

# PIRLS Category-specific Question Generation for Reading Comprehension

Anonymous ACL submission

## Abstract

001 According to the internationally recognized  
002 PIRLS (Progress in International Reading Literacy  
003 Study) assessment standards, reading comprehension  
004 questions should encompass all four  
005 comprehension processes: retrieval, inferencing,  
006 integrating and evaluation. This paper  
007 investigates whether Large Language Models  
008 can produce high-quality questions for each of  
009 these categories. Human assessment on a Chinese  
010 dataset shows that GPT-4o can generate  
011 usable and category-specific questions, ranging  
012 from 74% to 90% accuracy depending on the  
013 category.

## 1 Introduction

014 Given the importance of asking questions for effective  
015 learning (Dillon, 2006; Etemadzadeh et al.,  
016 2013; Kurdi et al., 2020), there has been extensive  
017 effort in developing automatic Question Generation  
018 (QG) models to produce high-quality questions for  
019 reading materials in educational systems (Heilman  
020 and Smith, 2010; Lindberg et al., 2013). Through  
021 automatic creation of pedagogical and assessment  
022 material, QG benefits teachers by reducing their  
023 workload. It also levels the playing field for students,  
024 providing them with instant and free access  
025 to questions for review and practice.  
026

027 According to PIRLS (Progress in International  
028 Reading Literacy Study), reading requires four  
029 comprehension processes: retrieval, inferencing,  
030 integrating and evaluation (Mullis and Martin, 2019)  
031 (Table 1). A balanced set of questions, involving all  
032 four processes, is therefore needed to assess reading  
033 comprehension. However, existing QG benchmarks  
034 such as SQuAD (Rajpurkar et al., 2016) mostly  
035 focus on factoid short-answer questions.  
036

037 This paper investigates question generation of  
038 the four PIRLS categories with Large Language  
039 Models (LLMs) using zero-shot, few-shot and fine-  
tuning approaches. Our contribution is two-fold. In

Process	Description
Retrieval	Focus on and Retrieve Explicitly Stated Information
Inferencing	Make Straightforward Inferences
Integrating	Interpret and Integrate Ideas and Information
Evaluation	Evaluate and Critique Content and Textual Elements

Table 1: Comprehension processes in reading according to PIRLS (Mullis and Martin, 2019)

040 this first attempt of QG based on PIRLS, an inter-  
041 nationally recognized standard for reading compre-  
042 hension assessment, we show that GPT-4o can gen-  
043 erate high-quality questions with category-specific  
044 prompts. Second, we contribute a dataset of Chi-  
045 nese passages and questions, annotated with PIRLS  
046 categories, that may serve as a benchmark for fu-  
047 ture research.

## 2 Previous work

048 Early QG approaches mostly relied on heuris-  
049 tics, linguistic templates and rules (Labutov et al.,  
050 2015; Mostow et al., 2016). With the avail-  
051 ability of large-scale datasets, QG began to be  
052 formulated as a sequence-to-sequence generation  
053 task. An encoder-decoder architecture with a  
054 global attention mechanism was found to be ef-  
055 fective (Du et al., 2017; Kim et al., 2019), but  
056 can be further improved with transformer-based  
057 approaches (Scialom et al., 2019), and fully fine-  
058 tuned language models (LM) (Xiao et al., 2021).  
059 Answer-agnostic QG can be performed via joint  
060 Question and Answer Generation (QAG) (Lewis  
061 et al., 2021). A QAG model based on fine-tuning  
062 encoder-decoder LMs produces high-quality ques-  
063 tions (Ushio et al., 2022), but has not been eval-  
064 uated in terms of question type. The most recent  
065 research has adopted LLMs. On a textbook dataset,  
066

Excerpt of input passage (in Chinese):

太阳和地球虽然相距1.5亿公里，但它却会提供光和热。除此以外，它还会给地球带来意想不到的“礼物”呢！其实太阳的表面常常发生爆炸，在最活跃的时候，更会把表面的物质抛射出去，形成太阳风暴。当太阳风暴经过地球时，不但会损毁人造卫星，干扰无线电通讯， ...  
Even though the Sun is 150 million kilometers away from Earth, it provides light and heat. Besides, it also gives a surprising ‘gift’ to Earth! There are frequent explosions on the surface of the Sun ... forming solar storms. When a solar storm passes by the Earth, it not only destroys satellites and interfere with wireless communication, ...

Type	Example Question
Retrieval: word-match	太阳和地球虽然相距一亿五千万公里，但它却会提供什么？ Even though the Sun is 150 million kilometers away from Earth, What does it provide?
Retrieval: paraphrase	文章提到太阳和地球之间的距离是多少？ What is the distance between the sun and the Earth, as mentioned in the passage?
Inferenc- ing	根据文章，太阳爆炸造成的“太阳风暴”会对地球造成哪些影响？ How is the Earth affected by the solar storms caused by explosions on the Sun?
Integrat- ing	文章中提到太阳常常发生爆炸会带来什么「礼物」？ According to the passage, what ‘gift’ is brought by the frequent explosions at the Sun?
Evaluat- ion	作者认为太阳的影响对地球有什么优势和缺陷 What does the author think are the Sun’s positive and negative impact on the Earth?

Table 2: Example input passage and output questions of each PIRLS question type (Section 4)

few-shot prompting with GPT-3 was able to generate human-like questions ready for classroom use (Wang et al., 2022). A fine-tuned version of ChatGPT was able to generate questions that are competitive with human ones in terms of readability, correctness, coherence and engagement (Xiao et al., 2023).

The research most closely to ours was reported by Elkins et al. (2024). GPT-3.5 was prompted to generate six kinds of questions based on Bloom’s taxonomy (Krathwohl, 2002). In an evaluation using Wikipedia passages on biology and machine learning, the generated questions were shown to be highly semantically relevant, fluent, and answerable. For questions generated with InstructGPT reported by Elkins et al. (2023), the accuracy in question category varies from only 36.1%-40.0% for the ‘creating’ category, but higher for the more objective categories such as 83.3% to 91.7% for the ‘remembering’ category. Our study uses PIRLS, a framework that focuses more specifically on grade-school reading comprehension than Bloom’s. Further, we reported the effect of fine-tuning LLMs and contribute a dataset in Chinese, which has more limited resources for QG.

### 3 Dataset

Existing reading comprehension datasets in Chinese, such as the Delta Reading Comprehension

Dataset<sup>1</sup> and DuReader<sup>2</sup>, are primarily drawn from newspapers, Wikipedia and user logs. Further, the questions are not annotated with their categories. We therefore constructed new datasets<sup>3</sup> using Chinese-language pedagogical materials:

**Training set** The fine-tuning data consists of 804 manually composed questions about 72 passages taken from published Chinese story books. There are 201 questions at each PIRLS category.<sup>4</sup> The average passage length is 1,131 Chinese characters.

**Test set** The test set consists of 50 passages from a public reading comprehension assessment<sup>5</sup>, with 25 passages from Grade 3 and 25 from Grade 6. The average passage length is 648 Chinese characters.

## 4 Annotation Scheme

According to the International Association for the Evaluation of Educational Achievement, a reading

<sup>1</sup><https://github.com/DRCKnowledgeTeam/DRCD>

<sup>2</sup><https://github.com/baidu/DuReader>

<sup>3</sup>The test set will be made available at <http://anonymous>. Due to copyright issues, the training set will be made available for research purposes upon contact with the last author.

<sup>4</sup>181 questions were used for training and 20 for validation.

<sup>5</sup>Downloaded from the website of the Territory-wide System Assessment (TSA) <https://www.bca.hkeaa.edu.hk/web/TSA/en/\PriPaperSchema.html>.

comprehension question should address the following comprehension processes, as defined in the PIRLS standards (Table 1):

**Retrieval** The answer is explicitly given in a text span in the passage.

**Inferencing** Answering the question requires inferences about ideas or information that is not explicitly stated.

**Integrating** Answering the question “requires comprehension of the entire text, or at least significant portions of it.” (Mullis and Martin, 2019)

**Evaluation** The answer “involves a judgement about some aspect of the text”, and is not necessarily found in the passage.

Example questions can be found in Table 2.<sup>6</sup>

## 5 Approach

The input is a Chinese text, without any specified answer span. We used two LLMs — GPT-4o<sup>7</sup> and *LLaMa-3* (Cui and Yao, 2024)<sup>8</sup> to generate questions<sup>9</sup> for the text, using the following prompts:

**Zero-shot** As shown in Table 6, for each of the four PIRLS category, a different prompt describing the requirements of the category is used.<sup>10</sup>

**Generic** This is the same as the zero-shot approach, except that the prompt does not specify the question category:

基於所提供的文章，請創作一個簡答題，並提供對應的答案。  
文章:<input>

[Translation: “Based on the given passage, create a short-answer question and provide a corresponding answer. Passage: <input>]

**Few-shot** The PIRLS category-specific prompt used in zero-shot above is accompanied with

<sup>6</sup>The Chinese passage is taken from a Chinese-language public examinations at [https://www.hkeaa.edu.hk/en/sa\\_tsa/](https://www.hkeaa.edu.hk/en/sa_tsa/)

<sup>7</sup><https://openai.com/index/hello-gpt-4o/>

<sup>8</sup>Chinese 8B Instruct-v1, downloaded from <https://huggingface.co/hfl/llama-3-chinese-8b-instruct>

<sup>9</sup>max\_tokens=200; temperature=0.6; top\_p=0.9

<sup>10</sup>In all experiments, if multiple questions were generated, only the first one was kept. Regardless of whether the question was without an answer or invalid, we kept the output, and none of the questions were regenerated.

Model	Unusable	Usable	
		minor rev.	wo/ rev.
Llama-3 (generic)	4%	24%	72%
Llama-3 (zero-shot)	4%	17.5%	78.5%
Llama-3 (few-shot)	14%	15%	71%
Llama-3 (fine-tuned)	15%	26.5%	58.5%
GPT-4o (generic)	2%	10%	88%
GPT-4o (zero-shot)	<b>0%</b>	<b>4%</b>	<b>96%</b>

Table 3: Evaluation results on usability using the scale defined in Section 6

an input passage and  $N$  sample questions, according to the template in Table 8 (Appendix B). We set  $N = 5$ , with a sample passage and five questions taken from the training set.

**Fine-tuned** We fine-tuned LLaMa-3, an open-source LLM, with the PIRLS category-specific prompts Table 6 on the training set (Section 3).<sup>11</sup>

For each passage in the test set, a question was generated from each prompt type described above.

## 6 Evaluation set-up

Four assessors, all native Chinese speakers with a bachelor’s degree, annotated each generated question on its *usability* and *PIRLS category*. The order of the questions was randomized to avoid bias. Each question was independently evaluated by two of the assessors. In case of disagreement, a PIRLS expert with a Master’s degree in Education, adjudicated the decision.

First, the assessors rated the quality of the question on the following three-point scale:

**Usable without revision** The question can be used as is: it is grammatical, fluent, and relevant for the input passage.

**Usable with minor revision** The question is relevant for the input passage, but requires improvement in its linguistic quality, e.g., correction of grammatical errors, better vocabulary choice or phrasing.

<sup>11</sup>The fine-tuning was performed for 1 epoch using the following hyperparameters: learning rate=1e-4; lora\_rank=64; lora\_alpha=128; lora\_dropout=0.05; batch\_size = 1; gradient\_accumulation\_steps=8; max\_seq\_length=3303. On an A100 GPU, the training took 4 minutes and 34 seconds.

Model	PIRLS category				Average
	Retrieval	Inferencing	Integrating	Evaluation	
Llama-3 (generic)	56%	32%	8%	0%	24%
Llama-3 (zero-shot)	78%	40%	22%	20%	40%
Llama-3 (few-shot)	82%	26%	10%	4%	30.5%
Llama-3 (fine-tuned)	68%	42%	10%	34%	38.5%
GPT-4o (generic)	54%	32%	12%	0%	24.5%
GPT-4o (zero-shot)	<b>86%</b>	<b>74%</b>	<b>78%</b>	<b>90%</b>	82%

Table 4: Accuracy in question category

Category	Retrieval	Infer.	Integr.	Eval.
Retrieval	<b>43</b>	6	1	0
Infer.	8	<b>37</b>	3	2
Integr.	0	3	<b>39</b>	8
Eval.	0	0	5	<b>45</b>

Table 5: Confusion matrix of the PIRLS category of the questions generated by GPT-4o (zero-shot)

**Unusable** The question is irrelevant for the passage, or cannot be understood.

Second, the usable questions (either without revision or with minor revision) were classified in terms of PIRLS question type (Section

## 7 Results

### 7.1 Question Usability

The four assessors agreed on 90% of questions on the usable vs. unusable classification, leading to a 0.499 Kappa score, a “moderate” level of agreement (Landis and Koch, 1977).

Among questions generated by Llama-3 with the generic prompt, 72% were usable without revision. The use of category-specific prompts, which supply more detailed instructions, increased the proportion of directly usable questions to 78.5%. Providing examples through few-shot and fine-tuning resulted in more unusable questions. Anecdotal examination suggests that the model was led to overly prefer the wording in the given samples, even if it results in unnatural questions.

On GPT-4o, the category-specific prompts also led to gains in usability over the generic one. Overall, GPT-4o attained substantially superior performance, with a vast majority of the generated questions (96%) annotated as directly usable.

### 7.2 Question category

Excluding the unusable questions, the assessors agreed on 55.17% of the generated questions on

the 4-way classification of PIRLS category. This yielded a 0.494 weighted kappa score, a “moderate” level of agreement (Landis and Koch, 1977).

The generic prompt produced mostly ‘retrieval’ questions on both Llama-3 (56%) and GPT-4o (54%). It would be highly inefficient, however, for users looking for ‘Inferencing’ and ‘Integrating’ questions. The category-specific (zero-shot) prompts improved the accuracy across all categories raising the average accuracy to 40% for Llama-3 and 82% for GPT-4o. This result suggests that both models were able to understand the instructions in the prompt.

On Llama-3, the few-shot approach improved the generation of ‘retrieval’ questions to 82%. The five samples, however, may not have been sufficient for the higher-order categories, resulting in lower accuracy. The larger quantity of training data likely enabled the fine-tuned model to perform better at two of the higher-order categories, namely ‘Inferencing’ and ‘Evaluation’. The overall accuracy, however, was still offset by the other two categories.

The GPT-4o zero-shot approach offers the best performance in all categories, with an average of 82% accuracy. As shown in the confusion matrix (Table 5), most errors were within one category above or below the target in the PIRLS scale.

## 8 Conclusion

In assessing reading comprehension, it is important to use questions that target various comprehension processes. This paper has presented the first study on automatic question generation for reading comprehension based on the four categories in the PIRLS framework. Experiments on Chinese passages show that zero-shot GPT-4o can produce questions belonging to the target category at 74% to 90% accuracy, outperforming both the zero-shot and fine-tuned LLaMA-3 model.

## 9 Limitations and Ethics Consideration

The evaluation has focused on the quality of the questions, but cannot show their pedagogical impact on the students. At the time of system deployment, users should be clearly informed that the automatically generated questions should be viewed only as a first draft, to minimize the risk that the teacher may fail to edit an unusable question and pass it to students.

The experimental and evaluation protocol was approved by the (Anonymous Grant) administered by the Department of Education, (Anonymous Country).

### Acknowledgements

Anonymous.

### References

Y. Cui and X. Yao. 2024. Rethinking LLM Language Adaptation: A Case Study on Chinese Mixtral. In *arXiv preprint arXiv:2403.01851*.

James T. Dillon. 2006. Effect of questions in education and other enterprises. In *Rethinking schooling*, page 145–174. Routledge.

Xinya Du, Junru Shao, and Claire Cardie. 2017. Learning to Ask: Neural Question Generation for Reading Comprehension. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (ACL)*.

Sabina Elkins, Ekaterina Kochmar, Jackie Chi Kit Cheung, and Iulian Vlad Serban. 2024. How teachers can use large language models and bloom’s taxonomy to create educational quizzes. In *Proc. 14th Symposium on Educational Advances in Artificial Intelligence*.

Sabina Elkins, Ekaterina Kochmar, Iulian Serban, and Jackie C. K. Cheung. 2023. How Useful Are Educational Questions Generated by Large Language Models? *AIED 2023, CCIS*, 1831:536–542.

Atika Etemadzadeh, Samira Seifi, and Hamid Roohbakhsh Far. 2013. The role of questioning technique in developing thinking skills: The ongoing effect on writing skill. *Procedia-Social and Behavioral Sciences*, 70:1024–1031.

Michael Heilman and Noah A. Smith. 2010. Good Question! Statistical Ranking for Question Generation. In *Proc. Human Language Technologies: The 2010 Annual Conference of the North American Chapter of the ACL (HLT-NAACL)*, page 609–617.

Yanghoon Kim, Hwanhee Lee, Joongbo Shin, and Kyomin Jung. 2019. Improving Neural Question Generation Using Answer Separation. In *Proc. 33rd AAAI Conference on Artificial Intelligence (AAAI-19)*.

D. R. Krathwohl. 2002. A revision of Bloom’s taxonomy: An overview. *Theory into practice*, 41(4):212–218.

Ghader Kurdi, Jared Leo, Bijan Parsia, Uli Sattler, and Salam Al-Emari. 2020. A systematic review of automatic question generation for educational purposes. *International Journal of Artificial Intelligence in Education*, 30:121–204.

I. Labutov, S. Basu, and L. Vanderwende. 2015. Deep questions without deep understanding. In *Proc. ACL*.

J. Richard Landis and Gary G. Koch. 1977. The Measurement of Observer Agreement for Categorical Data. *Biometrics*, 33:159–174.

Patrick Lewis, Yuxiang Wu, Linqing Liu, Pasquale Minervini, Heinrich Küttler, Aleksandra Piktus, Pontus Stenertorp, and Sebastian Riedel. 2021. PAQ: 65 million probably-asked questions and what you can do with them. *Transactions of the Association for Computational Linguistics*, 9:1098–1115.

David Lindberg, Fred Popowich, John Nesbit, and Phil Winne. 2013. Generating natural language questions to support learning on-line. In *Proceedings of the 14th European Workshop on Natural Language Generation*, page 105–114.

Jack Mostow, Yi ting Huang, Hyeju Jang, Anders Weinstein, Joe Valeri, and Donna Gates. 2016. Developing, evaluating, and refining an automatic generator of diagnostic multiple choice cloze questions to assess children’s comprehension while reading. *Natural Language Engineering*, 23(2):245–294.

Ina V. S. Mullis and Michael O. Martin. 2019. *PIRLS 2021 Assessment Frameworks*. International Association for the Evaluation of Educational Achievement.

Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. SQuAD: 100,000+ Questions for Machine Comprehension of Text. In *Proc. Conference on Empirical Methods in Natural Language Processing (EMNLP)*, page 2383–2392.

Thomas Scialom, Benjamin Piwowarski, and Jacopo Staiano. 2019. Self-Attention Architectures for Answer-Agnostic Neural Question Generation. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, page 6027–6032.

Asahi Ushio, Fernando Alva-Manchego, and Jose Camacho-Collados. 2022. Generative Language Models for Paragraph-Level Question Generation. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, page 670–688.

Z. Wang, J. Valdez, D. Basu Mallick, and R. G. Baraniuk. 2022. Towards Human-Like Educational Question Generation with Large Language Models. *Artificial Intelligence in Education. AIED 2022. Lecture Notes in Computer Science*, 13355.

353 Changrong Xiao, Sean Xin Xu, Kunpeng Zhang, Yufang  
354 Wang, and Lei Xia. 2023. Evaluating reading com-  
355 prehension exercises generated by llms: A showcase  
356 of chatgpt in education applications. In *Proc. 18th*  
357 *Workshop on Innovative Use of NLP for Building*  
358 *Educational Applications*, page 610–625.

359 Dongling Xiao, Han Zhang, Yukun Li, Yu Sun, Hao  
360 Tian, Hua Wu, and Haifeng Wang. 2021. Ernie-gen:  
361 an enhanced multi-flow pre-training and fine-tuning  
362 framework for natural language generation. In *Pro-*  
363 *ceedings of the Twenty-Ninth International Confer-*  
364 *ence on International Joint Conferences on Artificial*  
365 *Intelligence*, page 3997–4003.

## 366 **A Appendix: Instruction to Human** 367 **Assessors**

368 The human assessors gave consent to the data col-  
369 lection and were informed that the results would  
370 remain anonymous. They were shown the follow-  
371 ing instructions:

372 <passage>  
373 <question>

- 374 1. Is the question understandable and relevant  
375 for the passage?
- 376 2. Does the language quality of the question need  
377 to be improved?
- 378 3. If the answer to #1 is “Yes”, choose one of the  
379 categories for the question:
  - 380 • Retrieval (Focus on and Retrieve Explic-  
381 itly Stated Information)
  - 382 • Inferencing (Make Straightforward Infer-  
383 ences)
  - 384 • Integrating (Interpret and Integrate Ideas  
385 and Information)
  - 386 • Evaluation (Evaluate and Critique Con-  
387 tent Textual Elements)

## 388 **B Appendix: Few-shot prompt template**

389 The few-shot template is shown in Table 8.

Type	Prompt (in Chinese)
System prompt	你是一個能幹的閱讀理解問題生成器，始終遵循給定的說明和要求來生成問題。
Retrieval questions (PIRLS level 1)	基於所提供的文章，請創作一個屬於PIRLS第一層次的簡答題，並提供對應的答案。這個問題應著重於檢索文本中明確表述的信息，也就是資訊檢索型的問題。此類問題要求考生識別和回憶文本中明確提到的信息，如事件的順序、角色的特徵或進行比較等。 文章:{input passage}
Inferencing questions (PIRLS level 2)	基於所提供的文章，請創作一個屬於PIRLS第二層次的簡答題，並提供對應的答案。這個問題應鼓勵考生從文本中進行直接推理，進一步超越單純的信息提取，也就是需要進行簡單推理的問題。這類問題需要考生進行直接推理，例如理解因果關係或推測未明確陳述但可以從文本邏輯推導出的結果。 文章:{input passage}
Integrating questions (PIRLS level 3)	基於所提供的文章，請創作一個屬於PIRLS第三層次的簡答題，並提供對應的答案。這個問題應促使考生解釋想法並整合文本不同部分信息，也就是需要進行解釋及整合的問題。這類問題需要考生全面理解並能夠從文本的不同部分綜合信息，如解釋角色的感受和行為，並整合文本中的想法和信息。 文章:{input passage}
Evaluation questions (PIRLS level 4)	基於所提供的文章，請創作一個屬於PIRLS第四層次的簡答題，並提供對應的答案。這個問題應需要考生批判性地檢視和評估文本內容、語言和文本元素，也就是評鑒型的問題。這類問題是最高層次的問題，問題挑戰考生批判性地評估文本的內容、語言和文本元素，如對價值、期望和接受度作出判斷，或考慮他們如果處於某個角色的位置會如何反應。 文章:{input passage}

Table 6: LLM prompts for generating questions for each PIRLS category

Type	Prompt (in English)
System prompt	You are a capable reading comprehension question generator, always following the given instructions and requirements to generate questions.
Retrieval questions (PIRLS level 1)	Based on the article provided, please create a short answer question belonging to PIRLS level 1 and provide the corresponding answer. This question should focus on retrieving information explicitly stated in the text, i.e. an information retrieval type question. This kind of question requires candidates to identify and recall information explicitly mentioned in the text, such as the sequence of events, character traits, or making comparisons.  article:{input passage }
Inferencing questions (PIRLS level 2)	Based on the article provided, please create a short answer question belonging to PIRLS level 2 and provide the corresponding answer. This question should encourage candidates to make straightforward inferences from the article, moving further beyond information retrieval, i.e. a question requiring simple inferences. This type of question requires candidates to make straightforward inferences, such as understanding cause and effect relationships or inferring consequences that are not explicitly stated but can be logically deduced from the text.  article:{input passage }
Integrating questions (PIRLS level 3)	Based on the article provided, please create a short answer question belonging to the PIRLS level 3 and provide the corresponding answer. This question should prompt the candidate to interpret ideas and integrate information from different parts of the text, i.e. a question that requires interpretation and integration. This type of question requires candidates to have a comprehensive understanding and be able to integrate information from different parts of the text, such as explaining a character's feelings and actions, and integrating ideas and information across the text.  article:{input passage }
Evaluation questions (PIRLS level 4)	Based on the article provided, please create a short answer question belonging to PIRLS level 4 and provide the corresponding answer. This question should require candidates to critically examine and evaluate the text content, language, and textual elements, i.e. an evaluative question. This type of question is the highest-level question that challenges candidates to critically evaluate a text content, language, and textual elements, such as making judgments about value, desirability, and acceptability or considering how they would react if they were in a character's position.  article:{input passage }

Table 7: LLM prompts for generating questions for each PIRLS category (translated)



---

{category-specific prompt}  
範例文章及相應的範例問題(請參考範例來創作問題):{範例文章:{example passage}  
PIRLS第{required level}層次範例問題1:{example question-answer pair 1}  
...  
PIRLS第{required level}層次範例問題5:{example question-answer pair 5}}  
文章: {input passage}

---

Table 8: Prompt template for few-shot question generation