

# Unlocking Markets: A Multilingual Benchmark to Cross-Market Question Answering

Yifei Yuan<sup>1</sup>, Yang Deng<sup>2</sup>, Anders Søgaard<sup>1</sup>, Mohammad Aliannejadi<sup>3</sup>

<sup>1</sup>University of Copenhagen, <sup>2</sup>Singapore Management University, <sup>3</sup>University of Amsterdam  
{yiya,soegaard}@di.ku.dk  
ydeng@smu.edu.sg, m.aliannejadi@uva.nl

## Abstract

Users post numerous product-related questions on e-commerce platforms, affecting their purchase decisions. Product-related question answering (PQA) entails utilizing product-related resources to provide precise responses to users. We propose a novel task of Multilingual Cross-market Product-based Question Answering (MCPQA) and define the task as providing answers to product-related questions in a main marketplace by utilizing information from another resource-rich auxiliary marketplace in a multilingual context. We introduce a large-scale dataset comprising over 7 million questions from 17 marketplaces across 11 languages. We then perform automatic translation on the Electronics category of our dataset, naming it as McMarket. We focus on two subtasks: review-based answer generation and product-related question ranking. For each subtask, we label a subset of McMarket using an LLM and further evaluate the quality of the annotations via human assessment. We then conduct experiments to benchmark our dataset, using models ranging from traditional lexical models to LLMs in both single-market and cross-market scenarios across McMarket and the corresponding LLM subset. Results show that incorporating cross-market information significantly enhances performance in both tasks.

## 1 Introduction

Online shoppers on platforms such as Amazon post numerous questions related to specific products every day (McAuley and Yang, 2015). Product-related question answering (PQA) involves providing accurate and informative responses to these questions. By leveraging product-related information, such as reviews and product meta information, responses to product-related questions can be expanded, offering enhanced depth and authenticity for potential customers (Gupta et al., 2019).

The recent success in cross-market PQA underscores the capability to effectively leverage rele-

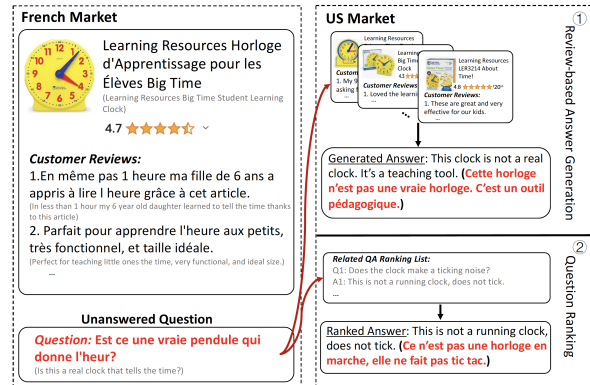


Figure 1: An example of enhancing product-related QA using cross-market data. ① depicts generating answers with cross-market reviews. ② depicts ranking-related cross-market questions to find the answer.

vant questions from a resource-rich marketplace to address questions in a resource-scarce marketplace (Shen et al., 2023; Ghasemi et al., 2023). In this work, we extend the hypothesis that leveraging knowledge from popular marketplaces can also enhance the quality of answers in less common marketplaces, even in a different language. As shown in Figure 1, for a question to a product in the French marketplace (denoted as **main marketplace**) asking if the clock is a real one, we can either address it by examining reviews of the same product or similar ones in the much larger US marketplace (denoted as **auxiliary marketplace**), or ranking related questions from both main and auxiliary marketplaces to find the answer. These multilingual reviews and related questions serve as valuable hints, by saying “it’s not a real clock.”, thereby providing crucial information for the pertinent question at hand.

We, therefore, propose a novel task of *Multilingual Cross-market Product-based Question Answering* (MCPQA). We define this task as producing the answer to a product-related question in an original marketplace, using information sourced from an auxiliary marketplace with richer resources, within a multilingual setting. To this

end, our initial goal is to address the following research question **RQ1**: *In a multilingual context, how can we utilize an auxiliary marketplace to enhance question-answering in the main marketplace by leveraging product-related resources (i.e., questions, reviews)?* To answer **RQ1**, we propose the first large-scale MCPQA dataset, covering 17 different marketplaces (including the **us** auxiliary marketplace and 16 main marketplaces) across 11 different languages from real Amazon product QA sources. Specifically, our dataset consists of over 7 million product-related questions with a total of 52 million product reviews. Different from existing PQA datasets, more diverse information is provided in the dataset, exploring the possible answers with both questions and reviews available. Additionally, we perform automatic translation on the Electronics category of the dataset, naming it McMarket. We then perform comprehensive data analysis on McMarket to address **RQ1**. We demonstrate a notable increase in the percentage of review-answerable questions across all marketplaces, with support from the auxiliary **us** marketplace.

Given the recent success of large language models (LLMs) in NLP tasks (Touvron et al., 2023a; OpenAI, 2023), their potential application to the MCPQA task prompts our second research question **RQ2**: *Can LLMs benefit the dataset construction in the MCPQA task?* Delving into **RQ2**, on McMarket, we create a subset by randomly selecting some questions from each marketplace and perform GPT-4 auto-labeling. Specifically, we focus on two widely-studied PQA subtasks under the multilingual cross-market settings, including **review-based answer generation (AG)** (Gao et al., 2019; Chen et al., 2019) and **product-related question ranking (QR)** (Rozen et al., 2021). For AG, we ask LLMs to judge whether a question is answerable from associated reviews and provide its corresponding answer. We denote the subset as  $\text{McMarket}_r$ . For QR, given two QA pairs, we ask LLMs to judge if one helps answer the other and denote the subset as  $\text{McMarket}_q$ . With these two subsets, we conduct human assessment to analyze LLM-generated results from multiple perspectives. Surprisingly, in  $\text{McMarket}_r$ , 61.8% LLM-generated answers are assumed ‘better’ than the human ground truth.

Finally, we are interested in answering the research question **RQ3**: *In the multilingual context, how can we effectively leverage the unique features of cross-market information to enhance*

*product-related question answering?* To this end, we perform experiments of models on AG and QR subtasks. For each task, we report the performance of state-of-the-art methods under single- and cross-market scenarios on both McMarket and the corresponding LLM-labeled subset. We benchmark methods ranging from traditional lexical models (i.e., BM25) to LLM-based approaches (i.e., LLaMA-2, Flan-T5). We demonstrate the superiority of cross-market methods against their single-market counterparts on both subtasks.

In conclusion, our contributions are as follows:

- We propose a novel task named MCPQA, where product-related information from an auxiliary marketplace is leveraged to answer questions in a resource-scarce marketplace in a multilingual setting. Specifically, we investigate two subtasks named AG and QR.
- We benchmark a large-scale real-world dataset to facilitate the research in the MCPQA task. We also collect two LLM-annotated subsets and adopt human assessment to analyze their characteristics.
- To provide a comprehensive evaluation of the task and verify the superiority of cross-market methods, experiments are performed under both single/cross-market scenarios.<sup>1</sup>

## 2 Related Work

**Product-related QA.** Product-related QA (PQA) seeks to address consumers’ general inquiries by utilizing diverse product-related resources such as customer reviews, or the pre-existing QA sections available on a retail platform (Yu et al., 2012; Deng et al., 2023). Among the existing literature in this area, retrieval-based methods have been a popular direction that retrieve related reviews for providing the right answer (Wan and McAuley, 2016; Zhang et al., 2019b; Yu and Lam, 2018; Zhang et al., 2020a,b,c,d). For example, McAuley and Yang (2015) propose a model that leverages questions from previous records for selecting the relevant review for the question. While most of these works assume there are no user-written answers available, Zhang et al. (2020b) rank answers for the given question with review as an auxiliary input. Another line of research (Gao et al., 2019;

<sup>1</sup>Data and code available in <https://github.com/yfyuan01/MCPQA>.

Chen et al., 2019; Gao et al., 2021; Feng et al., 2021; Deng et al., 2020, 2022) investigates answer generation grounding on retrieved product-related documents. More recently, Ghasemi et al. (2023) introduce a novel task of utilizing available data in a resource-rich marketplace to answer new questions in a resource-scarce marketplace. Building upon their research, we expand the scope to a multilingual scenario, exploring additional marketplaces with non-English content. Furthermore, we explore both questions and review information from the auxiliary marketplace.

**Cross-domain and cross-lingual QA.** Our work can be seen as a special format of cross-domain QA in E-commerce, which involves addressing questions that span different domains or fields of knowledge (Deng et al., 2018; Qu et al., 2020; Liu et al., 2019; Longpre et al., 2020; Yuan and Lam, 2021; Abbasiantaeb et al., 2023). For instance, Yu et al. (2017) propose a general framework that effectively applies the shared knowledge from a domain with abundant resources to a domain with limited resources. Also, cross-domain QA is often in close connection to cross-lingual QA in the sense that both involve transferring knowledge and understanding from one domain or language to another (Artetxe et al., 2019; Clark et al., 2020; Zhang et al., 2019a). Asai et al. (2020) expand the scope of open-retrieval question answering to a cross-lingual setting, allowing questions in one language to be answered using contents from another language. Most recently, Shen et al. (2023) introduce a multilingual PQA dataset called xPQA where cross-market information is also leveraged to aid the product-based question answering.

### 3 Problem Formulation

We investigate two subtasks of the MCPQA task, *review-based answer generation (AG)* and *product-related question ranking (QR)*, where answers to a product question are obtained by a generative or ranking way, respectively.

**Review-based answer generation.** In this task, we assume that the answer can be obtained from the reviews of the product (or similar products). Based on the setting in (Gupta et al., 2019), we define this task in a multilingual cross-market scenario. Given a question  $Q$  in the main marketplace  $M_T$ , we first retrieve and rank all the related reviews from similar items within both  $M_T$  and auxiliary marketplace  $M_A$ . Given the retrieved review set

$\Omega = \{R_1, \dots, R_k\}$ , we predict if  $Q$  is answerable from it by assigning a tag  $t$ . If yes, a generative function  $\Gamma$  is learned:  $A = \Gamma(Q, \Omega)$ , so that answer  $A$  is generated with both  $Q$  and  $\Omega$  as input.

**Product-related question ranking.** Following the problem setting in (Ghasemi et al., 2023), we assume that there are similar questions already asked about the product or similar products in other marketplaces. Therefore, given a *main marketplace* in language  $L_M$ , denoted as  $M_T$ , which usually suffers from resource scarcity of the number of knowledgeable users answers,  $M_T$  consists of several items  $\{I_1, \dots, I_m\}$ , where each  $I_k$  contains a set of question answering pairs  $\{QA_{k1}, \dots, QA_{kn}\}$ . Besides, there also exists a high-resource marketplace  $M_A$ , denoted as the *auxiliary marketplace* (the **us** marketplace in our case) in language  $L_A$  (note that in some cases  $L_A$  can be the same as  $L_M$ ). Similarly,  $M_A$  also includes several items  $\{I'_1, \dots, I'_z\}$ , where we can assume  $z \gg m$ . The task is defined as, for a given question  $Q$  in the main marketplace  $M_T$ , in a multilingual setting, we rank the questions from both  $M_T$  and  $M_A$  to take the corresponding answers of the top ranks as the possible answer to  $Q$ .

## 4 Data Collection & Analysis

We describe how we collect our dataset and perform several analysis to answer **RQ1** and **RQ2**.

### 4.1 Data collection

#### 4.1.1 Data preprocessing

We construct our dataset on top of an Amazon product dataset called XMarket (Bonab et al., 2021). XMarket includes authentic Amazon product metadata and user-generated reviews. Specifically, we sample 17 marketplaces covering 11 different languages from it. For each marketplace, we gather metadata and reviews for each product from XMarket. We also collect the question-answering pairs posed by the users by crawling the Amazon website. We then provide the corresponding English translation for the non-English contents of the Electronics category, naming it as McMarket. Specifically, we adopt the professional translation tool by DeepL<sup>2</sup> for all the QA translation and the pre-trained NLLB model (team et al., 2022) fine-tuned on each non-English language for review translation. To the best of our knowledge, this is the first multilingual cross-market QA dataset with questions and reviews in

<sup>2</sup><https://www.deepl.com/>

Name	# markets	# languages	# products	# questions	# reviews	Average QPM
xPQA (Shen et al., 2023)	12	12	16,615	18,000	-	1,500
XMarket-QA (Ghasemi et al., 2023)	2	1	34,100	4,821,332	-	2,410,666
semiPQA (Shen et al., 2022)	1	1	-	11,243	-	11,243
SubjQA (Bjerva et al., 2020)	1	1	-	10,098	10,098	10,098
ReviewRC (Xu et al., 2019)	1	1	-	2,596	959	2,596
AmazonQA (Gupta et al., 2019)	1	1	155,375	923,685	8,556,569	923,685
Amazon (McAuley and Yang, 2015)	1	1	191,185	1,447,173	13,498,681	1,447,173
Ours	17	11	143,068	7,268,393	52,469,322	427,552

Table 1: Comparison of our dataset with existing PQA datasets. QPM denotes question per marketplace.

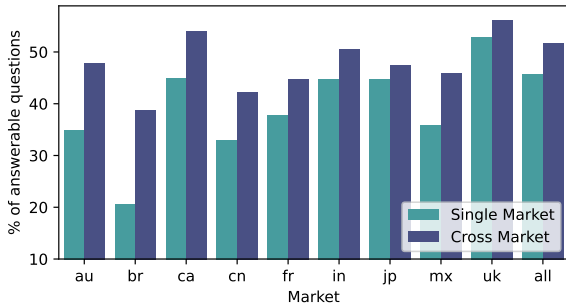


Figure 2: Portion of answerable questions in McMarket using single/cross-market review information.

the community (more information about privacy and license in Section 7).

#### 4.1.2 LLM annotation

For the two concerned subtasks, we both provide LLM-labeled data for supervised training. Specifically, we randomly select some data from McMarket and instruct GPT-4 to perform annotation. For AG, we randomly select 1000 questions per marketplace.<sup>3</sup> Then, we follow the typical top-K pooling technique (González and Gómez, 2007) and pool the top five retrieved reviews from a variety of retrieval methods. Next, we instruct GPT-4 to evaluate whether the question is answerable. If it is, GPT-4 generates an appropriate response using the question and reviews as input. If no, GPT-4 is instructed to output the reason and ‘no answer’. We denote this subset as McMarket<sub>r</sub>. For QR, we randomly select 200 questions from each marketplace. Employing the same strategy, we retrieve the top five related question-answering pairs from both the main and auxiliary marketplaces. Consequently, we acquire 1,000 question-answering pairs for each marketplace, with 9k pairs in total. Then, GPT-4 is instructed to determine if the retrieved QA pairs would be useful in answering the original question by assigning a score from 0–2, representing ‘Very useful’, ‘Partially useful’, and ‘Not useful’, respec-

<sup>3</sup>For the **au** marketplace, the total is 584 questions, so we sample all of them.

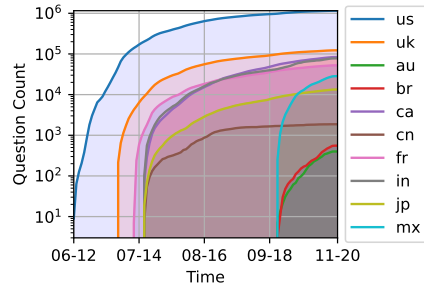


Figure 3: Temporal gap analysis.

tively. We denote this subset as McMarket<sub>q</sub>. More details of the subsets as well as the prompts we gave to GPT-4 are listed in Appendix A.

## 4.2 Data analysis

### 4.2.1 Dataset overview

Overall, our dataset covers marketplaces ranging from those with a small scale (*i.e.*, **au**, **br**) to those with rich resources (*i.e.*, **uk**, **us**). It contains over 7 million product-related questions, 52 million reviews, and 143k unique products in total.

We compare our dataset with existing PQA datasets. According to Table 1, our dataset exhibits advantages in various aspects: (1) **contains multiple languages** – we provide product, question, and review information in the original text of their respective marketplaces and additionally offer the corresponding English translations; (2) **supports cross-market QA** – our dataset is designed to facilitate question answering research across different marketplaces, enhancing its utility for cross-market analyses and evaluations; (3) **includes diverse information** – compared with existing multilingual PQA dataset, McMarket encompasses comprehensive question and review information, paving the way for more diverse research avenues and tasks in the future; (4) **is large in scale** – overall, McMarket surpasses most PQA datasets in terms of size, ensuring it comprises a substantial amount of data for experimentation and analysis.

	<b>au</b>	<b>br</b>	<b>ca</b>	<b>cn</b>	<b>fr</b>	<b>in</b>	<b>jp</b>	<b>mx</b>	<b>uk</b>	<b>us</b>	Total
Language	en	pt	en	cn	fr	en	jp	es	en	en	-
Question Num.	584	1,378	101,126	3,324	66,536	115,829	17,418	34,433	164,848	1,782,092	2,287,568
Review Num.	3,062	3,650	575,052	1,893	359,703	240,167	130,604	125,317	775,900	4,169,476	6,384,824
Product Num.	85	95	5,432	210	2,199	2,085	903	1,464	4,406	29,976	30,606
Mean ques. len	12.0±6.6	10.3±6.4	12.7±7.3	10.2±8.1	15.2±6.9	10.1±6.0	20.3±15.7	10.9±6.8	13.6±7.6	13.4±7.8	13.3±7.9
Medium ques. len.	10	8	10	8	14	8	14	9	11	11	11
Mean review len.	25.5±30.0	17.5±25.3	29.9±50.4	56.4±60.2	39.1±49.7	21.8±42.9	28.7±36.8	28.7±36.8	40.1±68.4	59.3±93.0	51.5±84.0
Medium review len.	15	10	14	39	26	10	46	21	20	30	26

Table 2: Overall statistics of the McMarket dataset. The length is reported on the token level.

	Very Bad	Bad	Good	Very Good
Correctness	2.5	0.9	8.5	88.1
Completeness	4.9	1.3	15.6	78.2
Relevance	3.5	2.7	13.4	80.4
Naturalness	0.8	0.9	5.4	92.9
<b>Better than Ground Truth</b>				<b>61.8</b>

Table 3: Human evaluation on McMarket<sub>r</sub>. All the numbers are shown in percentage.

#### 4.2.2 Dataset Statistics

Our full dataset contains product information from 17 different marketplaces, **au**, **br**, **ca**, **cn**, **de**, **es**, **fr**, **it**, **in**, **jp**, **mx**, **nl**, **sa**, **sg**, **tr**, **uk**, **us** respectively, covering 11 languages including **en**, **ar**, **cn**, **de**, **es**, **fr**, **it**, **nl**, **jp**, **pt**, **tr**. To reduce costs and facilitate baseline model training, we automatically translate the non-English contents in the Electronics category and abandon marketplaces with insufficient QA pairs. We name it as McMarket. Table 2 shows the detailed statistics of McMarket.

#### 4.2.3 Cross-market QA analysis

To answer **RQ1**, we compare the effect of product-related resources (*i.e.*, reviews) on question answering under both single- and cross-market scenarios. Figure 2 shows the comparison of answerable questions based on both single- and cross-market retrieved reviews in McMarket.<sup>4</sup> We notice that the portion of answerable questions gets raised in every marketplace with cross-market reviews, with a particularly significant uplift observed in low-resource marketplaces (*i.e.*, **br**). This verifies the transferability of knowledge across marketplaces and underscores the advantages of leveraging cross-market information in enhancing the performance of product QA models.

We further analyze the temporal characteristics of McMarket. Figure 3 illustrates the cumulative sum of the number of QA data available on all the items in all marketplaces. There are several notable

<sup>4</sup>We adopt the answerable question prediction model in (Gupta et al., 2019) to predict if a question is answerable or not given the review information.

	Incorrect	Partially correct	Correct
Portion	6.0	10.9	83.0
Overall Precision			98.2
Overall Recall			97.4
<b>Overall F1</b>			<b>97.6</b>

Table 4: Human evaluation on McMarket<sub>q</sub>. All the numbers are shown in percentage.

observations: 1) at the beginning, all marketplaces feature very few QA data. 2) At each timestep, the most resource-rich marketplace (*i.e.*, **us**) always dominates the number of QA data compared to other marketplaces by several orders of magnitude. 3) Over time, the resource intensity levels of different marketplaces continue to change. For example, the number of QA data in **mx** surpasses that in **cn** and **jp** after 2018/09. We further observe that, on average, over 70% of the questions in the main marketplace have already been answered in the **us** auxiliary marketplace under the same item, before the first question even receives an answer. These findings confirm the practicality and importance of exploring how auxiliary marketplaces can be utilized as valuable resources for PQA.

#### 4.2.4 LLM-generated data analysis

To assess the quality of LLM-generated data, we perform several analyses. On both McMarket<sub>r</sub> and McMarket<sub>q</sub>, we randomly select 50 questions from each marketplace, and hire 3 crowd-workers<sup>5</sup> to manually assess the GPT-4 labels.

**Review-based generation.** For McMarket<sub>r</sub>, we ask the crowd-workers to assess GPT-4-generated answers in terms of correctness, completeness, relevance, and naturalness. The detailed definitions of them are listed in Appendix C. For each metric, we asked them to assign a score from  $-2$  to  $+2$  to assess the answer quality, with  $-2$  representing ‘very bad’ and  $+2$  representing ‘very good.’ We also asked them to choose the better answer

<sup>5</sup>We hire the crowd-workers via a professional data management company named Appen (<https://appen.com/>).

Method	au		br		ca		cn		fr		in		jp		mx		uk		AVG		
	B	R	B	R	B	R	B	R	B	R	B	R	B	R	B	R	B	R	B	R	
Single	BM25	6.1	7.0	4.9	6.9	6.9	7.7	4.8	5.2	8.0	8.1	4.7	5.9	11.0	9.6	7.0	8.2	10.3	9.3	8.0	7.9
	BERT	7.4	7.3	9.0	5.3	7.3	6.8	5.4	5.0	8.5	7.2	5.1	4.8	10.6	8.7	9.4	7.7	9.5	8.2	7.9	7.0
	T5	15.5	11.4	14.3	12.6	16.4	12.1	13.5	10.7	16.5	11.5	12.8	9.9	22.6	15.6	20.2	14.4	18.9	13.3	16.9	12.2
	Llama-2*	10.2	14.7	16.4	17.1	15.9	13.1	14.8	13.6	18.3	14.2	13.5	13.1	26.6	19.7	22.3	16.6	20.1	18.3	17.8	15.4
Cross	BM25	10.6	7.9	9.0	6.1	7.8	7.9	4.6	5.4	9.0	8.2	5.6	6.1	11.3	9.5	9.9	9.1	10.4	9.2	8.9	8.0
	BERT	10.5	8.1	9.5	6.4	8.5	8.9	5.8	5.1	9.8	8.3	6.1	7.3	11.8	9.6	10.4	8.7	11.4	10.3	9.4	9.0
	Exact-T5	14.0	11.8	16.6	13.0	18.2	11.9	13.0	11.0	18.1	11.3	12.5	10.1	22.7	15.0	20.3	14.2	20.6	13.7	17.9	12.3
	T5	16.1	11.3	17.0	14.1	17.0	12.7	15.1	11.3	19.4	12.6	13.2	10.6	23.6	16.0	22.3	16.6	20.2	15.4	18.1	13.5
	Exact-Llama-2*	19.5	15.1	17.4	15.5	16.4	13.8	15.6	11.4	21.6	17.6	<b>16.9</b>	<b>15.1</b>	27.3	17.8	24.7	17.8	22.4	19.8	20.1	17.0
	Llama-2*	<b>21.4</b>	<b>20.6</b>	<b>18.9</b>	<b>19.5</b>	<b>19.5</b>	<b>14.4</b>	<b>17.6</b>	<b>15.5</b>	<b>22.0</b>	<b>19.0</b>	16.5	15.0	<b>29.5</b>	<b>18.6</b>	<b>25.7</b>	<b>19.2</b>	<b>25.0</b>	<b>22.7</b>	<b>21.7</b>	<b>18.3</b>

Table 5: Experimental results of AG on McMarket. Where B denotes BLEU-4, R denoted ROUGE-L. \* denotes LLM based methods. The best-performed model in the single-market setting is highlighted in light grey. The models in dark grey are highlighted to distinguish from their Exact- counterparts.

Method	au		br		ca		cn		fr		in		jp		mx		uk		AVG		
	B	R	B	R	B	R	B	R	B	R	B	R	B	R	B	R	B	R	B	R	
Single	BM25	10.3	11.7	10.7	12.5	8.3	13.0	8.5	10.1	11.6	15.7	11.7	14.3	12.8	12.1	13.3	13.6	12.4	14.7	10.7	13.3
	BERT	12.4	10.0	14.8	8.7	11.3	8.8	8.5	7.1	11.1	10.2	12.0	10.6	10.9	9.0	14.1	9.5	9.0	11.1	10.8	9.5
	T5	29.8	27.0	26.7	33.6	29.2	27.4	31.1	24.2	34.9	30.8	29.0	32.2	31.1	27.0	27.2	26.5	29.5	25.9	29.9	28.4
	Llama-2*	35.7	34.3	37.6	<b>40.8</b>	36.3	37.2	38.7	34.3	35.7	32.6	34.4	35.8	34.7	32.4	35.9	34.7	35.4	37.0	35.4	35.9
Cross	BM25	13.5	11.0	12.9	10.0	13.4	12.2	7.4	8.5	12.8	13.0	14.6	15.0	11.6	10.1	15.5	12.6	12.0	15.2	12.6	12.0
	BERT	15.8	10.6	15.7	11.0	14.4	9.8	6.8	8.1	12.2	14.2	13.0	12.1	13.8	11.3	15.7	11.1	10.1	13.1	12.9	11.3
	Exact-T5	30.9	28.2	30.1	29.0	29.3	30.7	29.8	26.7	34.7	31.7	31.8	30.3	30.0	24.6	27.3	28.0	29.1	25.9	30.3	28.4
	T5	<b>32.0</b>	<b>30.2</b>	<b>31.0</b>	<b>28.6</b>	<b>29.9</b>	<b>29.7</b>	<b>32.1</b>	<b>26.8</b>	<b>32.2</b>	<b>31.5</b>	<b>30.1</b>	<b>32.4</b>	<b>36.3</b>	<b>29.9</b>	<b>29.4</b>	<b>27.6</b>	<b>30.2</b>	<b>26.0</b>	<b>31.4</b>	<b>29.1</b>
	Exact-Llama-2*	<b>37.0</b>	34.6	34.1	32.6	38.0	39.9	33.0	35.2	<b>40.8</b>	<b>44.3</b>	36.2	40.2	38.0	34.7	38.4	37.8	35.2	37.9	36.7	37.3
	Llama-2*	35.9	<b>37.4</b>	<b>38.0</b>	37.9	<b>39.2</b>	<b>40.2</b>	<b>39.1</b>	<b>36.9</b>	39.6	41.7	<b>37.0</b>	<b>41.0</b>	<b>40.9</b>	<b>35.2</b>	<b>38.8</b>	<b>37.1</b>	<b>35.9</b>	<b>38.5</b>	<b>38.4</b>	<b>38.5</b>

Table 6: Experimental results of AG on McMarket<sub>r</sub>.

between the GPT-4 generated response and the human-provided ground truth, without disclosing the true category. From Table 3, we note that GPT-4 answers demonstrate reasonable performance in terms of every metric. Surprisingly, our findings reveal that in most cases, human assessors perceive GPT-4 results to be better than human-generated ground truth. Notably, GPT-4’s outcomes are derived solely from reviews, whereas human ground truth relies on both reviews and user experiences.

**Question ranking.** For McMarket<sub>q</sub>, we ask the workers to judge the GPT-4-generated question ranking quality, by assigning a score between 0–2 to each sample, where 0 denotes GPT-4 answers are not correct, 1 as partially correct, and 2 as completely correct. Furthermore, we instruct the annotators to provide their own judgment of the ranking score if they mark GPT-4 answers as 0 or 1. Table 4 shows that the quality of the generated question ranking results by GPT-4 is also deemed satisfactory, achieving over 93% correctness in question ranking pairs and an overall F1 score of 97.6%.

## 5 Experiments

### 5.1 Experimental setup

**Dataset.** We perform experiments on AG and QR. For each task, we report the single/cross-market

results on the whole dataset and its subset.

For AG, on the McMarket dataset, we first adopt the BERT classifier trained in (Gupta et al., 2019). It assesses each question based on the review information, categorizing them as either answerable or unanswerable. Subsequently, we employ it to filter out all answerable questions. We then split the training/validation/testing sets following the portion of 70/10/20%, resulting in 183,092/24,973/49,958 samples, respectively. On the McMarket<sub>r</sub> dataset, we also split the data into three sets with the same portions. Specifically, we adopt the GPT-4 generated answers as the ground truth. In the single-market setting, we retrieve the top  $K$  reviews from the main marketplace before generating the answers<sup>6</sup>. In the cross-market setting, we retrieve the reviews from both the main and auxiliary marketplaces. We report the generation performance of baselines on the testing set.

For QR, we first rank products, then among the top  $N$  products, we rank the top  $K$  questions<sup>7</sup>. Since McMarket does not come with any ground-truth ranking results, we perform unsupervised training and adopt GPT-4-labeled data, McMarket<sub>q</sub>, as the testing set. Besides, to further test the per-

<sup>6</sup>We choose  $K = 5$  in our case.

<sup>7</sup>Following (Ghasemi et al., 2023), we use  $N = 3$  and  $K = 50$ .

Method	au		br		ca		cn		fr		in		jp		mx		uk		AVG		
	M	P	M	P	M	P	M	P	M	P	M	P	M	P	M	P	M	P	M	P	
Single	BM25	24.5	16.9	15.2	18.3	31.5	28.7	22.0	28.7	21.0	34.7	44.4	46.0	23.8	31.5	28.9	38.7	38.4	40.2	27.7	31.5
	BERT	26.9	43.0	18.2	35.0	30.4	42.8	18.2	34.3	17.7	40.8	47.9	52.7	28.5	34.2	30.0	47.0	40.0	51.8	28.6	42.4
	UPR-m	30.4	46.0	21.9	39.3	31.9	48.0	36.2	45.5	36.3	43.7	25.7	56.3	34.7	43.3	39.5	54.2	32.5	52.7	32.1	47.7
	UPR-l*	38.9	48.8	27.8	43.3	36.5	49.7	38.1	48.3	42.5	47.3	35.2	59.8	43.3	47.2	49.0	57.2	38.9	55.5	38.9	50.8
Cross	BM25	51.2	45.2	47.4	40.0	51.0	47.5	50.2	46.8	50.8	44.3	58.0	57.5	54.6	45.5	59.0	54.3	50.8	57.5	52.6	48.7
	Exact-BERT	50.7	38.8	49.1	41.8	48.8	47.0	46.2	46.5	50.1	44.7	59.0	57.3	54.8	45.8	59.3	55.7	51.2	57.3	52.1	48.3
	BERT	52.3	45.7	49.7	42.8	50.4	48.8	49.3	44.2	49.4	43.5	60.5	58.3	55.9	46.0	59.7	57.0	52.5	59.3	53.3	49.5
	CMJim	57.5	56.7	52.4	49.3	53.3	57.7	54.0	50.5	56.9	54.3	62.9	66.8	58.4	53.2	64.9	63.8	52.9	62.7	57.0	57.2
	UPR-m	59.1	55.5	57.8	56.0	54.3	58.5	52.8	52.1	54.9	52.3	64.1	64.3	57.5	52.9	62.8	63.7	53.6	64.5	57.4	57.8
	Exact-UPR-l*	59.3	56.0	56.3	57.1	59.7	59.5	54.4	53.7	55.4	54.0	65.6	68.8	58.5	53.3	62.4	62.9	54.1	62.8	58.4	58.7
	UPR-l*	60.0	59.5	57.7	57.5	59.0	63.2	61.1	54.8	57.8	58.0	67.2	70.5	62.8	56.0	67.2	66.2	59.0	66.3	60.5	60.9

Table 7: Unsupervised experimental results of the QR on McMarket. Where M and P denote MRR and Precision@3, respectively. \* denotes LLM-based methods.

Method	au		br		ca		cn		fr		in		jp		mx		uk		AVG		
	M	P	M	P	M	P	M	P	M	P	M	P	M	P	M	P	M	P	M	P	
Single	BERT-f	32.7	44.4	25.8	48.9	30.0	42.2	31.7	35.6	45.8	47.8	46.2	64.4	51.1	48.9	46.4	58.9	54.4	61.1	40.5	50.2
	T5	29.4	42.2	23.3	41.1	31.7	38.9	31.3	30.9	42.0	45.1	43.8	58.4	49.7	47.8	44.4	54.1	53.9	56.4	38.8	46.1
	monoT5	30.1	44.4	23.1	41.1	31.3	43.2	31.4	31.1	43.2	46.7	49.4	63.3	53.5	49.9	47.8	54.4	53.4	58.9	40.4	48.1
	Flan-T5*	39.7	51.1	26.9	50.0	34.0	46.7	38.3	42.2	52.2	54.4	51.4	63.3	54.8	64.4	49.3	60.0	55.8	62.2	44.7	54.9
Cross	Exact-BERT-f	46.4	45.6	40.0	51.1	51.5	47.8	49.4	45.6	52.3	53.2	49.3	66.0	53.4	47.8	48.9	63.3	58.7	66.7	50.0	54.1
	BERT-f	58.6	54.4	52.3	54.4	55.3	53.3	56.2	46.7	53.9	55.6	65.8	70.0	56.0	52.2	63.2	71.1	59.6	70.0	57.9	58.6
	Exact-monoT5	52.6	48.9	50.7	53.8	54.6	55.6	54.4	44.9	53.2	53.1	63.1	71.0	56.9	52.1	62.8	67.8	59.3	66.8	56.4	57.1
	monoT5	52.9	53.3	51.4	52.2	54.1	56.7	56.8	44.4	52.8	52.2	68.1	75.6	56.8	53.3	62.9	68.9	58.2	67.8	57.1	58.3
	Exact-Flan-T5*	60.8	60.3	55.7	56.9	61.3	59.2	57.6	55.2	58.1	57.8	67.2	73.3	57.1	54.3	63.9	74.9	63.0	73.9	60.5	62.9
	Flan-T5*	63.6	62.2	56.9	55.6	62.9	61.1	59.7	57.8	60.8	61.1	69.7	76.7	60.4	56.7	64.3	75.6	63.6	72.2	62.4	64.3

Table 8: Supervised experimental results of QR using McMarket<sub>q</sub>.

formance of supervised methods on this task, we split McMarket<sub>q</sub> into three sets, with 1260/180/360 samples in each. We then train each model on the training set and report results on the testing set.

**Evaluation metrics.** We adopt several evaluation metrics to assess the performance of models on two tasks. For AG, we compare the model-generated answers with ground-truth user answers using BLEU-4 (Papineni et al., 2002) and ROUGE-L (Lin, 2004) scores. For QR, we report major information retrieval (IR) metrics, namely, mean reciprocal rank (MRR) and Precision@3 to evaluate the ranking performance of different methods.

## 5.2 Compared methods

For AG, we first directly rank and select a review as the answer with methods such as BM25 (Robertson and Zaragoza, 2009), BERT (Devlin et al., 2019). Besides, several generative methods such as T5 (Raffel et al., 2019), LLaMA-2 (Touvron et al., 2023b), are leveraged to train the model to generate the answer given the question and reviews. Specifically, under the cross-market scenario, Exact-model means that in the auxiliary marketplace, we only use reviews from the same item before performing answer generation.

For QR, on McMarket, we report ranking methods that do not involve any training (*i.e.*, BERT,

UPR (Sachan et al., 2022)) or methods that perform unsupervised training (*i.e.*, CMJim (Ghasemi et al., 2023)). On McMarket<sub>q</sub>, we adopt supervised fine-tuning methods (*i.e.*, BERT-f/monoT5 (Nogueira et al., 2020)), and report testing performance. Details of each method are listed in Appendix D.

## 5.3 Experimental results

### 5.3.1 Review-based answer generation

Tables 5 and 6 show the single/cross-market answer generation performance on McMarket and McMarket<sub>r</sub> datasets. We have the following observations: first of all, cross-market models have superior overall performance in all marketplaces compared with methods in the single-market setting. This result verifies **RQ1** from the model perspective, showing that external resources (*i.e.*, reviews), from auxiliary marketplaces, can significantly contribute to improved outcomes in the main marketplace. A clear advantage of LLMs over traditional methods is evident across various marketplaces. Notably, LLaMA-2 outperforms the overall cross-market McMarket dataset, with a notable ROUGE improvement from 13.5 in T5 to 18.3. Similarly, in McMarket<sub>r</sub>, the overall ROUGE score sees significant enhancement, rising from 29.1 to 38.5. This provides an answer for **RQ3**, offering insights into the efficacy and potential advancements of LLMs.

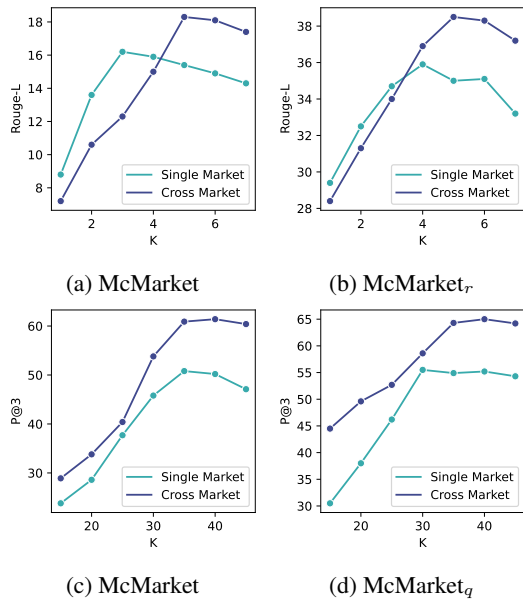


Figure 4:  $K$ -value analysis across different marketplaces on the best-performed model. The upper row is on AG, and the lower is QR.

### 5.3.2 Product-related question ranking

Tables 7 and 8 show the question ranking results within the single/cross-market scenario on two datasets. We notice that most observations from Section 5.3.1 still hold. For example, performance advantages persist in product-related question ranking compared to a single-market scenario. This shows that a large number of relevant questions in the auxiliary marketplaces help address similar questions in a low-resource marketplace. Furthermore, the performance boost is more obvious in marketplaces with a smaller scale (*i.e.*, **au**, **br**) compared with marketplaces with a larger scale (*i.e.*, **uk**). For instance, the P@3 BM25 performance exhibits an improvement 28.3 and 21.7 for **au** and **br** marketplaces, respectively, compared with 17.3 in **uk** on McMarket. We also find that in the cross-market setting, the Exact-models have a weaker overall performance than their original counterparts (*i.e.*, Exact-T5/Llama-2 v.s. T5/Llama-2). For example, on McMarket<sub>q</sub>, the cross-market Exact-Flan-T5 is 1.4 weaker in terms of overall P@3 compared with Flan-T5. This demonstrates that valuable information can be found within similar products from auxiliary marketplaces, even when they possess slightly different titles. We list some cases in Appendix E to elaborate on this.

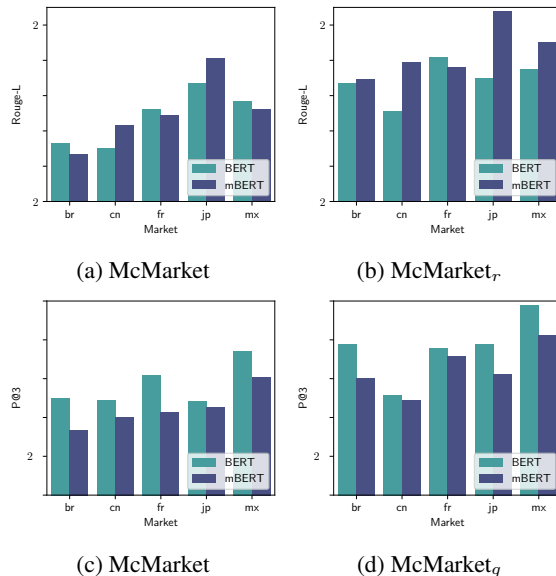


Figure 5: Multilingual analysis on non-English marketplaces. The upper row is on AG, the lower is QR.

## 6 External Analysis

### 6.1 Hyperparameter analysis

We investigate the effect of the number of retrieved product-related resources (*i.e.*, questions, reviews)  $K$  under both single/cross-market scenarios. We report the average performance among every marketplace on both McMarket and the corresponding subset in Figure 4. We observe that in AG, initially, the performance of Llama-2 in the cross-market setting is inferior to that in the single market. However, after increasing the value of  $K$ , the optimal  $K$  value in the cross-market scenario surpasses that in the single market. This tendency indicates that richer information is contained in the cross-market reviews. In QR, the ranking performance in the single-market scenario begins to decline when  $K$  is around 50. This indicates that some less relevant questions are retrieved, negatively impacting the results. Conversely, in the cross-market scenario, a greater number of relevant questions are accessible, helping to effectively mitigate this issue.

### 6.2 Multilingual analysis

We undertake a comparative analysis between translated and non-translated content to delve deeper into performance variations across non-English marketplaces. In particular, within the single-market scenario, we compare mBERT with BERT in 5 non-English marketplaces. Here, mBERT refers to a setup where all contents and the model itself are preserved and fine-tuned in their origi-



nal language without translation. The results are shown in Figure 5. We notice that in the AG task, concerning some non-Latin languages (*i.e.*, **cn**, **jp**), the performance of single-market mBERT without translation results in higher score compared with BERT on two datasets. However, we observe opposite results in some other non-English marketplaces (*i.e.*, **fr**). Besides, in the QR task, the performance of mBERT is inferior to the translated BERT model. This underscores a crucial future direction for this task: effectively enhancing performance in non-English marketplaces, an aspect that has been relatively underexplored.

## 7 Conclusions

We propose the task of Multilingual Cross-market Product-based Question Answering (MCPQA). We hypothesize that product-related information from a resource-rich marketplace can be leveraged to enhance the QA in a resource-scarce marketplace. To facilitate the research, we then propose a large-scale dataset, covering over 7 million questions across 17 marketplaces and 11 languages. Additionally, we perform automatic translation for the Electronics category, labeling it as McMarket. We also provide LLM-labeled subsets on McMarket for each of the two tasks, namely McMarket<sub>r</sub> and McMarket<sub>q</sub>. Specifically, we focus on two different tasks: AG and QR. We conduct experiments to compare the performance of models under single/cross-market scenarios on both datasets.

## Limitations

The task of PQA holds significant potential in improving user experiences on e-commerce platforms. However, there are several limitations and challenges associated. One major challenge is the quality and reliability of the information available for answering user questions. Even though we make sure all of the information comes from real user-generated data, the reviews and QA pairs might still contain biased or inaccurate information. Furthermore, language barriers and the availability of data in multiple languages add complexity to the task of product-related QA, particularly in cross-lingual scenarios. The limited availability of data in low-resource languages further exacerbates this challenge. To address them, continued research and development efforts are still under process which aim at improving data quality, handling language diversity, etc. We discuss it as our future work in

Appendix B.

## Ethics Statement

Our dataset is derived from the publicly available product question-answering dataset, XMarket (Bonab et al., 2021). We adhere to the policies throughout the creation and utilization of this dataset to ensure the protection of user privacy. When preparing the question-answering pairs, we strictly ensure that no personally identifiable information is exposed or utilized in any form during the processes. We prioritize user privacy and confidentiality to maintain the integrity and ethical standards of our dataset. We have licensed our data under CC0 1.0 DEED and will ask the users to sign an agreement such that the dataset will only be available for academic research purposes to further protect the users.

## References

- Zahra Abbasiantaeb, Yifei Yuan, E. Kanoulas, and Mohammad Aliannejadi. 2023. *Let the llms talk: Simulating human-to-human conversational qa via zero-shot llm-to-llm interactions*. *Proceedings of the 17th ACM International Conference on Web Search and Data Mining*.
- Mikel Artetxe, Sebastian Ruder, and Dani Yogatama. 2019. *On the cross-lingual transferability of monolingual representations*. In *Annual Meeting of the Association for Computational Linguistics*.
- Akari Asai, Jungo Kasai, J. Clark, Kenton Lee, Eunsol Choi, and Hannaneh Hajishirzi. 2020. *Xor qa: Cross-lingual open-retrieval question answering*. In *North American Chapter of the Association for Computational Linguistics*.
- Johannes Bjerva, Nikita Bhutani, Behzad Golshan, Wang Chiew Tan, and Isabelle Augenstein. 2020. *Subjqa: A dataset for subjectivity and review comprehension*. In *Conference on Empirical Methods in Natural Language Processing*.
- Hamed Bonab, Mohammad Aliannejadi, Ali Vardasbi, E. Kanoulas, and James Allan. 2021. *Cross-market product recommendation*. *Proceedings of the 30th ACM International Conference on Information & Knowledge Management*.
- Shiqian Chen, Chenliang Li, Feng Ji, Wei Zhou, and Haiqing Chen. 2019. *Review-driven answer generation for product-related questions in e-commerce*. In *Proceedings of the Twelfth ACM International Conference on Web Search and Data Mining, WSDM 2019, Melbourne, VIC, Australia, February 11-15, 2019*, pages 411–419. ACM.

- Hyung Won Chung, Le Hou, S. Longpre, Barret Zoph, Yi Tay, William Fedus, Eric Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, Albert Webson, Shixiang Shane Gu, Zhuyun Dai, Mirac Suzgun, Xinyun Chen, Aakanksha Chowdhery, Dasha Valter, Sharan Narang, Gaurav Mishra, Adams Wei Yu, Vincent Zhao, Yanping Huang, Andrew M. Dai, Hongkun Yu, Slav Petrov, Ed Huai hsin Chi, Jeff Dean, Jacob Devlin, Adam Roberts, Denny Zhou, Quoc V. Le, and Jason Wei. 2022. [Scaling instruction-finetuned language models](#). *ArXiv*, abs/2210.11416.
- J. Clark, Eunsol Choi, Michael Collins, Dan Garette, Tom Kwiatkowski, Vitaly Nikolaev, and Jenimaria Palomaki. 2020. [Tydi qa: A benchmark for information-seeking question answering in typologically diverse languages](#). *Transactions of the Association for Computational Linguistics*, 8:454–470.
- Yang Deng, Yaliang Li, Wenxuan Zhang, Bolin Ding, and Wai Lam. 2022. [Toward personalized answer generation in e-commerce via multi-perspective preference modeling](#). *ACM Trans. Inf. Syst.*, 40(4):87:1–87:28.
- Yang Deng, Ying Shen, Min Yang, Yaliang Li, Nan Du, Wei Fan, and Kai Lei. 2018. [Knowledge as A bridge: Improving cross-domain answer selection with external knowledge](#). In *Proceedings of the 27th International Conference on Computational Linguistics, COLING 2018, Santa Fe, New Mexico, USA, August 20-26, 2018*, pages 3295–3305. Association for Computational Linguistics.
- Yang Deng, Wenxuan Zhang, and Wai Lam. 2020. [Opinion-aware answer generation for review-driven question answering in e-commerce](#). In *CIKM '20: The 29th ACM International Conference on Information and Knowledge Management, Virtual Event, Ireland, October 19-23, 2020*, pages 255–264. ACM.
- Yang Deng, Wenxuan Zhang, Qian Yu, and Wai Lam. 2023. [Product question answering in e-commerce: A survey](#). *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics*.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [Bert: Pre-training of deep bidirectional transformers for language understanding](#). In *North American Chapter of the Association for Computational Linguistics*.
- Yue Feng, Zhaochun Ren, Weijie Zhao, Mingming Sun, and Ping Li. 2021. [Multi-type textual reasoning for product-aware answer generation](#). In *SIGIR '21: The 44th International ACM SIGIR Conference on Research and Development in Information Retrieval, Virtual Event, Canada, July 11-15, 2021*, pages 1135–1145. ACM.
- Shen Gao, Xiuying Chen, Zhaochun Ren, Dongyan Zhao, and Rui Yan. 2021. [Meaningful answer generation of e-commerce question-answering](#). *ACM Trans. Inf. Syst.*, 39(2):18:1–18:26.
- Shen Gao, Zhaochun Ren, Yihong Eric Zhao, Dongyan Zhao, Dawei Yin, and Rui Yan. 2019. [Product-aware answer generation in e-commerce question-answering](#). In *Proceedings of the Twelfth ACM International Conference on Web Search and Data Mining, WSDM 2019, Melbourne, VIC, Australia, February 11-15, 2019*, pages 429–437. ACM.
- Negin Ghasemi, Mohammad Aliannejadi, Hamed Bonab, E. Kanoulas, Arjen P. de Vries, James Allan, and Djoerd Hiemstra. 2023. [Cross-market product-related question answering](#). *Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval*.
- José Luis Vicedo González and Jaime Gómez. 2007. [Trec: Experiment and evaluation in information retrieval](#). *J. Assoc. Inf. Sci. Technol.*, 58:910–911.
- Mansi Gupta, Nitish Kulkarni, Raghuveer Chanda, Anirudha Rayasam, and Zachary Chase Lipton. 2019. [Amazonqa: A review-based question answering task](#). In *International Joint Conference on Artificial Intelligence*.
- Chin-Yew Lin. 2004. [Rouge: A package for automatic evaluation of summaries](#). In *Annual Meeting of the Association for Computational Linguistics*.
- Jiahua Liu, Yankai Lin, Zhiyuan Liu, and Maosong Sun. 2019. [Xqa: A cross-lingual open-domain question answering dataset](#). In *Annual Meeting of the Association for Computational Linguistics*.
- S. Longpre, Yi Lu, and Joachim Daiber. 2020. [Mkqa: A linguistically diverse benchmark for multilingual open domain question answering](#). *Transactions of the Association for Computational Linguistics*, 9:1389–1406.
- Julian McAuley and Alex Yang. 2015. [Addressing complex and subjective product-related queries with customer reviews](#). *Proceedings of the 25th International Conference on World Wide Web*.
- Rodrigo Nogueira, Zhiying Jiang, Ronak Pradeep, and Jimmy J. Lin. 2020. [Document ranking with a pre-trained sequence-to-sequence model](#). In *Findings*.
- OpenAI. 2023. [Gpt-4 technical report](#).
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. [Bleu: a method for automatic evaluation of machine translation](#). In *Annual Meeting of the Association for Computational Linguistics*.
- Yingqi Qu, Yuchen Ding, Jing Liu, Kai Liu, Ruiyang Ren, Xin Zhao, Daxiang Dong, Hua Wu, and Haifeng Wang. 2020. [Rocketqa: An optimized training approach to dense passage retrieval for open-domain question answering](#). In *North American Chapter of the Association for Computational Linguistics*.
- Colin Raffel, Noam M. Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2019. [Exploring the](#)

- limits of transfer learning with a unified text-to-text transformer. *J. Mach. Learn. Res.*, 21:140:1–140:67.
- Stephen E. Robertson and Hugo Zaragoza. 2009. [The probabilistic relevance framework: Bm25 and beyond](#). *Found. Trends Inf. Retr.*, 3:333–389.
- Ohad Rozen, David Carmel, Avihai Mejer, Vitaly Mirkis, and Yftah Ziser. 2021. [Answering product-questions by utilizing questions from other contextually similar products](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2021, Online, June 6-11, 2021*, pages 242–253. Association for Computational Linguistics.
- Devendra Singh Sachan, Mike Lewis, Mandar Joshi, Armen Aghajanyan, Wen tau Yih, Joëlle Pineau, and Luke Zettlemoyer. 2022. [Improving passage retrieval with zero-shot question generation](#). In *Conference on Empirical Methods in Natural Language Processing*.
- Xiaoyu Shen, Akari Asai, Bill Byrne, and A. Gispert. 2023. [xpqa: Cross-lingual product question answering in 12 languages](#). *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics*.
- Xiaoyu Shen, Gianni Barlacchi, Marco Del Tredici, Weiwei Cheng, and A. Gispert. 2022. [semipqa: A study on product question answering over semi-structured data](#). *Proceedings of The Fifth Workshop on e-Commerce and NLP (ECNLP 5)*.
- Nllb team, Marta Ruiz Costa-jussà, James Cross, Onur cCelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, Anna Sun, Skyler Wang, Guillaume Wenzek, Alison Youngblood, Bapi Akula, Loïc Barrault, Gabriel Mejia Gonzalez, Prangthip Hansanti, John Hoffman, Semarley Jarrett, Kaushik Ram Sadagopan, Dirk Rowe, Shannon L. Spruit, C. Tran, Pierre Yves Andrews, Necip Fazil Ayan, Shruti Bhosale, Sergey Edunov, Angela Fan, Cynthia Gao, Vedanuj Goswami, Francisco Guzm'an, Philipp Koehn, Alexandre Mourachko, Christophe Ropers, Safiyyah Saleem, Holger Schwenk, and Jeff Wang. 2022. [No language left behind: Scaling human-centered machine translation](#). *ArXiv*, abs/2207.04672.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023a. [Llama: Open and efficient foundation language models](#). *ArXiv*, abs/2302.13971.
- Hugo Touvron, Louis Martin, Kevin R. Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Daniel M. Bikel, Lukas Blecher, Cristian Cantón Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony S. Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel M. Kloumann, A. V. Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, R. Subramanian, Xia Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zhengxu Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023b. [Llama 2: Open foundation and fine-tuned chat models](#). *ArXiv*, abs/2307.09288.
- Mengting Wan and Julian McAuley. 2016. [Modeling ambiguity, subjectivity, and diverging viewpoints in opinion question answering systems](#). *2016 IEEE 16th International Conference on Data Mining (ICDM)*, pages 489–498.
- Hu Xu, Bing Liu, Lei Shu, and Philip S. Yu. 2019. [Bert post-training for review reading comprehension and aspect-based sentiment analysis](#). In *North American Chapter of the Association for Computational Linguistics*.
- Jianfei Yu, Minghui Qiu, Jing Jiang, Jun Huang, Shuangyong Song, Feng Ji, Wei Chu, and Haiqing Chen. 2017. [Modelling domain relationships for transfer learning on retrieval-based question answering systems in e-commerce](#). *Proceedings of the Eleventh ACM International Conference on Web Search and Data Mining*.
- Jianxing Yu, Zheng-Jun Zha, and Tat-Seng Chua. 2012. [Answering opinion questions on products by exploiting hierarchical organization of consumer reviews](#). In *Proceedings of the 2012 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning, EMNLP-CoNLL 2012, July 12-14, 2012, Jeju Island, Korea*, pages 391–401. ACL.
- Qian Yu and Wai Lam. 2018. [Review-aware answer prediction for product-related questions incorporating aspects](#). *Proceedings of the Eleventh ACM International Conference on Web Search and Data Mining*.
- Yifei Yuan and Wai Lam. 2021. [Conversational fashion image retrieval via multiturn natural language feedback](#). *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*.
- Rui Zhang, Caitlin Westerfield, Sungrok Shim, Garrett Bingham, Alexander R. Fabbri, Neha Verma, William Hu, and Dragomir R. Radev. 2019a. [Improving low-resource cross-lingual document retrieval by reranking with deep bilingual representations](#). *ArXiv*, abs/1906.03492.

Shiwei Zhang, Jey Han Lau, Xiuzhen Zhang, Jeffrey Chan, and Cécile Paris. 2019b. [Discovering relevant reviews for answering product-related queries](#). *2019 IEEE International Conference on Data Mining (ICDM)*, pages 1468–1473.

Shiwei Zhang, Xiuzhen Zhang, Jey Han Lau, Jeffrey Chan, and Cécile Paris. 2020a. [Less is more: Rejecting unreliable reviews for product question answering](#). In *Machine Learning and Knowledge Discovery in Databases - European Conference, ECML PKDD 2020, Ghent, Belgium, September 14-18, 2020, Proceedings, Part III*, volume 12459 of *Lecture Notes in Computer Science*, pages 567–583. Springer.

Wenxuan Zhang, Yang Deng, and Wai Lam. 2020b. [Answer ranking for product-related questions via multiple semantic relations modeling](#). *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval*.

Wenxuan Zhang, Yang Deng, Jing Ma, and Wai Lam. 2020c. [Answerfact: Fact checking in product question answering](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing, EMNLP 2020, Online, November 16-20, 2020*, pages 2407–2417. Association for Computational Linguistics.

Wenxuan Zhang, Wai Lam, Yang Deng, and Jing Ma. 2020d. [Review-guided helpful answer identification in e-commerce](#). In *WWW '20: The Web Conference 2020, Taipei, Taiwan, April 20-24, 2020*, pages 2620–2626. ACM / IW3C2.

## A LLM annotation details

We employ GPT-4 as the base LLM to perform automatic annotation. Specifically, gpt-4-1106-preview is adopted in our setting. For review-based answer generation, we pass the question, related reviews into the model, and ask GPT-4 to generate if the corresponding answer can be produced from the given information and write the answer if possible. We also instruct GPT-4 to provide the corresponding reason. We use the following prompt:

- *In this task, you will be given a product question, and some reviews. You should judge if the reviews are helpful for answering the question. If yes, please write the corresponding answer and the reason. If no, please give the corresponding reason and provide the answer as no answer. Please output the answer format as: Judgement:yes/no, Reason: , Answer:*

In our setup for product-related question ranking, we follow the annotation setting outlined in (Ghasemi et al., 2023). Here, we utilize GPT-4

to evaluate the relevance of other question-answer pairs. The model is presented with two question-answer pairs from distinct products along with their respective product titles. Its task is to assess whether the QA pair associated with the second product proves useful in addressing the questions posed for the first product. Similarly, the model is also requested to provide the reason for making the judgment. The prompt is given as follows:

- *In this task, you will be given two different products, namely, Product A and B, respectively. Each product is associated with a question-answer pair. You should judge if the question-answer pair to Product B is useful for answering the question to Product A. You should assign a score from 0–2, as 0 represents not useful, 1 represents partially useful, and 2 represents very useful. Please also give the corresponding reason for making the decision. Please output the answer format as: Judgement:[score], Reason:*

## B Future Directions

Future directions for the MCPQA task could involve several areas of exploration. First of all, more efforts could be put in the continued advancement and refinement of multilingual models capable of understanding and generating text across multiple languages. Furthermore, as a substantial portion of our dataset remains in its original, untranslated form, we are actively researching how models perform when fine-tuned on this untranslated data. Our focus lies particularly on assessing their question-answering performance in multilingual contexts. Based on that, investigation of cross-lingual transfer learning techniques to facilitate knowledge transfer and adaptation between languages could also be a promising direction in this task. This includes exploring approaches for transferring knowledge from high-resource to low-resource languages and vice versa.

## C Human evaluation metrics

- *Correctness* aims to judge whether GPT-4 answers accurately serve as correct answers to the question, based on the given information. For example, if the question is not answerable from the reviews, GPT-4 should make the corresponding judgment. Otherwise, GPT-4 should first classify the question as answerable, and then give the corresponding answer.

- *Completeness* is designed to determine whether the GPT-4 generated answers are complete and cover all aspects of the question.
- *Relevance* is designed to determine whether the GPT-4 answers are relevant to the question, and whether contain hallucination that does not correspond to the original question.
- *Naturalness* aims to determine whether the GPT-4 answers are smooth and natural. Whether there are obvious language errors and inconsistencies.

## D Baseline details

We provide a detailed explanation of the baseline models we implement.

**Review-based answer generation.** In this task, we report performance on McMarket and McMarket<sub>r</sub>. In contrast to utilizing human answers in McMarket, in McMarket<sub>r</sub>, we employ the GPT-4 generated results as the ground truth. For each dataset, we split the training/validation/testing set with the portion 70/10/20% and report the results on the testing set. The detailed information of each baseline is as follows:

- BM25 (Robertson and Zaragoza, 2009) retrieves the top 5 reviews and adopts the top one directly as the answer.
- BERT (Devlin et al., 2019) adopts a BERT ranker to re-rank the reviews retrieved by the top 100 BM25 results. Then the top 1 review is selected as the answer.
- T5 (Raffel et al., 2019) takes the BM25 top 5 reviews as input and is fine-tuned to generate the corresponding answer.
- Exact-T5 (Ghasemi et al., 2023) is an answer generation model based on T5, wherein we initially identify the exact same item in the auxiliary marketplace and exclusively utilize the top 5 reviews among them as input.
- LLaMA-2 (Touvron et al., 2023b) is in a similar setting as T5 but adopts LLaMA-2 as the backbone.
- Exact-LLaMA-2 is in a similar setting as Exact-T5 but adopts LLaMA-2 as the backbone.

**Product-related question ranking.** In this task, we also report results on McMarket and McMarket<sub>q</sub>. Given that the McMarket<sub>q</sub> subset is the only portion in McMarket that contains ranking labels, Table 7 exclusively showcases unsupervised methods that leverage the remaining McMarket as the training set and subsequently present results on the McMarket<sub>q</sub> subset. Besides, to show the performance of supervised methods in this task, Table 8 splits McMarket<sub>q</sub> as the training/validation/testing set following the same portion as before. Performance is then reported on the testing set.

We first provide details for the unsupervised methods in Table 7:

- BM25 (Robertson and Zaragoza, 2009) reports the top-50 BM25 ranking results.
- BERT (Devlin et al., 2019) performs BERT re-rank on BM25 top results.
- UPR-m (Sachan et al., 2022) is an unsupervised ranking method where we use a PLM to compute the probability of the input question conditioned on a related question. We use T5-base as the backbone.
- UPR-l adopts the same structure as UPR-m but uses T0-3B as the backbone.
- CMJim (Ghasemi et al., 2023) is an unsupervised method that ranks products and their corresponding questions across marketplaces.
- Exact-{BERT/UPR-l} ranks the questions of the item from the main marketplace as well as the exact same item in the auxiliary marketplace.

We then detail the supervised methods in Table 8:

- Bert-f (Devlin et al., 2019) fine-tunes the Bert ranker on the training set.
- T5 is trained to generate the sequence of the ranked questions.
- monoT5 (Nogueira et al., 2020) is another ranking method that takes T5 as backbone. We fine-tune the model on the training set and report the results on the testing portion.

Table 9: Real cases of the single/cross-market question answering results. The first two are review-based generated answers. The rest shows answers obtained by product-related question ranking. The information in the bracket shows the retrieved related review/question. All the information is shown in translated English.

Idx	Question	Single-market ans.	Cross-market ans.	Gold ans.
1	Which Dell Inspiron N4010 model can I buy and be compatible with?	(The express delivery was very fast. It arrived the next day.) I have no idea what to answer.	(I bought this to double my RAM on my new Dell Inspiron 7559.) The Dell Inspiron 7559 model is compatible with this product.	This model is low voltage...
2	Sir is it compatible with x-T3?	(The picture quality is very sharp. Very sharp indeed.) The given review information cannot provide the answer.	(Bokeh is awesome quality of photos is amazing paired with the xt3) Yes, the lens is compatible with xt3.	Yes
3	Will this work for a small fan that's 18 watts?	(Will this work for charging an electric razor?) It provides power out the same as a wall power point. So yes it should charge an electric razor.	(Will a travel fan work with this?) ... A travel fan should work this voltage converter. The converter has a maximum of 200 watts...	I think yes ...
4	Does it have an external microphone input?	(Does it have an audio and hdmi output?) Good morning. It has a micro HDMI output.	(Can you put an external microphone?) No microphone input	Doesn't have it!

Table 10: Examples of data samples in McMarket.

Market	Product title	Question	Reviews	Answer
br	Sony - HDRCX405 HD Video Recording Handycam Camcorder (black)	É compatível com eos 80d?	Objetiva com desempenho muito bom. Estabilização de imagem (IS) funciona muito bem para uso sem tripé. STM com foco silencioso. Cumpre o que promete.	Bom dia, é totalmente compatível.
cn	AKG Pro Audio K612 PRO Over-Ear, Open-Back, Premium Reference Studio Headphones	akg品控真有那么差吗还是一群职业黑?	一言难尽。买了十几天刚煲开右耳时响时不响。现在退货中	没有问题，还可以
fr	ViewSonic VG2439SMH 24 Inch 1080p Ergonomic Monitor with HDMI DisplayPort and VGA for Home and Office, Black	Sur écran webcam il y a t'il du son ? fait t'il webcam et micro en même temps?	Après réception; et déballage : produit simple et mise en marche facile. J'ai commandé deux écrans pour une station de travail. l'utilisateur est à l'aise	Pas le microphone. Webcam ok Son ok
jp	SanDisk Ultra 64GB USB 3.0 OTG Flash Drive With micro USB connector For Android Mobile Devices(SDDD2-064G-G46)	A1954に多用できますか	小さすぎて使いにくい (笑) 商品は、ゆうメールですすぐに配達されました。	A1954とは、何ですか? キーボードは、英語配列です。
mx	ZOTAC GeForce GT 730 1GB DDR3 PCI Express 2.0 x1 Graphics Card (ZT-71107-10L)	hola, es compatible con Lenovo TS-140?	Excelente producto y buen desempeño. Muy recomendable.	No conozco este equipo, solo se puede instalar en interfaces PIC x16.
uk	Peachtree Audio Deepblue2 High Performance Wireless Bluetooth Music System (Black)	Can you play music through this speaker as a wired device from an mp3 player.	Reluctant to pay so much but couldn't be more happy. Amazing sound quality and worth every penny. You will be blown away.	Yes, with the supplied cable plugged into the headphone jack on the MP3 player and the auxiliary input on the deepblue2.

- Flan-T5 (Chung et al., 2022) adopts the same structure as the monoT5 method but replaces the backbone to the Flan-T5-XL LLM.
- Exact- $\{BERT-f/monoT5/Flan-T5\}$  (Ghasemi et al., 2023) ranks the questions of the item from the main marketplace as well as the exact same item in the auxiliary marketplace.

## E Case study

Table 9 demonstrates four real cases concerning single/cross-market question answering. We see that the absence of useful information, such as related reviews or questions, within a single marketplace leads to inaccurate answers. For instance, in case 1, the retrieved reviews fail to provide suffi-

Market	Question	Reviews	GPT-4 ans.	Gold Ans.	GPT-4 reason
au	Can you use this with Xbox one?	... These however do not work as an Xbox one gamechat since there is no mic ...	No, you cannot use these headphones with Xbox One as there is no mic	Yes.	Review 1 mentions that these headphones do not work as an Xbox One gamechat as there is no mic.
ca	Will it work with the Nikon D3200	I have the lens fixed on to my Nikon D3200 almost permanently...	Yes, it will work with the Nikon D3200.	Works on my d3100.	The reviews mention that the product is compatible with the Nikon D3200 and they are using it effectively.
in	Does it has diopter correction ring?	Extremely well built. Great lens and use it mostly for bird watching during my hikes...	No answer.	Yes.	None of the reviews mention a diopter correction ring for the product.
jp	Does it come with a conversion plug to connect to an amplifier?	...They don't come with an adapter that can adjust to the standard plug, so be careful.	No, it does not come with a conversion plug to connect to an amplifier.	N/A.	Review 1 states explicitly that the headphones do not come with an adapter to adjust to the standard plug...
mx	How quiet is it?	...Good performance, very quiet; Nice and quiet easily installed and works well...	The product is very quiet.	It is very quiet, reliable, highly recommended	Reviews 2, 3, 4, and 5 directly address the noise level of the product by stating it is 'very quiet' and 'nice and quiet'

Table 11: Examples of data samples in McMarket<sub>r</sub>. All the data is translated into English.

Market	Product A	Product A QA	Product B	Product B QA	tag	GPT-4 reason
au	Newer 48 Macro LED Ring Flash Bundle with LCD Display Power Control...	Will this work with fuji x-t3 and x-t20? -> As long as they have a hot shoe, it will work. There is several lens ring adaptors for various lens sizes (talking about changeable lenses of course).	Newer 48 Macro LED Ring Flash Bundle with LCD Display Power Control...	Is this compatible with FujifilmX-T3? -> As long as you have a hotshoe it should work.	2	Both Product A and Product B are the same Newer 48 Macro LED Ring Flash Bundle, and the questions for both are concerning the compatibility with Fujifilm X-T3...
cn	Kingston Digital Multi-Kit/Mobility Kit 16 GB ...	Hello, what is the writing speed of this micro sdxc? -> Write: 14Mo/s   Read: 20Mo/s ...	Kingston Digital Multi-Kit/Mobility Kit 16 GB...	Speed of the card? -> Class 4 IE 4MB/sec.	1	The answer to Product B provides the class rating of a microSDHC card, though different from Product A...
fr	iPad Air New iPad 9.7 inch 2017 Case...	Good evening, is this case compatible with an iPad 2? Thank you -> Yes, no problem.	iPad Air New iPad 9.7 inch 2017 Case...	Does this case fit the ipad air 2? -> Hi, This case is not compatible with the iPad Air 2.	0	Product A is asking about iPad 2, while Product B is about compatibility with an iPad Air 2...
in	AmazonBasics USB 2.0 ...	Is it compatible with Nintendo switch? -> Dono but working good nice product.	AmazonBasics USB 2.0 ...	Is this compatible with MacOS? -> Yes.	0	The answer to Product B's question does not provide information for A...
uk	HDMI Media Player, Black Mini 1080p Full-HD Ultra...	Is it possible to power this through a usb cable? -> It has to be plugged in using the power lead...	MDN HD1080B 1080p Full-HD Ultra Portable Digital Media Player...	Can it be powered by a USB cable? I see on the pictures that power cable is USB on one end -> The USB port is for an external drive.	2	The question for both Product A and Product B pertains to the power source of the media players and whether they can be powered through a USB cable...

Table 12: Examples of data samples in McMarket<sub>q</sub>. All the data is translated into English.

cient information, resulting in a generated answer of “I have no idea what to answer.” In contrast, relevant and useful information is more likely to be available in the larger auxiliary marketplace. For

instance, in case 4, the model successfully retrieves a similar question, “Can you put an external microphone?” from the **us** marketplace, aligning the answer more closely with the ground-truth answer.

## F Dataset Examples

We show some examples from McMarket to provide a more comprehensive view of our data. Table 10 shows some examples from McMarket. For each example, we show the title of a product, a random review, and a question-answer pair of the product.

To provide a more comprehensive understanding of our dataset and task, we also show some examples of the GPT-4 annotated McMarket<sub>r</sub> (Table 11) and McMarket<sub>q</sub> (Table 12), respectively.