INVITE: a Testbed of Automatically Generated Invalid Questions to Evaluate Large Language Models for Hallucinations

Anil Ramakrishna Rahul Gupta Jens Lehmann Morteza Ziyadi

Amazon Alexa AI

{aniramak, gupra, jlehmnn, mziyadi}@amazon.com

Abstract

Recent advancements in Large language models (LLMs) have enabled them to hold free form conversations over multiple turns, but they exhibit a tendency to make unfounded and incorrect statements, commonly known as hallucinations. In particular, LLMs hallucinate frequently when given invalid questions, i.e. ones with incorrect assumptions. The most common approach to evaluate LLMs on hallucinations is to test them on Question Answering (QA) test sets such as TruthfulQA. However, LLMs are increasingly pretrained on massive text corpora scraped from the Internet, which may inevitably expose these test sets to the model during training, leading eventually to an overestimation of model performances on these test sets. In this work, we present an alternative framework to address this risk and to foster further research towards making LLMs robust against invalid questions. We name our framework INVITE: a testbed of automatically generated INValId questions to evaluaTE large language models for hallucinations. In each instantiation, our framework is set up to create a fresh batch of invalid questions by distorting valid facts in which subjects or objects are replaced by similar entities. We evaluate several state of the art LLMs against a testset generated by our framework and highlight its capacity to trigger hallucinations in these models.

1 Introduction

Despite their recent success, LLMs have long been known to exhibit several patterns of concern (Weidinger et al., 2021) such as generating statements which may be toxic (Ousidhoum et al., 2021), biased (Ferrara, 2023), unfair (Ramesh et al., 2023) and factually incorrect (Azamfirei et al., 2023). The last pattern of generating factually incorrect yet seemingly confident statements is commonly labeled as *hallucinations* in the literature (Ji et al., 2023). It is an important area of study since the confident tone of these generated statements can lead

to end users accepting them as accurate without any subsequent validation.

Model hallucinations occur in a variety of textual generative applications such as NLG, MT, QA, dialog systems, data to text systems, etc. It is believed to be caused by discrepancies in data used to train the models, or in the model training itself (Ji et al., 2023). Hallucinations are also believed to be caused by the supervised fine tuning process in which the model may learn to make factually ungrounded connections within its parametric memory in order to accurately answer the current question it is being trained on, which can trigger new ungrounded responses as hallucinations during inference.

Typical approaches to evaluate newly developed models for hallucinations have been to test them on Question Answering datasets such as TruthfulQA (Lin et al., 2021), which provides a curated set of challenging questions with valid answers, against which the model generated responses are compared. However, this approach of using a fixed test set with LLMs is inherently limited; the typical development cycle of a new LLM release involves pretraining on large text corpora regularly scraped from the Internet, and any new challenge dataset may eventually get scraped into this pre-training corpus. Given that LLMs have been shown to memorize training data (Carlini et al., 2022), this form of data leakage can lead to a false sense of improvement on the challenge test set in subsequent model releases. To address this risk, in this work, we instead propose to test LLM models using an evaluation framework which uses carefully crafted rules to create new challenge questions in each round. We call our framework INVITE: a testbed of automatically generated INValId questions to evaluaTE large language models for hallucinations. INVITE leverages valid facts from knowledge bases to create new invalid questions which may not have an answer. Our framework can be used to evaluate

new LLM release candidates on their robustness against invalid questions which can trigger specific forms of hallucinations, as well as when developing new algorithms to mitigate hallucinations in existing models. The key contributions of our work are as follows:

- We create a new framework to create invalid questions to evaluate robustness of LLMs against hallucinations¹.
- We test our framework on several latest LLMs, exploring different model sizes and training datasets.
- We conduct a pilot human evaluation study on the generated responses for these questions, and highlight the effectiveness of the test sets in triggering hallucinations in the models being evaluated.

2 Related Work

Question Answering datasets A number of QA datasets are available in literature to test LLMs for hallucinations, including TruthfulQA (Lin et al., 2021), SQUAD (Rajpurkar et al., 2016), TriviaQA (Joshi et al., 2017), among others. While challenging and effective, all of these test sets suffer from the previously described risk of possibly getting scraped and consumed in model training.

Adversarial QA datasets The QA datasets listed above test models in their ability to retain facts, and new models are tested by comparing their response to these questions against an expected answer. However, we argue that this strategy alone would not test a model against all possible failure modes related to hallucinations since this issue stems from the models' ability to concoct new facts as statements. Even if a model were to learn the correct answer to a particular QA question, it can still retain the general tendency to hallucinate. Hence, a more robust strategy would be to also test these models using adversarial questions with invalid assumptions. A model capable of avoiding hallucinations would identify that it does not have an answer to the question, or detect that the question itself is not plausible and that it cannot generate an answer, and hence choose to disengage.

Notable adversarial datasets in NLP literature include (Jia and Liang, 2017), where the authors

Source Dataset	Question Categories		
DBpedia	almaMater, associatedBand, author, award, birthPlace, city, commander, country, musicComposer, office, party, position, predecessor, pub- lisher, spouse, successor, team, writer		
TriviaQA	InvalidDate, FutureDate		

Table 1: Question Categories

used a rule based framework to add adversarial statements to passages from reading comprehension task in order to confuse target models. Subsequently, Rajpurkar et al. (2018) developed SQUAD 2, a richer set of unanswerable questions using human annotators for the same reading comprehension task. While rich in diversity and volume, these datasets still suffer from the same risk of getting consumed in model training noted above. Further, their datasets are limited to reading comprehension tasks and hence do not necessarily test the full boundary of a model's knowledge.

Automated Testset Generation Automatically created unit tests have been explored in deterministic applications such as software testing (Chen et al., 2022; Schäfer et al., 2023). With machine learning models such as LLMs, we can leverage a variety of generative models to create new datasets (Duan et al., 2017; Nikolenko, 2019), but to the best of our knowledge, no prior works have tried to generate questions with invalid assumptions. Our proposed approach addresses this gap by setting up a framework to automatically create a new test set of challenge questions with verifiably invalid assumptions, which are likely to trigger hallucinations in the target model.

3 The INVITE Framework

To create new test questions, we collect valid facts from a knowledge base and distort these to create new unanswerable questions. We describe our process in more detail below, and create a test set using this framework.

3.1 Creating Invalid Questions

We use the DBpedia knowledge base (Lehmann et al., 2015) as a source of valid facts to create our questions. The choice of knowledge base here is arbitrary and can be replaced by an alternate appli-

¹Full code available at https://github.com/amazon-science/invite-llm-hallucinations.

cation specific knowledge base as necessary. DBpedia extracts structured factual information from Wikipedia, the world's largest encyclopedia. It contains a large volume of facts which are stored in the Resource Description Framework (RDF) format of subject-predicate-object triples. The most recent release of DBpedia contains over 850 million such factual triples (Holze), making it a decidedly rich source of information to create new test questions for our task. Of these, we use a subset of 42 million triples containing facts about objects and literals extracted from the Wikipedia Infoboxes, which are reported to be of higher quality because of their standardized format. For operational simplicity, we limit our scope here to the 100 most frequent predicate types by volume from this subset. We further discard noisy predicate types which contain ambiguous entries after manual inspection (for example, we discard the nationality predicate type since it contains answers of the form country-name as well as citizen/people of country-name, making facts of this type difficult to fit in a consistent question template). The exact list of predicates selected in our dataset creation is listed in Table 1.

To create new questions, we first curated over 300 predicate specific template questions which were manually crafted by annotators on Amazon Mechanical Turk, and further denoised by the authors. Next, we further refined a subset of these to create high quality question templates by posing the questions on a search engine, and iterated this process until the responses were unambiguous. We also created template answers for these selected high quality question templates, which we use in our subsequent experiments reported below. The specific prompts we used in our experiments are listed in Appendix A^2 .

For each new question generation, given a predicate, we first sample a fact triple from this predicate type and create a valid question using the corresponding template. Next, to create the invalid question, we create an invalid fact triple by sampling new subjects or objects found in facts from the same predicate type. We verify that this new triple does not exist as factual predicate in the dataset; if such a triple exists, then our created fact is actually valid, so we discard the same and repeat the sampling process above until we have an invalid triple. Given this (invalid) triple, we use our template for

Model	Hallucination Rate
GPTNeo-2.7B	83%
GPTJ-6B	82%
Open-LLaMA-7B	88%
RedPajama-7B	81%
GPT3.5-Turbo	17%
GPT4	6%

Table 2: Model specific hallucination rates on a test set of invalid questions (results sorted by model size).

this predicate type to create a new invalid question and a corresponding answer, subsequently adding both to our test set.

Questions with Invalid Dates In addition to the questions extracted above, we create two more categories containing questions with invalid dates. Using regular expressions, we sample questions containing dates and years from the TriviaQA dataset's test set (Joshi et al., 2017) and create various distortions before adding these questions to our test set. Specifically, we distort full dates containing months by randomly selecting a new date beyond valid dates of the month (for example: March 32nd, 2023) and replace the old date. Similarly, we distort years by randomly sampling a new year from [2025, 2100] and replace the old year.

4 Experiments

Testing model responses for hallucinations is a challenging task which needs a comprehensive fact verification system for automated evaluations. We instead use human verification to test for hallucinations in the generated responses. To evaluate the efficacy of our proposed framework, we first created a pilot test set of 100 questions, sampling uniformly from each category listed in Table 1. Next we generated responses to each of these questions using the models described below, leading to a total set of 600 generated responses. Finally, we manually examine these generations and label them for hallucinations, utilizing a search engine for additional validation of model responses. While manually labeling samples, we only treat responses which explicitly make an inaccurate statement as hallucinations, treating all others (including empty or degenerate responses) as non-hallucinations.

4.1 Models

We evaluate the test set described in Section 4 on a list of open source and proprietary large language

²Our curated question templates, along with model responses with labels can be downloaded from https://github.com/amazon-science/invite-llm-hallucinations.

Model	BLEU	METEOR	ROUGE		BERTScore	AlianScore
			ROUGE-1	ROUGE-L	DERISCOL	Alighiscole
GPTNeo-2.7B	0.0106	0.1909	0.0925	0.0896	0.4249	0.2073
GPTJ-6B	0.0173	0.2336	0.1134	0.1099	0.4309	0.3781
Open-LLaMA-7B	0.0301	0.3311	0.2448	0.2361	0.5415	0.4503
RedPajama-7B	0.0024	0.0688	0.0388	0.0361	0.3739	0.2699
GPT3.5-Turbo	0.0711	0.4784	0.3362	0.3207	0.6460	0.7008
GPT4	0.0362	0.3748	0.2510	0.2381	0.5999	0.7795

Table 3: Automated Metrics between generated responses and references.

models described below. We chose a diverse set of models with varied size, and training datasets for a detailed evaluation of our test set. All open source models were downloaded from Huggingface and evaluated on Nvidia A100 Tensor Core GPUs, while the proprietary GPT models were evaluated using OpenAI APIs³. We ran inference without decoder sampling to further reduce the models' tendency for hallucinations, and stopped inference after 150 tokens.

GPT-Neo-2.7B GPT-Neo (Black et al., 2021) is a 2.7 billion parameter model developed by EleutherAI, and it follows the architecture of GPT-3. It was trained on the Pile (Gao et al., 2020), a large-scale dataset curated by EleutherAI for this task, which spans diverse tasks.

GPT-J-6B GPT-J-6B is 6 billion parameter model trained using Mesh Transformer JAX (Wang, 2021), and also trained on the Pile dataset from EleutherAI.

Open-LLaMA-7b-Open-Instruct This is an instruction tuned, open sourced release of the 7 billion parameter LLaMA model (Touvron et al., 2023), trained on the Open-Instruct-v1 dataset which consists of 63000 instruction training samples.

RedPajama-INCITE-7B-Instruct The RedPajama models were developed by a team of open source developers from several organizations. The base model was trained on the RedPajama dataset, a 1T token open-source implementation of the LLaMA dataset. Several model variants were available at the time of writing, and we used the 7B instruction tuned version of the model in our evaluations.

GPT models We also ran evaluations on OpenAI's GPT3.5-Turbo (OpenAI, a) and GPT4 (OpenAI, b) models. GPT3.5-Turbo is a text only model which supports understanding and generation of natural language and code, while GPT4 is OpenAI's most powerful LLM at the time of this writing. We used model snapshots from June 13, 2023 in both cases. These models are first pretrained on a large corpus and subsequently aligned using Supervised Fine Tuning (SFT) and Reinforcement Learning from Human Feedback (RLHF) (Ouyang et al., 2022). Additional details have not been disclosed by OpenAI.

5 Results

We present our hallucination rates for the models described above in Table 2. As observed in the table, all the open source models exhibit significant rates of hallucinations triggered by invalid assumptions in the test questions. Examples of the generated hallucinations from each model are shown in Table 4. Among all the open source models, RedPajama-7B exhibited the least hallucination rate at 81% while LLaMA-7B demonstrated the highest rate of hallucinations at 88%, closely followed by GPTNeo-2.7B (83%) and GPTJ-6B (82%).

On the other hand, OpenAI's GPT3.5-Turbo and GPT4 models showed considerably lower (yet nontrivial) hallucination rates at 17% and 6% respectively. Since OpenAI does not release internal details of these models, it is unclear if they incorporate any form of external knowledge retrieval which would provide an effective solution to handle questions with invalid assumptions. However, given the growing popularity of LLMs and the demand to deploy them in various low latency regimes, smaller ungrounded open source models such as those evaluated above are likely production candidates in many applications, and INVITE provides a valuable framework to evaluate these for hallucinations, without relying on a fixed test set which runs the risk of getting ingested in future training data crawls. Further, even in GPT4 and GPT3.5-Turbo,

³Accessed on August 27, 2023

Model	Hallucinated Examples				
CDTI (D	q: What was the name of Heather Sears's character in Gooische Vrouwen?				
UF IJ-0D	<i>a</i> : The name of Heather Sears's character in Gooische Vrouwen is Mieke.				
CPTNeo 27B	q: When did Michael Stephenson write Northern Touch?				
UF INCO-2.7D	<i>a</i> : Northern Touch is a novel by Michael Stephenson, published in 2007. It is the				
	first of a trilogy, and the second book in the series is The Last of the Breed. The				
	novel is set in the fictional town of Ballymote, in County Mayo, Ireland. The town				
	is a fictionalised version of the town of Ballymote, County Mayo, in Ireland. The				
	novel is set in the late 1950s, and is the first of a trilogy, and the second book in the				
	series is The Last of the Breed.				
Open-LLaMA-	q: When did Terry Sabo succeed Fernand Roberge?				
7B	a: Terry Sabo was elected to the Maine House of Representatives in 1988 and				
	served until 1994. In 1992, he was elected to the Maine Senate and served until				
	1996. In 1994, he was elected to the U.S. House of Representatives and served				
	until 2002. Fernand Roberge served as the U.S. Representative for Maine's 1st				
	congressional district from 1973 to 1987. Terry Sabo succeeded him in the U.S.				
	House of Representatives in 1987.				
RedPajama-7B	q: When did Rolf Presthus join Democratic Party?				
	<i>a</i> : Rolf Presthus joined the Democratic Party in the Norwegian Parliament in 2005.				
GPT3.5-Turbo	q: When did Sherard Cowper-Coles succeed Jack Marshall?				
	<i>a</i> : Sherard Cowper-Coles succeeded Jack Marshall as the British Ambassador to				
	Saudi Arabia in 2003.				
GPT4	q: When did Edmund Blampied marry Pamela Lincoln?				
	<i>a</i> : Edmund Blampied married Pamela Lincoln in 1920.				

Table 4: Examples of hallucinated model generations on test questions with invalid assumptions.

identifying the 6% and 17% examples which were hallucinated highlights existing gaps in these models and can provide valuable insights needed to make them safe for deployment.

5.1 Automated Evaluation Metrics

Since human labeling is slow and expensive, they may not always be feasible to obtain. As an alternative, we also computed various automated metrics commonly employed in evaluating natural language generations. To estimate these metrics, we compare model generations to the invalid questions against reference answers created using the category wise answer templates listed in Appendix A. As additional references, we also included two disengagement answers: I don't know and I can't answer that for comparison. Table 3 lists estimates for these metrics. From our experiments, n-gram count based metrics such BLEU (Papineni et al., 2002), METEOR (Banerjee and Lavie, 2005) and ROUGE (Lin, 2004) do not correlate well with human labeled estimates of hallucination rate. Model based metrics such as BERTScore (Zhang* et al., 2020) and AlignScore (Zha et al., 2023) perform relatively better than n-gram based metrics as shown in Table 3, but they still do not perfectly align with gold standard labels from human labeling, which appears to be the most reliable estimate of whether a model response is hallucinated.

6 Conclusion

We developed a new framework called INVITE to evaluate large language models for hallucinations, in which new test questions are automatically generated in each round, thereby avoiding reliance on fixed test sets which carry the risk of getting ingested in future training corpora. Our framework creates a diverse (in both domains and entities) set of questions, obtained by distorting valid factual triples from a knowledge base. It is also flexible and easily extensible to new knowledge bases and predicate types. We evaluate an example test set generated by our framework against several state of the art LLMs, establishing the challenging nature of questions generated by our framework. Implementation of our framework, along with the curated question templates and labeled model responses are being released with the paper.

Limitations and Future Work

Any test set of limited size would not cover the entire possible space of invalid questions. Instead, we chose to sample a random subset of this space and obtain an empirical estimate of the model performance.

We define an invalid fact triple/relationship as one which does not exist in the knowledge base and this assumption maybe violated in boundary cases where facts may not have been entered into Wikipedia; however, we expect this to be marginal.

Generating invalid questions from a fixed set of templates may lead to limited diversity in questions. We used annotators from Amazon Mechanical Turk to address this but our coverage was limited to a relatively small set of predicates, which we will expand on in future work.

In current version of the framework, we only create single hop questions with two entities for simplicity. We will expand on this by creating multi-hop questions which necessitate complex reasoning in future work.

Finally, in future work we can also leverage this framework in LLM training by using such questions to teach the boundary of plausible knowledge to the model.

References

- Razvan Azamfirei, Sapna R Kudchadkar, and James Fackler. 2023. Large language models and the perils of their hallucinations. *Critical Care*, 27(1):1–2.
- Satanjeev Banerjee and Alon Lavie. 2005. METEOR: An automatic metric for MT evaluation with improved correlation with human judgments. In Proceedings of the ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization, pages 65–72, Ann Arbor, Michigan. Association for Computational Linguistics.
- Sid Black, Gao Leo, Phil Wang, Connor Leahy, and Stella Biderman. 2021. GPT-Neo: Large Scale Autoregressive Language Modeling with Mesh-Tensorflow. If you use this software, please cite it using these metadata.
- Nicholas Carlini, Daphne Ippolito, Matthew Jagielski, Katherine Lee, Florian Tramer, and Chiyuan Zhang. 2022. Quantifying memorization across neural language models. *arXiv preprint arXiv:2202.07646*.
- Bei Chen, Fengji Zhang, Anh Nguyen, Daoguang Zan, Zeqi Lin, Jian-Guang Lou, and Weizhu Chen. 2022. Codet: Code generation with generated tests. arXiv preprint arXiv:2207.10397.

- Nan Duan, Duyu Tang, Peng Chen, and Ming Zhou. 2017. Question generation for question answering. In Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing, pages 866–874, Copenhagen, Denmark. Association for Computational Linguistics.
- Emilio Ferrara. 2023. Should chatgpt be biased? challenges and risks of bias in large language models. *arXiv preprint arXiv:2304.03738*.
- Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason Phang, Horace He, Anish Thite, Noa Nabeshima, Shawn Presser, and Connor Leahy. 2020. The pile: An 800gb dataset of diverse text for language modeling.
- Julia Holze. Dbpedia snapshot 2022-12 release. Accessed: 2023-06-18.
- Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea Madotto, and Pascale Fung. 2023. Survey of hallucination in natural language generation. ACM Computing Surveys, 55(12):1–38.
- Robin Jia and Percy Liang. 2017. Adversarial examples for evaluating reading comprehension systems. *arXiv preprint arXiv:1707.07328*.
- Mandar Joshi, Eunsol Choi, Daniel S Weld, and Luke Zettlemoyer. 2017. Triviaqa: A large scale distantly supervised challenge dataset for reading comprehension. *arXiv preprint arXiv:1705.03551*.
- Jens Lehmann, Robert Isele, Max Jakob, Anja Jentzsch, Dimitris Kontokostas, Pablo N Mendes, Sebastian Hellmann, Mohamed Morsey, Patrick Van Kleef, Sören Auer, et al. 2015. Dbpedia–a large-scale, multilingual knowledge base extracted from wikipedia. *Semantic web*, 6(2):167–195.
- Chin-Yew Lin. 2004. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain. Association for Computational Linguistics.
- Stephanie Lin, Jacob Hilton, and Owain Evans. 2021. Truthfulqa: Measuring how models mimic human falsehoods. *arXiv preprint arXiv:2109.07958*.
- Sergey I Nikolenko. 2019. Synthetic data for deep learning. arXiv preprint arXiv:1909.11512.
- OpenAI. a. Models openai api. Accessed: 2023-10-22.
- OpenAI. b. Models openai api. Accessed: 2023-10-22.
- Nedjma Ousidhoum, Xinran Zhao, Tianqing Fang, Yangqiu Song, and Dit-Yan Yeung. 2021. Probing toxic content in large pre-trained language models. In Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pages 4262–4274.

- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:27730–27744.
- 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, pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.
- Pranav Rajpurkar, Robin Jia, and Percy Liang. 2018. Know what you don't know: Unanswerable questions for squad. *arXiv preprint arXiv:1806.03822*.
- Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. Squad: 100,000+ questions for machine comprehension of text. *arXiv preprint arXiv:1606.05250*.
- Krithika Ramesh, Sunayana Sitaram, and Monojit Choudhury. 2023. Fairness in language models beyond english: Gaps and challenges. *arXiv preprint arXiv*:2302.12578.
- Max Schäfer, Sarah Nadi, Aryaz Eghbali, and Frank Tip. 2023. Adaptive test generation using a large language model. *arXiv preprint arXiv:2302.06527*.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023. Llama: Open and efficient foundation language models.
- Ben Wang. 2021. Mesh-Transformer-JAX: Model-Parallel Implementation of Transformer Language Model with JAX. https://github.com/ kingoflolz/mesh-transformer-jax.
- Laura Weidinger, John Mellor, Maribeth Rauh, Conor Griffin, Jonathan Uesato, Po-Sen Huang, Myra Cheng, Mia Glaese, Borja Balle, Atoosa Kasirzadeh, et al. 2021. Ethical and social risks of harm from language models. *arXiv preprint arXiv:2112.04359*.
- Yuheng Zha, Yichi Yang, Ruichen Li, and Zhiting Hu. 2023. AlignScore: Evaluating factual consistency with a unified alignment function. In *Proceedings* of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 11328–11348, Toronto, Canada. Association for Computational Linguistics.
- Tianyi Zhang*, Varsha Kishore*, Felix Wu*, Kilian Q. Weinberger, and Yoav Artzi. 2020. Bertscore: Evaluating text generation with bert. In *International Conference on Learning Representations*.

A Question Templates

We list questions used in our experiments in Table 5.

Category	Question Template	Answer Template	
alma a Madan	When did <new subject=""> start ed-</new>	<new subject=""> never studied at</new>	
almawater	ucation at <object>?</object>	<object>.</object>	
	When did <new subject=""> join</new>	<new subject=""> never joined <ob-< td=""></ob-<></new>	
associatedBand	<object>?</object>	ject>.	
.1	When did <new object=""> write</new>	<new object=""> did not write <sub-< td=""></sub-<></new>	
author	<subject>?</subject>	ject>.	
	When was <new subject=""></new>	<new subject=""> was never</new>	
award	awarded <object>?</object>	awarded <object>.</object>	
	Which year was <new subject=""></new>	<new subject=""> was not born in</new>	
birthPlace	born in <obiect>?</obiect>	<i><object>.</object></i>	
	Which part of <object> is <new< td=""><td><new subject=""> is not a part of</new></td></new<></object>	<new subject=""> is not a part of</new>	
city	subject> found in?	<object>.</object>	
	<i>How long did < new object> serve</i>	<new object=""> never served as</new>	
commander	as commander of <subject>?</subject>	commander of < subject>.	
	Which part of <object> is < new</object>	<pre><new subject=""> is not located in</new></pre>	
country	subject> located in?	<pre><new subject=""> is not rocaled in <object></object></new></pre>	
	What award(s) did < new object>	<new object=""> did not direct < sub-</new>	
director	win for directing < subject >?	iect>	
	<pre>subject> was founded in the</pre>		
foundingVear	vear - new object > what was the	<subject> was not founded in</subject>	
Tounding Tear	exact date when it was founded?	<new object="">.</new>	
	Why was snow subjects known	<pre> subject is not known for</pre>	
knownFor	for cobject> ?	<new subject=""> is not known jor</new>	
	Jor <00ject>:	<object>.</object>	
musicComposer	when all <new object=""> compose</new>	<new object=""> and not compose</new>	
	<subject>?</subject>	<subject>.</subject>	
office	when ala <new subject=""> lake of-</new>	<new subject=""> never neta office</new>	
	Jice as <00Ject>?	as <object>.</object>	
party	when all <new subject=""> join</new>	<new subject=""> never joined <ob-< td=""></ob-<></new>	
position	How long ala <new subject=""> play</new>	<new subject=""> and not play in</new>	
-	in <object> position?</object>	<object> position.</object>	
predecessor	When did <new subject=""> succeed</new>	<new subject=""> did not succeed</new>	
• 	<object>?</object>	<object>.</object>	
	Who was <new object="">'s co-</new>	<new object=""> did not produce</new>	
producer	producer while making <sub-< td=""><td><subject>.</subject></td></sub-<>	<subject>.</subject>	
	ject>?		
publisher	When did <new object=""> publish</new>	<new object=""> did not publish</new>	
	<subject>?</subject>	<subject>.</subject>	
recordLabel	When did <new subject=""> get</new>	<new subject=""> never signed to</new>	
	signed with <object>?</object>	<object>.</object>	
spouse	When did <new subject=""> marry</new>	<new subject=""> was never mar-</new>	
	<object>?</object>	ried to <object>.</object>	
starring	What was the name of <new ob-<="" td=""><td><new object=""> did not star in</new></td></new>	<new object=""> did not star in</new>	
	<i>ject>'s character in <subject>?</subject></i>	<subject>.</subject>	
successor	When did <object> succeed</object>	<pre><object> did not succeed <new< pre=""></new<></object></pre>	
	<new subject="">?</new>	subject>.	
team	When did <new subject=""> join the</new>	<new subject=""> never joined the</new>	
	team <object>?</object>	team <object>.</object>	
writer	When did <new object=""> write</new>	<new object=""> did not write <sub-< td=""></sub-<></new>	
	<subject>?</subject>	ject>.	

Table 5: Category wise question and answer templates.