depicted in 188 red and green. Those synthetic questions that have high similarity with $\{q_k\} \bigcup H_k$, are depicted in red. As it can be seen two questions "Did she have any 191 children" and "How many children did they have" 192 have high similarity with the question "Did they 193 have children", and thus, they're discarded. The 194 effectiveness of question filtering is approved in Section 9.4. In addition, we need to set a distribu-196 tion to guide the selection of S number of generated 197 questions, for which we adopt a uniform distribu-198 tion. More details on the distribution selection are available in Section 9.5. 200

4.2 Question Answering Module

201

204

For each question q_k , as illustrated in Figure 1, we feed q_k , H_k , and D to the QA network $(QA_{\theta'})$ to compute the answer distribution. In parallel, we

feed q_k , H_k^{\star} , and D to the QA network to compute another answer distribution. As mentioned in Section 4, we need to impose the condition outlined in Equation 2. To achieve this, we employ KL-Divergence between the answer distributions. Additionally, we use cross-entropy loss to train the QA network for answer prediction. The losses are calculated as per Equation 3, where L_{CE} , L_{Cons} , and L_T represent the cross-entropy loss, consistency loss, and total loss. λ is a hyperparameter used to determine the ratio of the two losses.

$$L_{CE} = CE(QA_{\theta'}(q_k, H_k, D), a_k^{gold})$$

205

206

207

208

210

211

212

213

214

215

216

217

218

219

220

222

223

224

226

227

229

231

 $L_{Cons} = D_{KL}(QA_{\theta'}(q_k, H_k, D), \qquad (3)$ $QA_{\theta'}(q_k, H_k^{\star}, D))$

$$L_T = L_{CE} + \lambda L_{Cons}$$

Furthermore, we acknowledge that augmenting the history for all questions may not be optimal, as initial questions in a dialog, due to their little historical context, may not require augmentation for robust reasoning. In this case augmenting their history might add unnecessary noise, potentially degrading performance. Thus, we introduce a threshold named τ and only augment the history of q_k if $k \ge \tau$. According to Miyato et al. (2019), we only pass the gradients through one network. As shown in the Figure 1, the symbol \times is used to denote gradient cut. It should be noted that our method is

Model Name	F1	HEQ-Q	HEQ-D	Unrealistic Settings
Bert-HAE-Real (Siblini et al., 2021)	53.5	-	-	
HAM-Real (Siblini et al., 2021)	54.2	-	-	
Bert (Our Model)	58.9	52.9	5.3	
CoTaH-Bert (Our Model)	60.7	55.3	5.9	
Bert-HAE (Qu et al., 2019a)	62.4	57.8	5.1	\checkmark
HAM (Qu et al., 2019b)	64.4	60.2	6.1	\checkmark
Reinforced Backtracking (Qiu et al., 2021)	66.1	62.2	7.3	\checkmark

Table 1: Comparison of our methods with other benchmarks on the test set

model-agnostic, and any architecture could be used as the QA network.

5 Setup

232

240

241

243

246

247

250

251

259

264

268

We utilize the QuAC dataset (Choi et al., 2018), to conduct our experiments on, and data splitting is described in 9.1. We utilize Bert (Devlin et al., 2019) as our base model to conduct experiments following the previous research. For question generation, we adopt Bart-Large (Lewis et al., 2020). Following Choi et al. (2018), we use F1, HEQ-Q, and HEQ-D as our evaluation metrics. F1 measures the overlap between a_k^{gold} and a_k^{pred} . HEQ-Q and HEQ-D are the ratio of questions and dialogs, for which the model performs better than human (Choi et al., 2018). In Section 9.3, the process of choosing all other hyperparameters and their analysis is described. For all of our models, we concatenate the question with history questions, feeding them to the network. More details on reproducibility are presented in Section 9.7.

6 Results

In Table 1, we have depicted our results on the test set division in comparison to previous relevant models. It should be noted that our test set is different from previous methods, but it has been drawn from the same distribution. As stated before, Bert-HAE and HAM leverage the gold answers of history. Their re-implementations by Siblini et al. (2021) are shown in the Table as Bert-HAE-Real and HAM-Real, which indicate a significant drop in performance. In this scenario where common baselines experience a substantial decrease, we examine a basic Bert model with history concatenation as the baseline, as its performance is superior. Our model outperforms this baseline by 1.8% in the F1 score. According to Figure 2, this improvement is mostly due to an improvement in

the performance of questions with a large amount of history. This confirms that our intuition is valid that our method enhances the base model's ability to answer questions with a large historical context. Moreover, while Bert-HAE outperforms CoTaH-Bert in terms of F1 score, CoTaH-Bert exhibits superior performance in HEQ-D. This highlights the better consistency of our model to maintain its performance throughout the entire dialog, which is achieved through superiority in answering the questions in the latter turns. Additionally, we include the results of the history backtracking model (Qiu et al., 2021) in the table. Since this model's code is not publicly available, we have been unable to reimplement it with the correct settings and perform a meaningful comparison. However, it's worth noting that this model utilizes unrealistic settings in two stages: once for history selection and once for question answering, potentially exacerbating the modeling issues even further. We have used "Unrealistic Settings" as a term to indicate that a model uses gold answers of history in its modeling.

269

270

271

273

274

275

276

277

278

279

281

282

283

285

286

290

291

292

293

294

295

296

297

298

299

300

301

302

303

304

305

306

7 Conclusions

In this paper, we introduced a novel model-agnostic method to make the reasoning of conversational question-answering models robust to irrelevant history. We coped with this issue by augmenting the history and training the model with consistency training. In our experiments, we didn't follow the wrong modeling of past research in using the gold answers of history. We examined our method with Bert which exhibited a 1.8% performance boost compared to the baseline model. It was demonstrated that this improvement is primarily attributed to the enhancement of the model's performance on questions with a substantial historical context, suggesting that our method has been successful in making the reasoning robust for these questions.

8 Limitations

307

323

327

328

329

330

331

335

336

337

338

339

341

342

343

345

346

347

351

352

353

354

359

Our model requires a phase of question generation. For synthetic question generation, the history augmentation module could be slow and the speed is directly correlated to the number of questions that one opts to generate. However, question generation 312 313 is trained only once and all questions are generated in a single run, and all of other experiments are 314 conducted by only training the QA module. Moreover, although our model doesn't need any further computation during evaluation than merely running 317 the QA network, we need two forward passes during the training phase, which makes the training of 319 the QA network a bit more time-consuming than 321 training the baseline model.

References

- Eunsol Choi, He He, Mohit Iyyer, Mark Yatskar, Wentau Yih, Yejin Choi, Percy Liang, and Luke Zettlemoyer. 2018. Quac: Question answering in context. In Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Brussels, Belgium, October 31 November 4, 2018, pages 2174–2184. Association for Computational Linguistics.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: pre-training of deep bidirectional transformers for language understanding. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers), pages 4171–4186. Association for Computational Linguistics.
- Fangxiaoyu Feng, Yinfei Yang, Daniel Cer, Naveen Arivazhagan, and Wei Wang. 2022. Languageagnostic BERT sentence embedding. In Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2022, Dublin, Ireland, May 22-27, 2022, pages 878–891. Association for Computational Linguistics.
- Jing Gu, Mostafa Mirshekari, Zhou Yu, and Aaron Sisto. 2021. Chaincqg: Flow-aware conversational question generation. In Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume, EACL 2021, Online, April 19 - 23, 2021, pages 2061–2070. Association for Computational Linguistics.
- Hsin-Yuan Huang, Eunsol Choi, and Wen-tau Yih. 2019. Flowqa: Grasping flow in history for conversational machine comprehension. In 7th International Conference on Learning Representations,

ICLR 2019, New Orleans, LA, USA, May 6-9, 2019. OpenReview.net.

- Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. BART: denoising sequence-to-sequence pretraining for natural language generation, translation, and comprehension. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, ACL 2020, Online, July 5-10, 2020*, pages 7871–7880. Association for Computational Linguistics.
- Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*, pages 74–81.
- Takeru Miyato, Shin-ichi Maeda, Masanori Koyama, and Shin Ishii. 2019. Virtual adversarial training: A regularization method for supervised and semisupervised learning. *IEEE Trans. Pattern Anal. Mach. Intell.*, 41(8):1979–1993.
- OpenAI. 2023. GPT-4 technical report. CoRR, abs/2303.08774.
- Boyuan Pan, Hao Li, Ziyu Yao, Deng Cai, and Huan Sun. 2019. Reinforced dynamic reasoning for conversational question generation. In *Proceedings of the 57th Conference of the Association for Computational Linguistics, ACL 2019, Florence, Italy, July 28- August 2, 2019, Volume 1: Long Papers*, pages 2114–2124. Association for Computational Linguistics.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics, July 6-12, 2002, Philadelphia, PA, USA, pages 311–318. ACL.
- Minghui Qiu, Xinjing Huang, Cen Chen, Feng Ji, Chen Qu, Wei Wei, Jun Huang, and Yin Zhang. 2021. Reinforced history backtracking for conversational question answering. In Thirty-Fifth AAAI Conference on Artificial Intelligence, AAAI 2021, Thirty-Third Conference on Innovative Applications of Artificial Intelligence, IAAI 2021, The Eleventh Symposium on Educational Advances in Artificial Intelligence, EAAI 2021, Virtual Event, February 2-9, 2021, pages 13718–13726. AAAI Press.
- Chen Qu, Liu Yang, Minghui Qiu, W. Bruce Croft, Yongfeng Zhang, and Mohit Iyyer. 2019a. BERT with history answer embedding for conversational question answering. In *Proceedings of the 42nd International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR* 2019, Paris, France, July 21-25, 2019, pages 1133– 1136. ACM.
- Chen Qu, Liu Yang, Minghui Qiu, Yongfeng Zhang, Cen Chen, W. Bruce Croft, and Mohit Iyyer. 2019b.

395

396

397

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

360

361

362

- 416 417 418
- 419 420
- 421
- 422 423
- 424
- 425 426 497 428
- 429 430 431
- 432
- 433 434 435 436
- 437

438

- 439 440
- 441 442
- 443
- 444 445
- 446 447
- 448 449
- 450
- 451 452

453

- 454
- 455

- 456
- 457 458

459

460

461

462

463

464

465

tion answering. In Proceedings of the 28th ACM International Conference on Information and Knowledge Management, CIKM 2019, Beijing, China, November 3-7, 2019, pages 1391–1400. ACM. Siva Reddy, Danqi Chen, and Christopher D. Manning.

2019. Coqa: A conversational question answering challenge. Trans. Assoc. Comput. Linguistics, 7:249-266.

Attentive history selection for conversational ques-

- Wissam Siblini, Baris Sayil, and Yacine Kessaci. 2021. Towards a more robust evaluation for conversational question answering. In Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing, ACL/IJCNLP, Virtual Event, pages 1028–1034.
- Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. 2020. Bertscore: Evaluating text generation with BERT. In 8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020. OpenReview.net.

9 Appendix

9.1 Data Splitting

Since the test set of QuAC is not publicly available, we divide the development (dev) set into dev/test sets randomly, such that the number of questions in dev and test sets is almost equal. The total number of dev and test questions is 3678 and 3676 respectively after splitting. In our splitting, each dialog, with all of its questions, is either attributed to the dev set or the test set, in order to prevent test data leakage. Further, according to Choi et al. (2018), original dev set of QuAC contains unique documents, meaning that a single document will not be shared among the final dev and test sets, potentially preventing test data leakage.

9.2 **Question Generation Results**

The results of question generation are evaluated in Table 2. These scores are obtained from the dev data. Bleu-1,4 (Papineni et al., 2002), Rouge (Lin, 2004), and Bert-Score (Zhang et al., 2020) are used for criteria. We use the evaluate library¹ to implement these metrics. Gu et al. (2021) reports better results for the question generation, yet we didn't aim to optimize Bart-Large meticulously as the generated questions have a good quality for our task. The point is that in this research, we only utilize questions alone without considering answers. Thus, if the generated questions have less

¹https://github.com/huggingface/evaluate

correlations with answers, it's tolerable as they are 466 still relevant questions considering the overall flow 467 of conversation. it should be noted that if a future 468 research wants to incorporate predicted answers in 469 its modeling, it should be more cautious about the 470 quality of the question generation to ensure that the 471 right synthetic questions are generated concerning 472 their answers. 473

Table 2: Question generation results on the dev set

Bleu-1	Bleu-4	Rouge-L	Bert-Score
33.6	9.5	29.0	90.5

Hyperparameter Selection & Sensitivity 9.3 Analysis

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

Initially, we determine M and γ by assessing some examples of the dev data, setting M = 10 and $\gamma = 0.8$ based on our appraisal. Next, we determine the values of S, λ , and τ by conducting experiments on the dev set. In Table 3, we evaluate the effects of the model's two main hyperparameters, S and λ , through a grid search with the following values: $S \in \{1, 2, 3\}$ and $\lambda \in \{1.0, 1.5, 2.0\}$. Firstly, it is evident that the model performs better when $S \in 1, 2$ compared to when S = 3 overall. This suggests that S = 3 introduces too much noise, which could be detrimental for performance. Furthermore, when $\lambda \in 1.5, 2.0$, the performance is better compared to $\lambda = 1.0$, indicating that the introduction of λ is helpful, as simply adding L_{CE} and L_{KL} (or equally setting $\lambda = 1.0$) produces inferior performance. For the remaining experiments, we set S = 2 and $\lambda = 2.0$ as these settings yield the best F1 and HEQ-Q scores.

Table 3: The effect of S and λ on the dev set

		F1	HEQ-Q	HEQ-D
S = 1	$\begin{vmatrix} \lambda = 1.0 \\ \lambda = 1.5 \end{vmatrix}$	58.6	53.5 54.8	4.8 5.5
<i>D</i> – 1	$\lambda = 1.0$ $\lambda = 2.0$	59.0	54.2	4.4
	$\lambda = 1.0$	57.9	52.7	4.0
S=2	$\lambda = 1.5$	58.2	53.5	4.2
	$\lambda = 2.0$	59.4	54.8	5.1
	$\lambda = 1.0$	58.3	53.5	5.1
S=3	$\lambda = 1.5$	58.6	53.5	5.0
	$\lambda = 2.0$	58.8	54.1	4.2

After setting the right amount for S and λ , we 495 opt to examine whether the introduction of τ is 496 effective. Thus, we conduct experiments on three 497 different amount of this hyperparameter. In Table 498 4, it's evident that the right amount of τ has a considerable effect on the performance, confirming our 500 intuition about the functionality of τ . For all tested 501 values of τ within the set $\{5, 6, 7\}$, performance has increased compared to the base settings with $\tau = 0$ (or equivalently, using no threshold). No-504 tably, the maximum performance improvement is 505 observed when $\tau = 6$. 506

Table 4: The effect of τ on the dev set

	F1	HEQ-Q	HEQ-D
$\tau = 0$	59.4	54.8	5.1
$\tau = 5$	59.6	55.2	5.5
$\tau = 6$	59.9	55.2	5.5
$\tau = 7$	59.5	54.9	5.1

9.4 Question Filtering Effect

508

509

510

511

512

513

After determining the optimal τ , the effectiveness of the question-filtering, as discussed earlier, is examined. The results in Table 5 demonstrate that this filtering leads to a considerable additional performance boost by filtering out degenerate questions.

Table 5: The effect of question filtering on the dev set

Filtering Type	F1	HEQ-Q	HEQ-D
No Filtering	59.9	55.2	5.5
Similarity Filtering	60.9	56.3	5.3

9.5 Synthetic Question Selection Distribution

Although we select synthetic questions using a uniform distribution, we have conducted experiments 515 using two distributions: uniform and linear. In 516 the uniform setting, the generated questions are se-517 lected with the same probability. For the linear, if 518 q^{syn} is located between q_j and q_{j+1} , its probability of being selected $(P(q^{syn}))$ is $P(q^{syn}) \propto k - j$. We opt for the linear distribution, as we believe that closer synthetic questions to the original question 522 might contribute to greater robustness, as questions 524 that are further away are likely less relevant. The results are shown in Table 6. We observe a relatively 1% drop in both F1 and HEQ-Q scores with the linear distribution, concluding that our hypothesis has not been true. Given the superiority of the 528

uniform distribution, we choose to continue with it.

Table 6: The effect of question selection distribution on the dev set

Q-Selection Dist.	F1	HEQ-Q	HEQ-D
Uniform	60.9	56.3	5.3
Linear	59.9	55.2	5.9

9.6 Additional Results

In Figure 2, a comparison between the F1 scores of questions for each turn in Bert and CoTaH-Bert on the test set is presented. The score for the k-th turn represents the average F1 score for all questions in the k-th turn across all dialogs in the test set. Questions with a considerable amount of historical context are answered more effectively with our method. For $0 \le k \le 1$, the performances of both Bert and CoTaH-Bert are nearly equal, which is sensible as these questions contain little historical context, and thus, they have little irrelevant history. However, for most of k > 1 dialog turns, CoTaH-Bert outperforms Bert or it has on par performance with Bert. The performance upgrade is especially evident towards the end of dialogs, where questions contain significant historical context. This finding indicates the superiority of CoTaH-Bert over Bert in establishing greater robustness in answering these questions, by identifying and ignoring the irrelevant history turns.



Figure 2: The F1 score of the test set dialog turns

9.7 Reproducibility

The seed for all experiments, except the training of CQG_{θ} , is 1000. All of the experiments to train

529 530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554



Figure 3: A comparison between Bert and CoTaH-Bert extracted answers to a question, showing that CoTaH-Bert has been able to successfully ignore the irrelevant history by extracting the correct answer. However, the Bert model has been confused and returned a wrong answer.

the $QA_{\theta'}$ are conducted on a single RTX 3070 Ti with 8GB memory, on which each experiment takes approximately 6 hours. CQG_{θ} is trained on a single Tesla T4 from Google Colab. For each model, Bert or CoTaH-Bert, the hyperparameters are optimized on the dev set, and a final model will be trained on the train set with the optimized hyperparameters. Subsequently, a single result on the test set will be reported as depicted in Table 1.

9.8 Case Study

In Figure 3, a document sample with its corresponding dialog in the dev set is depicted. In the figure, ninth turn question, q_9 , with its history, H_9 , are shown. The answers of Bert and CoTaH-Bert to q_9 are compared, showing that CoTaH-Bert has been successful to answer this question with a full F1 score, while Bert has been unsuccessful. q_9 asks about the release date of the album stated in q_2 . This is a suitable sample for our context, as there are significant irrelevant history turns between q_9 and q_2 . We observe that CoTaH-Bert has been successful in identifying the relevant history by answering the question correctly. However, the Bert model has mistakenly reported another date which is wrong. As Bert has returned a span containing the word "mixing", it's possible that Bert has incorrectly identified the previous turn question, q_8 , as

570

574

575

578

579

580

567

555

582	relevant, and has returned a span by text matching
583	encompassing the word "mixing", and containing
584	merely some random dates.