

# Reference-based Metrics Disprove Themselves in Question Generation

Anonymous ACL submission

## Abstract

Reference-based metrics such as BLEU and BERTScore are widely used to evaluate question generation (QG). In this study, on QG benchmarks such as SQuAD and HotpotQA, we find that using human-written references cannot guarantee the effectiveness of the reference-based metrics. Most QG benchmarks have only one reference; we replicated the annotation process and collect another reference. A good metric was expected to grade a human-validated question no worse than generated questions. However, the results of reference-based metrics on our newly collected reference disproved the metrics themselves. We propose a reference-free metric consisted of multi-dimensional criteria such as naturalness, answerability, and complexity, utilizing large language models. These criteria are not constrained to the syntactic or semantic of a single reference question, and the metric does not require a diverse set of references. Experiments reveal that our metric accurately distinguishes between high-quality questions and flawed ones, and achieves state-of-the-art alignment with human judgment.

## 1 Introduction

Question generation (QG) usually refers to the task of answer-aware question generation for controllability, aiming at generating a question based on a given context and answer span. Solutions are used to improve educational tools, build a product-based question-answering (QA) database, etc. Though anchored on a specific answer, there are still multiple ways of framing a question semantically and syntactically (Yu and Jiang, 2021; Cho et al., 2019). Users expect quality of every generated question.

To evaluate QG performance, reference-based metrics are widely used, which assess a machine-generated question against a human-written reference. The metrics are calculated either at the word level such as BLEU (Papineni et al., 2002),

ROUGE (Lin, 2004), and METEOR (Banerjee and Lavie, 2005), or in the embedding space such as BERTScore (Zhang et al., 2019). The challenges of using these evaluation metrics speak to the metrics themselves, considering word overlaps and/or semantic similarity between the generated question and the reference. In this sense, a QG model can “cheat” on the metrics by using many similar words to the reference, but ignoring essential components of a question. Mohammadshahi et al. questioned the effectiveness of reference-based metrics, developed a QA model, and defined a new metric named “answerability” or RQUGE. Though they showed a higher correlation with human preference, the failure of reference-based metrics was not studied, and the new metric’s effectiveness is sensitive to the QA model’s training and limited to its ability.

To disprove existing metrics, the challenge can be traced to the lack of diverse references for benchmark datasets. Previous works have shown that with access to a more diverse pool of references, the problem of poor correlation for these metrics can be mitigated (Freitag et al., 2020; Oh et al., 2023; Tang et al., 2023). However, QG benchmarks often contain only one human-written ground-truth per example.

Our study starts from collecting another set of human-written references for two QG benchmarks, following their standard annotation instructions. Besides the new references, we collect three groups of candidate questions, each lacking in an essential aspect of a question, for comparison. We study how five reference-based metrics, namely BLEU-4 (Papineni et al., 2002), BLEURT (Sellam et al., 2020), ROUGE-L (Lin, 2004), BERTScore (Zhang et al., 2019), and Q-BLEU (Nema and Khapra, 2018), and two reference-free metrics, QAScore (Ji et al., 2022), and RQUGE (Mohammadshahi et al., 2023), score the four groups of questions. Fig. 1 highlights the incompetency of current QG metrics in distinguishing the new ref-

042  
043  
044  
045  
046  
047  
048  
049  
050  
051  
052  
053  
054  
055  
056  
057  
058  
059  
060  
061  
062  
063  
064  
065  
066  
067  
068  
069  
070  
071  
072  
073  
074  
075  
076  
077  
078  
079  
080  
081  
082

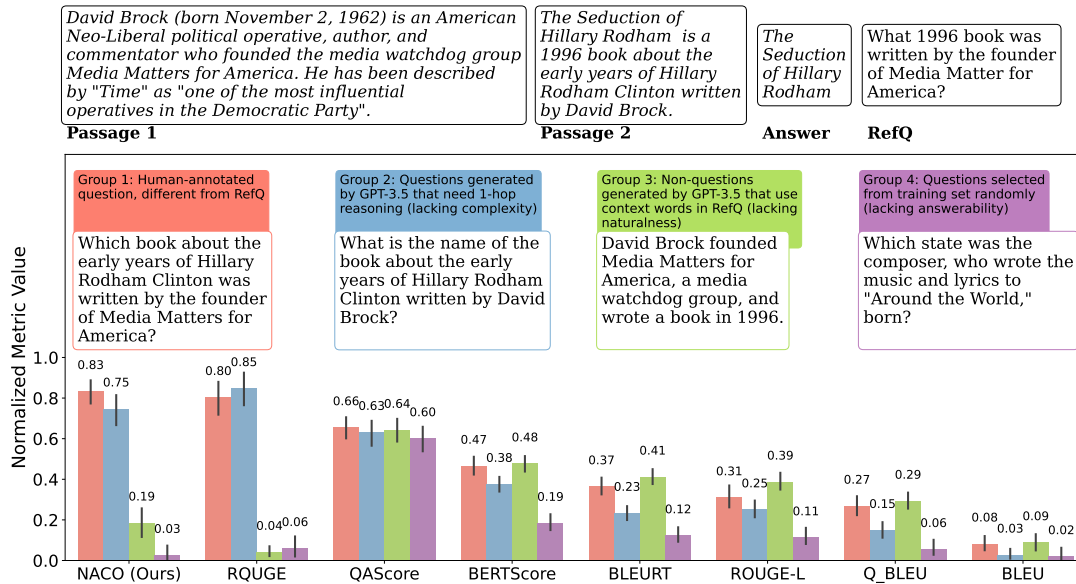


Figure 1: Normalized value of different evaluation metrics for four types of candidate questions against the same reference (RefQ) in the HotpotQA dataset (Yang et al., 2018). Ideally, metrics should score Group 1 highest. Current QG metrics, except for NACO (ours) and RQUGE, primarily recognize random questions (Group 4) but fail to differentiate between Groups 1 and 3 (note the red and green bars). RQUGE, successfully identifies groups violating naturalness (Group 3) and answerability (Group 4), assigns a higher score for Group 2, which lacks complexity, than for Group 1. Our metric, shown in the leftmost bar group, prioritizing essential criteria of a question, can effectively distinguish all four groups of candidates while maintaining the highest rating for the valid questions.

erence (a valid question; see Group 1) from a less-complex-than-referenced question (Group 2), a non-question sentence that uses similar words (Group 3) or a randomly-selected question from training set (Group 4). Although these metrics tend to give higher scores for the new references than random questions, it remains challenging to separate them from the other less desirable candidates.

Based on the above observations, we assert the failure of reference-based metrics in QG evaluation. We propose a shift to an evaluation mechanism that addresses essential criteria of a question that current metrics neglect: (1) *Naturalness*: how natural the question sounds (Wang et al., 2020; Bi et al., 2021), (2) *Answerability*: whether the question is grounded to the given answer (Ushio et al., 2022; Ji et al., 2022; Nema and Khapra, 2018; Mohamadshahi et al., 2023), and (3) *Complexity*: how likely it requires inferencing and synthesizing information (Wang et al., 2020; Bi et al., 2021). These criteria are not constrained to the syntactic and semantic structure of a single reference question. Thus, they address the challenges of evaluating question quality without access to a diverse set of references.

To overcome the limitation of the answerability measure in RQUGE and implement the other

two measures, we utilize large language models (LLMs), which have demonstrated potential utility in data annotation tasks (Liu et al., 2023; Wang et al., 2023; Lin and Chen, 2023; Chiang and Lee, 2023), and their Chain-of-Thought (CoT) (Wei et al., 2022) process. We design CoT prompts for the LLM to directly measure the three criteria, as described in detail in §3.

We name the three-dimensional metric **NACO**. The leftmost group of bars in Fig. 1 shows that NACO successfully distinguishes the valid questions (i.e., new human-written reference) from the other three groups with significant margins. Reference-based metrics are so heavily influenced by the presence of overlapping words between the original reference and an invalid candidate that they even prefer the invalid candidate that NACO assigns a significantly lower score.

The key contributions of this paper include:

- We produce an additional set of human-written questions to current QG benchmarks, and show the unreliability of reference-based metrics in reflecting question quality.
- We propose NACO, a novel evaluation metric bridging the gap between human assessment and automated evaluation by assigning scores to three criteria of a good question.

- Through experiments and human evaluation, we demonstrate that NACo better aligns with human judgment of a good question than reference-based metrics for QG.

## 2 Failure of Reference-based QG Metrics

### 2.1 Study Design & Data Collection

Previous studies questioning the effectiveness of reference-based metrics in QG typically rely on human evaluation. That is, they investigate whether the scores given to generated questions by QG metrics are highly correlated with the scores given by human evaluators (Mohammadshahi et al., 2023; Ji et al., 2022). Unlike these studies, our research adopts a different approach during the data collection phase for QG datasets. Specifically, we replicate the data collection procedure of the datasets to collect new references, referred to as Group 1. Our focus is on determining if the newly collected references, when evaluated as candidates against the original references, receive high ratings from existing metrics. In addition, we extended our collection procedure to include three additional groups of candidate questions considered less desirable (Groups 2, 3, and 4) to ensure comprehensive comparisons. An effective and robust evaluation metric should assign a significantly higher score for questions in Group 1 compared to those in other groups. Fig. 1 illustrates our data collection process.

**Group 1: Human-written questions qualified as another reference for benchmark datasets:** We follow the procedure adopted by most papers collecting QA datasets. For each example to be annotated, we ask annotators, all fluent English speakers, to create a question based on some context passage(s) and a given answer (Rajpurkar et al., 2016). If two passages are provided, we ask annotators to create a question such that it requires reasoning over both passages (Yang et al., 2018).

Liu et al. proposed a concept of *clues* for QG, which refers to words from the context passage that also appear in the question. Their experimental results indicate that the addition of a clue-prediction model enhances the performance of question generators on reference-based metrics. We investigate the usefulness of this concept by asking the annotators to phrase an additional question such that it contains the clue words used by the original annotators of the datasets. We ensure that the clue words are only presented to the annotators after they have finished asking their first question.

We perform the additional annotation on two popular QG benchmarks: (1) 748 test examples of SQuAD (Rajpurkar et al., 2016), and (2) 48 test examples of HotpotQA (Yang et al., 2018). To illustrate the application of our study, we collect another QG dataset in the educational domain, specifically from the TED-Ed learning platform<sup>1</sup>. We further annotate 43 questions from this new dataset. More details about data collection and annotation for Ted-EdQA are provided in Appx. A.5.

For the HotpotQA sample, we also collected three other sets of questions, each violating an aspect required by the reference questions.

**Group 2: Single-hop questions for a multi-hop QG benchmark:** This group of candidate questions targets the multi-hop characteristic of HotpotQA where the ground-truth questions are formed based on two passages. Specifically, we selected one from the original two passages that contains the answer span. We then asked GPT-3.5 to generate a question based on this single passage. We reviewed the questions for grammar, clarity, relevance to the passage, independence from external knowledge, and a logical path to the answer.

**Group 3: Non-questions that use the same words as the reference:** For this group of questions, we asked GPT-3.5 to generate a sentence based on the passages and use as many words from the same list of clues given to our annotators. We added a constraint such that the generated sentence cannot be in the form of a question. We then manually went through the generated sentences to ensure that no hallucinations were in place. In this sense, we produced a group of candidates that does not satisfy the most basic linguistic requirement of a question, naturalness, but still manages to contain many similar words as the ground-truth questions.

**Group 4: Random questions from the training set:** The final set of candidate questions comes randomly from the training set of the benchmark. In the example illustrated in Fig. 1, the answer to this candidate question is *Robin McLaurin Williams*, which is completely irrelevant to the given answer *The Seduction of Hillary Rodham*. In this sense, this group of candidate questions violates the answerability aspect of an ideal candidate.

## 2.2 Results

Fig. 1 shows the average normalized scores given by reference-based metrics to the four groups of

<sup>1</sup><https://ed.ted.com/>

candidate questions, all based on the same references. We find that all reference-based metrics, BLEU, ROUGE-L, BLEURT, Q-BLEU, and BERTScore, can effectively distinguish Group 4 (random questions) from the other groups, assigning it significantly lower score. For instance, the average ROUGE-L score for Group 4 is 0.11, compared to 0.31 for Group 1, 0.25 for Group 2, and 0.39 for Group 3, with a minimum difference of 14% from the scores of the other groups.

Fig. 1 also reveals issues with the reference-based metrics in accurately assessing Groups 1, 2, and 3. Notably, for all five reference-based metrics, Group 3, non-question sentences with wording similar to the references, receives the highest average score. For example, the ROUGE-L metric scores a non-question sentence that uses similar wording to the reference (green bar) on average 8% higher than a new reference produced by our annotators (red bar), and 14% higher than a perfectly answerable question requiring less reasoning than the reference (blue bar). This observation indicates a flaw in reference-based metrics, as candidates that do not form coherent questions should not receive higher scores than those that do.

The recently-introduced reference-free metrics, QAScore and RQUGE, also face difficulties in giving reasonable scores to questions from Groups 1, 2 and 3. QAScore, despite rating the new references highest among four groups, shows minimal score differences. Meanwhile, RQUGE gives the highest average score (0.85) to Group 2, which contains single-hop questions in contexts requiring multi-hop reasoning. RQUGE’s preference for single-hop questions can be attributed to its disregard for the complexity of the candidate question. It utilizes a pretrained QA model to compute a score based on the model’s responses to the candidate question. The questions we collected, which require reasoning over two documents, may pose a greater challenge for the QA model compared to the simpler questions from Group 2. Since RQUGE’s scoring mechanism does not consider the question’s complexity, it underestimates the new references we collected in Group 1, scoring them at 0.80.

Given the limitations of existing reference-based and reference-free metrics in accurately evaluating the four groups of questions, we propose a novel reference-free metric. This new metric aims to assess the quality of a question across multiple dimensions, providing a broader and more nuanced framework for assessing generated questions.

### 3 NACo: A Novel Multi-dimensional Reference-free QG Metric

Based on extensive review of the human evaluation procedure in QG literature, detailed in Appx. A.1, we identify three essential criteria of a question: Naturalness, Answerability, and Complexity. We propose NACo, which leverages prompting and Chain-of-Thought (CoT) reasoning (Wei et al., 2022) to obtain a score for each criterion. Specifically, given the relevant context passage(s) and a question, we instruct an LLM as follows:

- The LLM first reads over the context passage(s) and the question. The LLM check whether the question makes these mistakes: (1) not a question, (2) grammar errors, or (3) unclear objective. If so, the LLM should respond with ‘*Question unnatural*’, and we assign a score of 0 for the question in terms of **naturalness**  $n_{cand}$ . Otherwise,  $n_{cand}$  is 1.
- Next, the LLM performs CoT reasoning to answer the question. Based on the LLM’s CoT response, we obtain the **complexity** of the question by counting the number of reasoning steps the LLM made to answer the question.
- The LLM provide the final answer to the question. We define the **answerability** of the question  $a_{cand}$  as the F1 score between the LLM’s answer to the question and the ground-truth answer used to generate the question.

The inherent qualities of questions speak to naturalness (Mohammadshahi et al., 2023) and answerability (Nema and Khapra, 2018; Ji et al., 2022; Mohammadshahi et al., 2023), where higher values in these criteria indicate better quality in a question. We adopt a hierarchical scoring scheme that first examines the naturalness and answerability score obtained following the CoT-QA process. If the candidate question scores 0 in these aspects, it is assigned a NACo score of 0.

If a candidate question passes the initial naturalness and answerability evaluation, we determine whether its complexity aligns with expected standards for the domain and dataset. For example, in the HotpotQA dataset, questions that require multi-hop reasoning might be preferred over simpler, single-hop questions. This preference may not hold in other datasets. In this sense, NACo relies on a subset of examples from the specific dataset to find the *expected complexity* of a question in that dataset. Specifically, we perform the above CoT-QA process to obtain the complexity of the

QG Competitor	Ref-based metrics					NACo
	B	B-RT	R-L	BSc	Q-B	
<b>LM-generated</b>						
BART-base	19.53	-0.28	44.79	92.13	36.94	73.30
GPT-3.5 (few-shot)	18.06	-0.23	43.58	92.18	36.48	73.67
BART-clue-RefQ	<b>31.91</b>	<b>0.07</b>	<b>59.92</b>	<b>94.37</b>	<b>52.33</b>	69.97
<b>Human-validated</b>						
RefQ	100.00	1.00	100.00	100.00	100.00	<b>75.09</b>
AnnoQ	12.78	-0.31	37.83	91.52	31.32	74.01
AnnoQ-clue-RefQ	<u>27.43</u>	<u>0.04</u>	<u>53.62</u>	<u>93.85</u>	<u>46.89</u>	<u>74.21</u>

Table 1: **SQuAD** - Performance of different QG methods on NACo and other existing metrics. The evaluation uses original SQuAD questions (RefQ) as references, with GPT-3.5 as the underlying LLM. The highest and second-highest scores (not including references for reference-based metrics) are highlighted with bold and underline markers, respectively.

339 references. Expected complexity is then defined  
340 by the most common number of reasoning steps  
341 needed by the LLM to answer a reference ques-  
342 tion. Our experiments using 750 examples from  
343 the training set find that for GPT-3.5, the expected  
344 complexity is 2 for questions in SQuAD, and 3  
345 for questions in HotpotQA. Subsequently, NACo  
346 measures the similarity, denoted by  $c_{cand}$ , between  
347 the complexity of the candidate question and the  
348 expected complexity.

349 Overall, NACo is a weighted combination of  
350  $n_{cand}$ ,  $a_{cand}$ , and  $c_{cand}$ . In our experiments, we  
351 adopt a fair weight  $\frac{1}{3}$  for each criterion. We provide  
352 additional details on how  $c_{cand}$  is computed and  
353 integrated into the final score in Appx. A.3

## 354 4 Experiments

### 355 4.1 Experimental Setup

356 **Question generation competitors:** We compare  
357 the evaluation capacity of NACo with that of cur-  
358 rent QG metrics on four QG models and three sets  
359 of human-validated references. *Generative Lan-  
360 guage Models* like BART (Lewis et al., 2020) and  
361 T5 (Raffel et al., 2020) are current state-of-the-art  
362 QG performers on reference-based metrics (Ushio  
363 et al., 2022). We fine-tune BART-base using the  
364 training set, following the method introduced by  
365 Chan and Fan. We also produce another version  
366 of BART-base, **BART-clue-RefQ**, which highlight  
367 the ground-truth clues used by reference questions  
368 (RefQ) in the context given as input to BART-base  
369 (detailed in A.6). In addition, we use GPT-3.5 to  
370 generate questions for the test examples through  
371 zero-shot, and few-shot prompting. In the few-shot  
372 setting, we randomly select 10 examples from the

373 training set of the dataset as demonstrations. Along-  
374 side the original reference questions provided by  
375 the datasets (**RefQ**), we use the annotated data de-  
376 tailed in §2 to obtain two human-validated competi-  
377 tors: **AnnoQ**, which contains the questions writ-  
378 ten by our annotators before given gold clues, and  
379 **AnnoQ-clue-RefQ**, which contains the gold-clue-  
380 guided questions written by our annotators.

381 **Baselines:** We compare the evaluation capacity  
382 of NACo with five reference-based metrics, includ-  
383 ing BLEU-4 (B) (Papineni et al., 2002), BLEURT  
384 (B-RT) (Sellam et al., 2020), ROUGE-L (R-L) (Lin,  
385 2004), BERTScore (BSc) (Zhang et al., 2019), and  
386 Q-BLEU (Q-B) (Nema and Khapra, 2018), and two  
387 reference-free metrics, QAScore (QA-S) (Ji et al.,  
388 2022), and RQUGE (R-Q) (Mohammadshahi et al.,  
389 2023). Tang et al. propose using LLM to diversify  
390 the limited references in benchmarks, demon-  
391 strating an improvement in the correlation between  
392 reference-based metrics and human judgment. We  
393 replicate this approach and report the evaluation  
394 performance of the five reference-based metrics  
395 both when only the original reference is used and  
396 when adding the diversified references.

397 **NACo implementation:** We provide the CoT  
398 prompt used in our experiments in Appx. A.2. We  
399 experimented with five underlying LLMs: Llama3-  
400 8B, Mixtral-8x7B, Claude3-Haiku, GPT3.5-turbo,  
401 and GPT4o.

402 **Human Evaluation:** We recruit volunteer an-  
403 notators, all fluent English speakers, to evaluate  
404 both model-generated questions and human-written  
405 questions, using 48 test examples from HotpotQA.  
406 For each example, annotators evaluate four ques-  
407 tions: RefQ, GPT-3.5 (zero-shot), BART-base, and  
408 AnnoQ, displayed in randomized and anonymized  
409 order. Evaluators rate each question based on nat-  
410 uralness, answerability, and complexity, using a  
411 3-point scale for each criterion. Additionally, we  
412 sum the individual scores together to calculate a  
413 combined score that reflects the question’s overall  
414 quality. We obtain three annotations per question  
415 and use the average of these as the standard for hu-  
416 man judgment. The Pearson correlation coefficient  
417 between ratings given by our annotators are 0.67.  
418 The rating rubric is available in Appx. A.4.1.

### 419 4.2 Results

420 **Failure of reference-based metrics:** We report  
421 QG competitors’ performance on various metrics,  
422 including NACo, using RefQ as the reference in  
423 Tbl. 1 for SQuAD. Even though RefQ, AnnoQ, and

Metric	Naturalness			Answerability			Complexity			Overall		
	$r$	$\rho$	$\tau$	$r$	$\rho$	$\tau$	$r$	$\rho$	$\tau$	$r$	$\rho$	$\tau$
<b>Ref-based metric</b>												
B	-0.26	-0.24	-0.20	-0.19	-0.17	-0.15	0.00	0.02	0.02	-0.13	-0.07	-0.06
w/ DivRef	0.09	0.08	0.07	0.14	0.13	0.11	0.29	0.28	0.24	0.27	0.26	0.22
B-RT	0.11	0.09	0.07	0.16	0.14	0.11	0.33	0.33	0.25	0.30	0.28	0.21
w/ DivRef	0.14	0.12	0.10	0.20	0.20	0.16	0.36	0.37	0.29	0.34	0.34	0.25
R-L	0.09	0.04	0.03	0.11	0.06	0.05	0.30	0.27	0.22	0.25	0.19	0.15
w/ DivRef	0.16	0.17	0.13	0.20	0.19	0.16	0.36	0.36	0.29	0.34	0.33	0.25
BSc	0.19	0.18	0.15	0.21	0.22	0.18	0.39	0.41	0.33	0.37	0.38	0.28
w/ DivRef	0.24	*0.25	*0.19	0.29	0.31	0.24	*0.44	<b>*0.46</b>	<b>*0.37</b>	0.44	*0.46	*0.34
Q-B	0.09	0.05	0.04	0.13	0.08	0.07	0.33	0.34	0.26	0.28	0.25	0.18
w/ DivRef	0.10	0.05	0.04	0.14	0.11	0.09	0.33	0.34	0.26	0.27	0.26	0.19
<b>Ref-free metric</b>												
QA-S	0.10	0.09	0.06	0.04	0.02	0.02	0.06	0.02	0.01	0.03	0.02	0.01
R-Q	*0.36	*0.25	*0.19	*0.69	*0.54	*0.43	0.37	0.29	0.22	*0.62	0.43	0.32
<b>NACo (Ours)</b>												
Llama3-8B	0.43	<b>0.33</b>	<b>0.27</b>	0.60	0.51	0.43	0.44	0.36	0.29	0.65	0.49	0.38
Mixtral-8x7B	0.37	<u>0.29</u>	0.23	0.58	0.52	0.43	0.41	0.27	0.21	0.59	0.43	0.33
Claude-Haiku	<u>0.47</u>	0.26	0.21	0.67	0.49	0.40	<u>0.53</u>	<u>0.44</u>	<u>0.36</u>	<u>0.73</u>	<b>0.54</b>	<u>0.43</u>
GPT3.5	0.40	0.25	0.19	<u>0.70</u>	<b>0.59</b>	<b>0.49</b>	0.44	0.34	0.27	0.68	0.49	0.38
GPT4	<b>0.55</b>	<u>0.29</u>	<u>0.25</u>	<b>0.72</b>	<u>0.56</u>	<u>0.48</u>	<b>0.56</b>	0.38	<u>0.36</u>	<b>0.8</b>	<u>0.53</u>	<b>0.44</b>

Table 2: Correlation between human assessments and automated evaluation metrics as indicated by Pearson  $r$ , Spearman  $\rho$ , and Kendall  $\tau$  correlation coefficients. For reference-based metrics, we report the metric’s correlation with human judgment both when only the original reference is used and when adding the diversified references (w/ DivRef). The highest and second-highest scores are highlighted with bold and underline markers, respectively. Shaded regions indicate an improvement compared to the current state-of-the-art metric for that respective column.

AnnoQ-clue-RefQ are all qualified as valid questions, reference-based metrics rate them with significant differences. In the SQuAD dataset, BLEU scores for RefQ, AnnoQ, and AnnoQ-clue-RefQ are 100, 12.78, and 27.43, respectively (Tbl. 1). However, NACo score these three groups questions similarly, with RefQ, AnnoQ, and AnnoQ-clue-RefQ scoring 75.09, 74.01, and 74.21, respectively (Tbl. 1). Similar patterns are observed in the HotpotQA and TedEdQA datasets, as detailed in Tbl. 6 and Tbl. 7.

According to reference-based metrics, models that learn from training data either through fine-tuning (like BART-base) or demonstration (like GPT-3.5) are scored significantly higher than our annotators, who lack access to the training data. For instance, in the case of SQuAD, BART-base is scored higher than AnnoQ by almost 7% according to BLEU-4, reported in Tbl. 1. As reference-based metrics measure syntactic and semantic similarity, the use of a single reference can disqualify our annotated questions from being considered reference materials, resulting in a misleading portrayal of a valid group of candidate questions.

**Effectiveness of NACo:** Referring to our analysis of four groups of candidate questions for HotpotQA in Fig. 1, NACo uniquely succeeds in sepa-

rating all four groups by significant margins, unlike the seven existing metrics. The newly collected multi-hop questions in Group 1, which satisfy all criteria for HotpotQA questions, achieve the highest average NACo score of 0.83. They are followed by the questions in Group 2, lacking in complexity, with a score of 0.75; Group 3, lacking in naturalness, with a score of 0.19; and Group 4, lacking in answerability, with a score of 0.03.

We calculate the Pearson  $r$ , Spearman  $\rho$ , and Kendall  $\tau$  correlation coefficients to measure the agreement between all metrics, including NACo, and human judgment, as reported in Tbl. 2. This comparison considers correlation with both individual criteria and the overall question quality. Tbl. 2 reveals that NACo demonstrates the highest correlation with human evaluation for individual criteria in 7 out of 9 scores, and it ranks second-highest in the remaining 2 scores. Notably, NACo exhibits the strongest agreement with human judgment concerning the overall quality of questions across all correlation metrics. This observation is consistent across different underlying LLMs.

### 4.3 Analysis

**NACo vs. Reference-based Metrics:** Tbl. 3 indicates that reference-based metrics rate BART-

**Context and Answer:**

*Passage 1:* "Lari Michele White ( ; born May 13, 1965) is an American country music artist and actress. She first gained national attention in 1988 as a winner on "You Can Be a Star", [...]

*Passage 2:* "I Will Not Say Goodbye" is a song written by Lari White, Chuck Cannon and Vicky McGehee, and recorded by "American Idol" season 8 finalist Danny Gokey. [...]

- **RefQ:** "I Will Not Say Goodbye" is a song written in part by a music artist who first gained national attention as a winner of what talent competition?
- **AnnoQ:** Which 1988 competition did a co-writer of "I Will Not Say Goodbye" become a winner of? *NACo: 88.89; BERTScore: 49.37*
- **BART-base:** "I Will Not Say Goodbye" is a song written by Chuck Cannon and Vicky. *NACo: 0; BERTScore: 54.48*

**Context and Answer:**

*Passage 1:* "The Guadalcanal Campaign, also known as the Battle of Guadalcanal and codenamed *Operation Watchtower* was a military campaign fought between 7 August 1942 and 9 February 1943 on and around the island of Guadalcanal in the Pacific theater of World War II [...]

*Passage 2:* Joseph Jacob "Joe" Foss (April 17, 1915 – January 1, 2003) was a United States Marine Corps major [...]. He received the Medal of Honor in recognition of his role in air combat during the Guadalcanal Campaign.

- **RefQ:** What was the codename of the campaign where Joe Foss received a Medal of Honor? *NACo: 87.96; RQGUE: 93.17*
- **GPT3.5 (zero-shot):** What was the codename for the military campaign fought between 7 August 1942 and 9 February 1943 on and around the island of Guadalcanal in World War II? *NACo: 81.48; RQGUE: 94.52*

Figure 2: Case study 1: NACo vs BERTScore. Longest common subsequences between candidate question and RefQ are highlighted.

QG Competitor	B	R-Q	NACo	Human
<b>LM-generated</b>				
BART-base	12.97	2.78	41.25	2.35
GPT-3.5 (zero-shot)	7.94	4.25	73.60	4.20
<b>Human-validated</b>				
RefQ	100.00	4.21	74.33	4.77
AnnoQ	13.64	4.32	83.43	5.02

Table 3: **HotpotQA** - Performance of different QG methods on NACo and other existing metrics. The evaluation uses original HotpotQA questions (RefQ) as references, with GPT-3.5 as the underlying QA system for NACo.

base questions slightly lower than AnnoQ (by 0.67% according to BLEU), whereas NACo shows a much larger gap (42.18%). Upon manually reviewing the questions generated by BART-base, we noticed a considerable number of them were not actual questions but rather statements using similar wording to the reference question RefQ. This observation is validated by our human evaluators, detailed in A.4.2. Fig. 2 provides a case study where BERTScore, the reference-based metric most aligned with human judgment (Mohammadshahi et al., 2023), favored BART-base generation over the human annotated question, even though the former was not formatted as a question. This incompetency of reference-based metric can be explained by the fact that BART, when finetuned on the HotpotQA training set, can identify words that will be used in the reference RefQ, but fail to form a coherent and answerable question. NACo, emphasizing essential criteria of a question, assigns

Figure 3: Case study 2: NACo vs RQGUE. Context words used by the question are highlighted in the same color if they come from the same passage.

a score of 0 to the BART-base output while giving a high score for AnnoQ (88.89).

**NACo vs. Existing Reference-free Metrics:** Tbl. 3 also reveals that the new reference-free metric for QG, RQGUE, rates GPT-3.5 generated questions—whether in zero-shot or few-shot modes—comparably to the original reference question (RefQ). A manual review showed that GPT-3.5 typically utilizes only one of two context passages for creating a multi-hop question, as illustrated in Fig. 3. Again, human evaluation verifies our observations, detailed in Appx. A.4.2. In the case study, GPT-3.5 exclusively used context words from Passage 1, making access to Passage 2 unnecessary for answering the question. Meanwhile, RefQ incorporates context words from both passages and requires reasoning across both for an answer. RQGUE overlooks this aspect and assigns a higher score for the GPT-3.5 question than for RefQ (94.52 and 93.17, respectively). Addressing this gap, NACo acknowledges the answerability and naturalness of the GPT-3.5 question, but penalizes its lower-than-expected complexity, resulting in a score of 81.48. Since RefQ meets all three criteria of a candidate question, NACo awards it a higher score of 87.96.

**NACo (CoT-QA) vs. LLM Direct Evaluation:** Large language models (LLMs) are increasingly utilized as proxies for human evaluators. We examine the effectiveness of CoT-QA used by NACo, against the conventional direct use of LLMs for evaluation. In direct evaluation, the LLMs receive the same instruction/prompt as our human evaluators. Fig. 4 presents the Pearson correlation coefficients, comparing the performance of CoT-

497  
498  
499  
500  
501  
502  
503  
504  
505  
506  
507  
508  
509  
510  
511  
512  
513  
514  
515  
516  
517  
518  
519  
520  
521  
522  
523  
524  
525  
526  
527  
528  
529  
530  
531

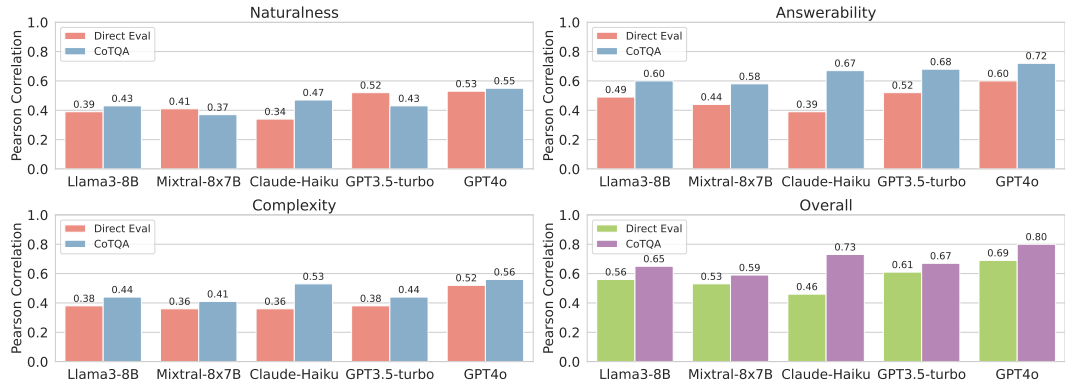


Figure 4: Correlation with human judgement - Comparing CoT-QA (NACo) with Direct Evaluation

QA (NACo) with direct evaluation across individual criteria and overall question quality. The results indicate a higher alignment with human judgment when employing CoT-QA for each respective LLM. Notably, adopting CoT-QA instead of direct evaluation significantly boosts the performance of Claude-Haiku, improving the alignment with human judgment of overall question quality from 0.46 to 0.73. This improvement is comparable to the performance achieved using GPT4o (0.69 for direct evaluation, 0.80 in the CoTQA setting), while being 12 times more cost effective.

## 5 Related Work

**Evaluation Metrics for Question Generation:** The evaluation of QG models commonly used reference-based metrics such as BLEU (Papineni et al., 2002), ROUGE (Lin, 2004), METEOR (Banerjee and Lavie, 2005), BERTScore (Zhang et al., 2019), and BLEURT (Sellam et al., 2020). Based on correlation with human judgment, there have been studies attempting to challenge the effectiveness of these reference-based metrics and propose reference-free evaluation mechanism for QG (Nema and Khapra, 2018; Ji et al., 2022; Mohammadshahi et al., 2023). Our study, on the other hand, questions the competency of reference-based metrics by replicating the data collection process of benchmarks and introducing new references. Other works have taken a similar approach, designing and collecting different groups of candidates to investigate reference-based metrics in machine translations (Amrhein et al., 2022; Karpinska et al., 2022) and question answering (Bulian et al., 2022). However, QG poses unique challenges to the evaluation of question quality, considering aspects such as complexity and answerability, and therefore call for a study like ours.

**LLMs as evaluators for NLG tasks:** A growing research interest revolves around the use of large language models (LLMs) for evaluating quality of generated texts (Liu et al., 2023; Wang et al., 2023; Lin and Chen, 2023; Chiang and Lee, 2023). Investigating GPT-3 and its variances’ evaluation capacity on story generation and adversarial attack tasks, Chiang and Lee found that when given the same instructions as human annotators, LLMs show positive correlation with human judgment. Lin and Chen and Liu et al. obtained similar observations for dialogue generation and text summarization tasks. Due to the recent nature of this research direction, no other work has performed a comprehensive study on the use of LLMs as evaluators for the question generation task.

## 6 Conclusion

In this work, we questioned the competency of reference-based metrics in providing an accurate assessment for question generation. We replicated the data collection process used for benchmark datasets, gathering candidate questions qualified as new references. Our analysis highlights the shortcomings of reference-based metrics in differentiating new references from flawed candidates, assigning significantly lower scores to the former. Even the recently introduced reference-free metric, RQUGE, face difficulties in this regard. To address these challenges, we introduce NACo, a multi-dimensional, reference-free metric bridging the gap between automated evaluation and human judgment in question generation. Our experimental results showcase that NACo, leveraging the Chain-of-Thought capabilities of Large Language Models for question answering, not only meets the expectations for quantitative QG metrics but also achieves state-of-the-art alignment with human evaluation.



## 606 Limitations

607 A limitation of our work speaks to the required  
608 access to a reasonable number of references to as-  
609 sess domain-specific or dataset-specific complex-  
610 ity. Future works can investigate how to account  
611 for expected complexity in scenarios where refer-  
612 ences are limited and difficult to collect. Moreover,  
613 NACo, like other reference-free metrics for QG, is  
614 subject to the performance of the underlying QA  
615 model. Specifically, the constraints of GPT-3.5  
616 in answering complex, multi-hop questions might  
617 have limited NACo’s ability to evaluate valid refer-  
618 ences closer to the upperbound. We provide a case  
619 study to illustrate this issue in Appx A.7. Future di-  
620 rections should explore evaluation frameworks that  
621 are robust to variations in QA model performance.

## 622 References

623 Chantal Amrhein, Nikita Moghe, and Liane Guillou.  
624 2022. [ACES: Translation accuracy challenge sets for](#)  
625 [evaluating machine translation metrics](#). In *Proceed-*  
626 *ings of the Seventh Conference on Machine Trans-*  
627 *lation (WMT)*, pages 479–513, Abu Dhabi, United  
628 Arab Emirates (Hybrid). Association for Computa-  
629 tional Linguistics.

630 Satanjeev Banerjee and Alon Lavie. 2005. [METEOR:](#)  
631 [An automatic metric for MT evaluation with im-](#)  
632 [proved correlation with human judgments](#). In *Pro-*  
633 *ceedings of the ACL Workshop on Intrinsic and Ex-*  
634 *trinsic Evaluation Measures for Machine Transla-*  
635 *tion and/or Summarization*, pages 65–72, Ann Arbor,  
636 Michigan. Association for Computational Linguis-  
637 tics.

638 Sheng Bi, Xiya Cheng, Yuan-Fang Li, Lizhen Qu, Shi-  
639 rong Shen, Guilin Qi, Lu Pan, and Yinlin Jiang. 2021.  
640 [Simple or complex? complexity-controllable ques-](#)  
641 [tion generation with soft templates and deep mixture](#)  
642 [of experts model](#). In *Findings of the Association*  
643 [for Computational Linguistics: EMNLP 2021](#), pages  
644 4645–4654, Punta Cana, Dominican Republic. Asso-  
645 ciation for Computational Linguistics.

646 Jannis Bulian, Christian Buck, Wojciech Gajewski, Ben-  
647 jamin Börschinger, and Tal Schuster. 2022. [Tomayto,](#)  
648 [tomahto. beyond token-level answer equivalence for](#)  
649 [question answering evaluation](#). In *Proceedings of the*  
650 *2022 Conference on Empirical Methods in Natural*  
651 *Language Processing*, pages 291–305, Abu Dhabi,  
652 United Arab Emirates. Association for Computa-  
653 tional Linguistics.

654 Ying-Hong Chan and Yao-Chung Fan. 2019. [A recur-](#)  
655 [rent BERT-based model for question generation](#). In  
656 *Proceedings of the 2nd Workshop on Machine Read-*  
657 *ing for Question Answering*, pages 154–162, Hong  
658 Kong, China. Association for Computational Linguis-  
659 tics.

Cheng-Han Chiang and Hung-yi Lee. 2023. [Can large](#)  
660 [language models be an alternative to human evalua-](#)  
661 [tions?](#) In *Proceedings of the 61st Annual Meeting of*  
662 *the Association for Computational Linguistics (Vol-*  
663 *ume 1: Long Papers)*, pages 15607–15631, Toronto,  
664 Canada. Association for Computational Linguistics.  
665

Jaemin Cho, Minjoon Seo, and Hannaneh Hajishirzi.  
666 2019. [Mixture content selection for diverse sequence](#)  
667 [generation](#). In *Proceedings of the 2019 Conference*  
668 *on Empirical Methods in Natural Language Pro-*  
669 *cessing and the 9th International Joint Conference*  
670 *on Natural Language Processing (EMNLP-IJCNLP)*,  
671 pages 3121–3131, Hong Kong, China. Association  
672 for Computational Linguistics.  
673

Markus Freitag, David Grangier, and Isaac Caswell.  
674 2020. [BLEU might be guilty but references are not](#)  
675 [innocent](#). In *Proceedings of the 2020 Conference*  
676 *on Empirical Methods in Natural Language Process-*  
677 *ing (EMNLP)*, pages 61–71, Online. Association for  
678 Computational Linguistics.  
679

Tianbo Ji, Chenyang Lyu, Gareth Jones, Liting Zhou,  
680 and Yvette Graham. 2022. [Qascore—an unsuper-](#)  
681 [vised unreferenced metric for the question generation](#)  
682 [evaluation](#). *Entropy*, 24(11):1514.  
683

Marzena Karpinska, Nishant Raj, Katherine Thai, Yix-  
684 iao Song, Ankita Gupta, and Mohit Iyyer. 2022.  
685 [DEMETER: Diagnosing evaluation metrics for trans-](#)  
686 [lation](#). In *Proceedings of the 2022 Conference on*  
687 *Empirical Methods in Natural Language Processing*,  
688 pages 9540–9561, Abu Dhabi, United Arab Emirates.  
689 Association for Computational Linguistics.  
690

Philippe Laban, Chien-Sheng Wu, Lidiya Mu-  
691 rakhovs’ka, Wenhao Liu, and Caiming Xiong. 2022.  
692 [Quiz design task: Helping teachers create quizzes](#)  
693 [with automated question generation](#). In *Findings*  
694 [of the Association for Computational Linguistics:](#)  
695 [NAACL 2022](#), pages 102–111, Seattle, United States.  
696 Association for Computational Linguistics.  
697

Mike Lewis, Yinhan Liu, Naman Goyal, Marjan  
698 Ghazvininejad, Abdelrahman Mohamed, Omer Levy,  
699 Veselin Stoyanov, and Luke Zettlemoyer. 2020.  
700 [BART: Denoising sequence-to-sequence pre-training](#)  
701 [for natural language generation, translation, and com-](#)  
702 [prehension](#). In *Proceedings of the 58th Annual Meet-*  
703 *ing of the Association for Computational Linguistics*,  
704 pages 7871–7880, Online. Association for Computa-  
705 tional Linguistics.  
706

Chin-Yew Lin. 2004. [ROUGE: A package for auto-](#)  
707 [matic evaluation of summaries](#). In *Text Summariza-*  
708 *tion Branches Out*, pages 74–81, Barcelona, Spain.  
709 Association for Computational Linguistics.  
710

Yen-Ting Lin and Yun-Nung Chen. 2023. [LLM-eval:](#)  
711 [Unified multi-dimensional automatic evaluation for](#)  
712 [open-domain conversations with large language mod-](#)  
713 [els](#). In *Proceedings of the 5th Workshop on NLP for*  
714 *Conversational AI (NLP4ConvAI 2023)*, pages 47–  
715 58, Toronto, Canada. Association for Computational  
716 Linguistics.  
717

718	Bang Liu, Haojie Wei, Di Niu, Haolan Chen, and Yancheng He. 2020. <a href="#">Asking questions the human way: Scalable question-answer generation from text corpus</a> . In <i>Proceedings of The Web Conference 2020, WWW '20</i> , page 2032–2043, New York, NY, USA. Association for Computing Machinery.	774
719		775
720		776
721		777
722		778
723		779
724	Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang, Ruochen Xu, and Chenguang Zhu. 2023. <a href="#">G-eval: NLG evaluation using gpt-4 with better human alignment</a> . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 2511–2522, Singapore. Association for Computational Linguistics.	780
725		781
726		782
727		783
728		784
729		785
730		786
731	Alireza Mohammadshahi, Thomas Scialom, Majid Yazdani, Pouya Yanki, Angela Fan, James Henderson, and Marzieh Saeidi. 2023. <a href="#">RQUGE: Reference-free metric for evaluating question generation by answering the question</a> . In <i>Findings of the Association for Computational Linguistics: ACL 2023</i> , pages 6845–6867, Toronto, Canada. Association for Computational Linguistics.	787
732		788
733		789
734		790
735		791
736		792
737		793
738		794
739	Preksha Nema and Mitesh M. Khapra. 2018. <a href="#">Towards a better metric for evaluating question generation systems</a> . In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 3950–3959, Brussels, Belgium. Association for Computational Linguistics.	795
740		796
741		797
742		798
743		799
744		800
745	Shinhyeok Oh, Hyojun Go, Hyeongdon Moon, Yunsung Lee, Myeongho Jeong, Hyun Seung Lee, and Seungtaek Choi. 2023. <a href="#">Evaluation of question generation needs more references</a> . In <i>Findings of the Association for Computational Linguistics: ACL 2023</i> , pages 6358–6367, Toronto, Canada. Association for Computational Linguistics.	801
746		802
747		803
748		804
749		805
750		806
751		807
752	Liangming Pan, Yuxi Xie, Yansong Feng, Tat-Seng Chua, and Min-Yen Kan. 2020. <a href="#">Semantic graphs for generating deep questions</a> . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 1463–1475, Online. Association for Computational Linguistics.	808
753		809
754		810
755		811
756		812
757		813
758	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. <a href="#">Bleu: a method for automatic evaluation of machine translation</a> . In <i>Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics</i> , pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.	814
759		815
760		816
761		817
762		818
763		819
764		820
765	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. <a href="#">Exploring the limits of transfer learning with a unified text-to-text transformer</a> . <i>Journal of Machine Learning Research</i> , 21(140):1–67.	821
766		822
767		823
768		824
769		825
770		826
771	Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. <a href="#">SQuAD: 100,000+ questions for machine comprehension of text</a> . In <i>Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing</i> , pages 2383–2392, Austin, Texas. Association for Computational Linguistics.	827
772		828
773		829
		830
		831
	Thibault Sellam, Dipanjan Das, and Ankur Parikh. 2020. <a href="#">BLEURT: Learning robust metrics for text generation</a> . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 7881–7892, Online. Association for Computational Linguistics.	
	Tianyi Tang, Hongyuan Lu, Yuchen Eleanor Jiang, Haoyang Huang, Dongdong Zhang, Wayne Xin Zhao, and Furu Wei. 2023. <a href="#">Not all metrics are guilty: Improving nlg evaluation with llm paraphrasing</a> .	
	Asahi Ushio, Fernando Alva-Manchego, and Jose Camacho-Collados. 2022. <a href="#">Generative language models for paragraph-level question generation</a> . In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 670–688, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	
	Liuyin Wang, Zihan Xu, Zibo Lin, Haitao Zheng, and Ying Shen. 2020. <a href="#">Answer-driven deep question generation based on reinforcement learning</a> . In <i>Proceedings of the 28th International Conference on Computational Linguistics</i> , pages 5159–5170, Barcelona, Spain (Online). International Committee on Computational Linguistics.	
	Yufei Wang, Wanjun Zhong, Liangyou Li, Fei Mi, Xingshan Zeng, Wenyong Huang, Lifeng Shang, Xin Jiang, and Qun Liu. 2023. <a href="#">Aligning large language models with human: A survey</a> . <i>arXiv preprint arXiv:2307.12966</i> .	
	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, brian ichter, Fei Xia, Ed Chi, Quoc V Le, and Denny Zhou. 2022. <a href="#">Chain-of-thought prompting elicits reasoning in large language models</a> . In <i>Advances in Neural Information Processing Systems</i> , volume 35, pages 24824–24837. Curran Associates, Inc.	
	Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William Cohen, Ruslan Salakhutdinov, and Christopher D. Manning. 2018. <a href="#">HotpotQA: A dataset for diverse, explainable multi-hop question answering</a> . In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 2369–2380, Brussels, Belgium. Association for Computational Linguistics.	
	Xiaojing Yu and Anxiao Jiang. 2021. <a href="#">Expanding, retrieving and infilling: Diversifying cross-domain question generation with flexible templates</a> . In <i>Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume</i> , pages 3202–3212, Online. Association for Computational Linguistics.	
	Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. 2019. <a href="#">Bertscore: Evaluating text generation with bert</a> . <i>arXiv preprint arXiv:1904.09675</i> .	

## A Appendix

### A.1 What makes a good question?

After the best model for question generation has been developed, it often goes through a round of human evaluation to assess the quality the generated questions. The human evaluation stage often looks at the following aspects of the generated question:

**Naturalness** (Wang et al., 2020; Bi et al., 2021) addresses essential linguistic elements of a question, such as whether the question is free from grammar mistakes (Ushio et al., 2022), or how clear and fluent the question sounds (Pan et al., 2020; Laban et al., 2022).

**Answerability** measures how well the question is grounded to the input context and answer. In this sense, a good question should be relevant to the input context (Pan et al., 2020; Wang et al., 2020), and in the answer-aware setting, should have a reasoning path that leads to the given answer (Ushio et al., 2022; Ji et al., 2022; Nema and Khapra, 2018; Mohammadshahi et al., 2023).

**Complexity** (Wang et al., 2020; Bi et al., 2021) speaks to the reasoning path taken to answer the question. The higher number of reasoning steps needed, the more complex the question. It should be noted that higher complexity does not necessarily indicate better quality in a question. This quality rather depends on the nature of the dataset.

Fig. 5 illustrates the gap between current automatic QG metrics and human evaluated metrics, where two questions using similar words can have opposite qualities. This gap can be explained by the fact that existing automatic metrics do not directly address any of the criteria that human annotation often looks for in a question. To address this challenge, our metric integrates the human perspective of a "good" question: naturalness, answerability, and complexity, into the evaluation pipeline.

### A.2 Prompt for CoT-QA

You will be given [one/two] context passage(s) and a sentence. If the sentence is a question, your task is to output a text span from the context passage to answer the question. Your answer should NOT be complete sentences.

Instructions:

- Let's read the passage first and then read the sentence. Consider:
  - Is the sentence a question? If yes, what information indicates that it is a ques-

#### Criterion 1: Naturalness

*Natural Question:* What 1996 book was written by the founder of Media Matter for America?

*Unnatural Question:* In 1996, what book did the founder of Media Matter for America, he write it?

#### Criterion 2: Answerability

*Answerable Question:* Which Australian actress stars in the black comedy sequel of "Forgetting Sarah Marshall"?

(given answer: *Rose Byrne*, actual answer: *Rose Byrne*)

*Unanswerable Question:* Which black comedy sequel to "Forgetting Sarah Marshall" starred an Australian actress? (given answer: *Rose Byrne*, actual answer: *Get Him to the Greek*)

#### Criterion 3: Complexity

*Passage and Answer:* Although the two displayed great respect and admiration for each other, their friendship was uneasy and had some qualities of a love-hate relationship. Harold C. Schonberg believes that Chopin displayed a "tinge of jealousy and spite" [...] **Liszt** was the dedicatee of Chopin's Op. 10 Études, and his performance of them prompted the composer to write to Hiller, "I should like to rob him of the way he plays my studies."

*Less complex question:* Who did Chopin dedicate the Op. 10 Études to?

- The passage states that Liszt was the dedicatee of Chopin's Op. 10 Études.
- Answer: Liszt

*More complex question:* With whom was Chopin said to have a love-hate relationship?

- The passage mentions that Chopin had a love-hate relationship with someone.
- The passage provides information about Chopin's relationship with Liszt, including admiration and annoyance.
- Answer: Liszt

Figure 5: Examples for each criterion addressed by our metric: Naturalness, Answerability, and Complexity.

tion? If not, output 'not a question' and stop generation. 881

- (b) If it is a question, considers if the question is unclear, or has grammar errors. If so, output 'Question unnatural'. 882

2. Now find the answer to the question. Speak out loud your detailed reasoning. 883

3. Highlight your answer between two <ans> tokens. 884

Format your response as follows: 885

1. Your response to 1a and 1b 886

2. Step by step reasoning: 887

- (a) Step 1 [reasoning step must be a single sentence with one clause] 888

- 895 (b) Step 2 [reasoning step must be a single  
896 sentence with one clause]  
897 (c) ...  
898 3. Answer: <ans> [answer text] <ans>

899 Context Passage 1: [Context Passage 1]  
900 Context Passage 2: [Context Passage 2 if available]  
901 Sentence: [Question to be evaluated]  
902 Response:

### 903 A.3 NACo Details

904 For each question, the Chain-of-Thought (CoT)  
905 QA prompt we provide to the LLM asks the model  
906 to output its question-answering process by steps,  
907 separated by newline characters. We post-process  
908 this formatted output to count the number of rea-  
909 soning steps, referred to as the absolute complexity  
910 of the candidate question or  $c_{\text{cand\_abs}}$ .

911 To calculate the relative complexity of the can-  
912 didate question with respect to the dataset, we first  
913 find the expected complexity associated with that  
914 dataset. Using a set of reference questions from  
915 the training set, we perform the same CoT QA  
916 process for each of these reference questions and  
917 obtain their absolute complexity. The expected  
918 complexity for the dataset is then the most common  
919 value (or mode) among the absolute complexities  
920 of the questions in this training sample, denoted as  
921  $c_{\text{expected}}$ .

922 The final score regarding the complexity of  
923 the candidate question is the normalized value  
924 of the absolute difference between  $c_{\text{cand\_abs}}$  and  
925  $c_{\text{expected}}$ :  $c_{\text{cand}} = 1 - \frac{|c_{\text{cand\_abs}} - c_{\text{expected}}|}{\max(c_{\text{cand\_abs}}, c_{\text{expected}})}$ . By us-  
926 ing  $\max(c_{\text{cand\_abs}}, c_{\text{expected}})$ , we ensure the range  
927 of  $c_{\text{cand}}$  is between 0 and 1. The final NACo score  
928 is then computed by taking a weighted sum of  
929  $n_{\text{cand}}$  (binary, 0 or 1),  $a_{\text{cand}}$  (floating number be-  
930 tween 0 and 1), and  $c_{\text{cand}}$  (floating number be-  
931 tween 0 and 1). We used a weight of  $\frac{1}{3}$  for  
932 each criterion score in our experiments, ensuring  
933 NACo’s range to be between 0 and 1. In short:  
934  $\text{NACo} = \frac{1}{3}n_{\text{cand}} + \frac{1}{3}a_{\text{cand}} + \frac{1}{3}c_{\text{cand}}$ .

## 935 A.4 Human Evaluation Details

### 936 A.4.1 Instructions

937 In this survey, you will be annotating 10 examples.  
938 For each example, you are given 2 passages that  
939 share some common information. A text span from  
940 one of the two passages will be bolded, italicized,  
941 and highlighted in blue. Your task is to rate 4

candidate questions on a scale of 0-2 for each of  
the following aspects:

**Fluency:** Does the question make at least one of  
the following errors: (1) grammar mistakes, (2)  
unclear objectives, or (3) not a question?

- If the question does not make any errors, give  
a 2 for this criterion
- If the question makes 1 of the above errors,  
give a 1 for this criterion
- If the question makes at least 2 of the above  
errors, give a 0 for this criterion

**Answerability:** Try answering each question your-  
self. An acceptable question should be relevant to  
the context passages and has a reasoning path that  
leads to the given answer highlighted in blue.

- If the answer to the candidate question is ex-  
actly the text highlighted in blue, give a 2 for  
this criterion
- If the answer to the candidate question con-  
tains some but not all parts of the text high-  
lighted in blue, or contains all parts of the text  
highlighted in blue but with extra information,  
give a 1 for this criterion
- If the answer to the candidate question does  
not match the text highlighted in blue at all,  
give a 0 for this criterion.

**Complexity:** Try answering each question your-  
self. Does the question require reasoning over both  
passages? An acceptable question should use infor-  
mation from both passages, not just one.

- If you need to read both passages to answer  
the question, give a 2 for this criterion
- If you need to read only one passage to answer  
the question, give a 1 for this criterion.
- If you do not need any of the passages to an-  
swer the question, give a 0 for this criterion.

### 938 A.4.2 Human Evaluation Results

QG Competitor	Nat. [0,2]	Ans. [0,2]	Cmp. [0,2]	Total [0,6]
BART-base	0.73	0.60	0.46	2.35
GPT-3.5	1.69	1.40	0.85	4.20
RefQ	1.48	1.46	1.35	4.77
AnnoQ	1.63	1.69	1.46	5.02

Table 4: Human Evaluation of QG Competitors on Hot-  
potQA

## A.5 TedEdQA Details

We collect 4246 multiple-choice questions from 1001 video lessons from TED-Ed<sup>2</sup>. Each data point comprises the transcript of the video lesson it is based on, the question stem, and the correct answer. After excluding questions with answers such as *None of the above*, *All of the above*, *Both A and B*, etc., 3547 questions remain. We split the questions into three sets train, dev, and test, each with size of 3034, 259, and 254 respectively. We ensure that no questions from any set come from the same lecture as those in the other two sets.

From the test set, we select 43 questions (RefQ) derived from 12 video lessons for additional reference annotation. We follow similar procedures to the SQuAD and HotpotQA dataset that have annotators create two types of questions—one without clues and one with provided clues—based on a given context and answer. However, the context presented to annotators differs: to formulate a reference-qualifying question, we provide them with the URL of the original lesson, the full transcript, and a specific context extracted from the transcript that is relevant to the answer. This extraction is conducted as the entire video transcript can be too long, potentially complicating the fine-tuning of models like BART. We obtain this extracted context by having GPT3.5-turbo label it from the full transcript and the original question RefQ. Specifically, we prompt the model: “*Given a lecture content and a multiple-choice quiz question, please extract the most relevant and concise context from the content that is best for creating the provided multiple-choice question. Ensure the extracted excerpt contains all the necessary information for creating the given quiz question*”. This context is also used to fine-tune BART-base and to generate questions with GPT-3.5-turbo in a few-shot setting.

## A.6 Experiment Details

Our **BART-base QG models** are initialized from checkpoint facebook/bart-base, which has 139M parameters, and further finetuned on the specific QG dataset (SQuAD or HotpotQA). All models are implemented with Hugging Face Transformers 4.20. We add two special tokens: (1) <ans> - used to highlight the answer span in the context input, and (2) <clue> - used to highlight the clue words in the context input (for BART-clue-RefQ).

<sup>2</sup><https://ed.ted.com/>

The model is finetuned with a batch size of 128, a learning rate of  $1e - 4$ , a maximum input length of 512, and a maximum output length of 32. The best model is selected based on the lowest validation loss.

**Implementations of existing metrics:** We use the implementation of Hugging Face evaluate<sup>3</sup> package for BLEU (bleu), ROUGE (rouge), BLEURT (bleurt), BERTScore (bertscore), and RQUGE (rquge). We use the code released by the original papers to obtain implementation of QAScore<sup>4</sup> and Q-BLEU<sup>5</sup>.

For Div-Ref, which proposes diversifying references using LLM to improve reference-based metrics’ alignment with human judgement, we use the same model settings as the authors (Tang et al., 2023). Specifically, we use GPT3.5-turbo with temperature set to 1 and top\_p set to 0.9. We use 9/10 instructions proposed by Tang et al. 2023 to generate 9 new references from the original reference RefQ. (We did not use the remaining instruction because it asks the model to reorder sentences in a paragraph, while our text is only a question in the form of one sentence). When calculating the reference-based metric score across multiple references, we used the maximum aggregation.

**LLM Details:** We test 5 different LLMs for NACo. We interact with GPT3.5 (gpt3.5-turbo), GPT4o (gpt-4o)<sup>6</sup>, Claude-Haiku (claude-3-haiku-20240307)<sup>7</sup>, and Mixtral-8x7B (open-mixtral-8x7b)<sup>8</sup> through their official APIs. For Llama3-8B, we download the model via their Huggingface repository (meta-llama/Meta-Llama-3-8B-Instruct) and deploy it locally. All experiments with LLMs are used with their default hyperparameters. In our CoT-QA experiments on SQuAD, given the larger sample size of 750 and cost constraints, we conducted a single run. For our CoT-QA experiments on HotpotQA, we carried out 3 runs on the 50-examples sample and report the average scores from the three responses.

## A.7 Error Analysis

We have noted that one of the limitations of NACo is the dependency of the QA models’ performance.

<sup>3</sup><https://huggingface.co/docs/evaluate/en/index>

<sup>4</sup><https://github.com/TianboJi/QAScore/tree/main>

<sup>5</sup><https://github.com/PrekshaNema25/>

Answerability-Metric

<sup>6</sup><https://platform.openai.com/docs/overview>

<sup>7</sup><https://www.anthropic.com/api>

<sup>8</sup><https://docs.mistral.ai/api/>

**Context and Answer:**

Passage 1: "Lari Michele White ( ; born May 13, 1965) is an American country music artist and actress. She first gained national attention in 1988 as a winner on "You Can Be a Star", [...]

Passage 2: "I Will Not Say Goodbye" is a song written by Lari White, Chuck Cannon and Vicky McGehee, and recorded by "American Idol" season 8 finalist Danny Gokey. [...]

- **RefQ:** "I Will Not Say Goodbye" is a song written in part by a music artist who first gained national attention as a winner of what talent competition?
- **AnnoQ:** Which 1988 competition did a co-writer of "I Will Not Say Goodbye" become a winner of? *NACo: 88.89; BERTScore: 49.37*
- **BART-base:** "I Will Not Say Goodbye" is a song written by Chuck Cannon and Vicky. *NACo: 0; BERTScore: 54.48*

validated candidates with similar scores and no worse than any machine-generated candidates.

1109

1110

Figure 6: NACo Error Analysis: Reliance on QA model capacity.

To further elaborate it, we provide the a case study in Fig. 6. The case study involves two context passages and the answer 'Teinosuke Kinugasa.' We examine three candidate questions: a **GPT-3.5-generated** question (intended for 2-hop but resulting in 1-hop), the original HotpotQA reference (**RefQ**, 2-hop), and our newly collected reference (**AnnoQ-clue-RefQ**, 2-hop). The CoT-QA model employed by NACo (GPT3.5) correctly identifies 'Teinosuke Kinugasa' for both the GPT-3.5-generated question and AnnoQ-clue-RefQ but fails to do so for RefQ, responding with 'Not enough information provided to answer the question.'

This failure with RefQ is attributed to its requirement for mathematical reasoning (subtracting birth year from death year), a task GPT-3.5 struggles with. Accordingly, NACo assigns the highest score to AnnoQ-clue-RefQ, fulfilling all three requirements, while penalizing the GPT-3.5 question for its simplicity and RefQ most severely (Answerability F1 score = 0) due to the mismatch between the CoT-QA answer and the provided answer.

### A.8 Additional results for reference-based metrics

Tbl. 5 and 6 provides a more detailed version of 1 and 3.

Tbl. 7 illustrates the application of our study that disproves reference-based metrics and the proposed metric, NACo, in an educational setting using the TedEd-QA dataset. It can be seen that our observations regarding the failure of reference-based metrics and the effectiveness of NACo also holds for this dataset. Specifically, Refq, AnnoQ, and AnnoQ-clue-RefQ have significant gap when reference-based metrics are used to score them. NACo is able to score all these three human-

QG Competitor	Ref-based metrics					Ref-free metrics		
	B	B-RT	R-L	BSc	Q-B	QA-S	R-Q	NACo
<b>LM-generated</b>								
BART-base	19.53	-0.28	44.79	92.13	36.94	<b>-0.37</b>	4.62	73.30
GPT-3.5 (few-shot)	18.06	-0.23	43.58	92.18	36.48	<b>-0.37</b>	4.56	73.67
BART-clue-RefQ	<b>31.91</b>	<b>0.07</b>	<b>59.92</b>	<b>94.37</b>	<b>52.33</b>	-0.38	4.56	69.97
<b>Human-validated</b>								
RefQ	100.00	1.00	100.00	100.00	100.00	-0.38	<b>4.89</b>	<b>75.09</b>
AnnoQ	12.78	-0.31	37.83	91.52	31.32	<b>-0.37</b>	4.71	74.01
AnnoQ-clue-RefQ	<u>27.43</u>	<u>0.04</u>	<u>53.62</u>	<u>93.85</u>	<u>46.89</u>	-0.38	<u>4.76</u>	<u>74.21</u>

Table 5: **SQuAD** - Performance of different QG methods on NACo and other existing metrics. The evaluation uses original SQuAD questions (RefQ) as references, with GPT-3.5 as the underlying QA system for NACo. The highest and second-highest scores (not including references for reference-based metrics) are highlighted with bold and underline markers, respectively.

QG Competitor	Ref-based metrics					Ref-free metrics			Human
	B	B-RT	R-L	BSc	Q-B	QA-S	R-Q	NACo	
<b>LM-generated</b>									
BART-base	12.97	-0.79	34.38	87.88	27.67	-0.28	2.78	41.25	2.35
GPT-3.5 (zero-shot)	7.94	-0.95	26.05	87.38	17.00	-0.28	4.25	73.60	4.20
GPT-3.5 (few-shot)	10.53	-0.87	27.51	87.76	18.67	<b>-0.27</b>	4.21	71.09	-
BART-clue-RefQ	<u>31.90</u>	<u>-0.10</u>	<b>58.77</b>	<b>92.61</b>	<b>51.98</b>	-0.28	3.31	63.17	-
<b>Human-validated</b>									
RefQ	100.00	1.00	100.00	100.00	100.00	<b>-0.27</b>	4.21	74.33	4.77
AnnoQ	13.64	-0.70	30.94	89.22	26.95	<b>-0.27</b>	<u>4.32</u>	<b>83.43</b>	5.02
AnnoQ-clue-RefQ	<b>35.59</b>	<b>-0.06</b>	<u>52.95</u>	<u>92.15</u>	<u>51.76</u>	<b>-0.27</b>	<b>4.44</b>	<u>81.60</u>	-

Table 6: **HotpotQA** - Performance of different QG methods on NACo and other existing metrics. The evaluation uses original HotpotQA questions (RefQ) as references, with GPT-3.5 as the underlying QA system for NACo. The highest and second-highest scores (not including references for reference-based metrics) are highlighted with bold and underline markers, respectively.

QG Competitor	Ref-based metrics					Ref-free metrics		Our metric NACo
	B	B-RT	R-L	BSc	Q-B	QA-S	R-Q	
<b>LM-generated</b>								
BART-base	13.71	<u>0.16</u>	32.42	89.19	-	<b>-0.15</b>	3.74	79.29
GPT-3.5 (few-shot)	9.23	-0.28	28.77	88.6	-	-0.16	3.89	79.78
BART-clue-RefQ	<u>13.63</u>	0.01	<u>40.66</u>	<u>90.09</u>	-	<b>-0.15</b>	3.33	75.29
<b>Human-validated</b>								
RefQ	100.00	1.00	100.00	100.00	-	-0.16	<u>3.99</u>	82.85
AnnoQ	14.78	-0.43	34.16	89.61	-	-0.16	3.98	<b>84.67</b>
AnnoQ-clue-RefQ	<b>26.4</b>	<b>0.59</b>	<b>50.22</b>	<b>92.1</b>	-	-0.17	<b>4.23</b>	<u>83.00</u>

Table 7: **TedEdQA** - Performance of different QG methods using NACo and other existing metrics. The evaluation uses original TedEd questions (RefQ) as references, with GPT-3.5 as the underlying QA system for NACo. Some questions in this dataset are in the fill-in-the-blank form and do not contain question words like *what*, *where*, *etc.* which Q-BLEU heavily relies on for scoring (Nema and Khapra, 2018); thus, we do not report this metric for this dataset. The highest and second-highest scores (not including references for reference-based metrics) are highlighted with bold and underline markers, respectively.