

MPMQA: Multimodal Question Answering on Product Manuals

Liang Zhang¹, Anwen Hu¹, Jing Zhang², Shuo Hu², Qin Jin^{1*}

¹School of Information, Renmin University of China

²Samsung Research China - Beijing (SRC-B)

{zhangliang00,anwenhu,qjin}@ruc.edu.cn, {jing97.zhang,shuo.hu}@samsung.com

Abstract

Visual contents, such as illustrations and images, play a big role in product manual understanding. Existing Product Manual Question Answering (PMQA) datasets tend to ignore visual contents and only retain textual parts. In this work, to emphasize the importance of multimodal contents, we propose a Multimodal Product Manual Question Answering (MPMQA) task. For each question, MPMQA requires the model not only to process multimodal contents but also to provide multimodal answers. To support MPMQA, a large-scale dataset PM209 is constructed with human annotations, which contains 209 product manuals from 27 well-known consumer electronic brands. Human annotations include 6 types of semantic regions for manual contents and 22,021 pairs of question and answer. Especially, each answer consists of a textual sentence and related visual regions from manuals. Taking into account the length of product manuals and the fact that a question is always related to a small number of pages, MPMQA can be naturally split into two subtasks: retrieving most related pages and then generating multimodal answers. We further propose a unified model that can perform these two subtasks all together and achieve comparable performance with multiple task-specific models. The PM209 dataset is available at <https://github.com/AIM3-RUC/MPMQA>.




Introduction

Product manuals contain detailed descriptions of product features and operating instructions. They are often so long that it is not easy for users to efficiently find the information they are looking for. Therefore, Product Manual Question Answering (PMQA) (Nandy et al. 2021; Castelli et al. 2020) aims to build an AI agent on product manuals to conveniently answer user questions. PMQA leverages textual information in the manual, but ignores the visual contents, such as illustrations, tables and images, which are also important for solving user problems. As shown in Figure 1, the textual contents are insufficient to answer the question. In contrast, a multimodal answer containing both textual and visual contents can answer the question more clearly and precisely, from which users can grasp answers more effectively and efficiently. Existing Multimodal Question Answering tasks are designed to answer questions from a single

*Corresponding Author.

Gesture Mode

In Gesture Mode, the Mavic's Vision System recognizes gestures, allowing it to follow and capture selfies without a phone or a controller. Follow the steps below to use Gesture Mode:

Modes	Prompts	Front LEDs	Remarks
1. Confirm the subject		⋮..... Slow Red Flashing	Ensure the forward vision system is active and there is enough light. Tap the icon and move in front of the camera for the Mavic to recognize you.
2. Confirm the distance		⋮: x2..... Red Flashes Twice	Raise your arms and wave to the Mavic, the Front LED will blink red twice once it confirms the shooting distance.
3. Selfie Count Down		⋮:..... Fast Red Flashing	Put your fingers in front on your face as shown.

Question: How to take a selfie in the Selfie Count Down mode?

Textual-part Answer: You should put your fingers in front on your face as shown in the figure.

Visual-part Answer:



Figure 1: In this case, the gesture of the 'Selfie Count Down mode' is hard to describe using only plain text, but can be easily delivered with an image. In the MPMQA task, each question is answered with multimodal content: a textual-part answer and a visual-part answer.

web page (Chen et al. 2021; Tanaka, Nishida, and Yoshida 2021) or an infographic (Mathew et al. 2022), which are not suitable for product manual question answering, because product manuals always contain multiple pages and most of the pages are irrelevant to the question. Therefore, to fill the research gap in this area, we propose a challenging task namely **Multimodal Product Manual Question Answering (MPMQA)**. It requires the model to comprehend both the visual and the textual contents in an entire product manual and provide a multimodal answer for a given question.

We construct a large-scale dataset named PM209 with human annotations to support the research on the MPMQA task. It contains 22,021 QA annotations over 209 product manuals in 27 well-known consumer electronic brands. To support understanding of the multimodal content, we classify manual content into 6 categories (Text, Title, Product image, Table, Illustration, and Graphic). Each question is associated with a multimodal answer which is comprised of two parts: a textual part in natural language sentences, and a visual part containing regions from the manual. Table 1

shows the basic comparison between PM209 and existing PMQA datasets (Nandy et al. 2021). The scale of PM209 is larger than existing PMQA datasets in terms of brands, manual numbers, and QA pairs.

Considering most pages are irrelevant to a given question, it is natural to split the MPMQA task into two subtasks: firstly retrieving the most relevant pages and then generating answers with detailed information. Thus a straightforward solution for MPMQA is to apply two task-specific models. However, both page retrieval and answer generation require the model to correlate multimodal manual contents with the question. It is possible to have both subtasks benefit from each other. Therefore, we propose the **Unified Retrieval and Question Answering (URA)** model that performs these two steps with shared multimodal understanding ability. Specifically, URA uses a shared encoder to encode the multimodal page in the retrieval and question-answering tasks. Based on multitask learning, the URA model achieves comparable performance with multiple task-specific models.

Our contributions are summarized as follows:

- We propose the novel MPMQA task, which requires the model to understand multimodal content in the product manual, and answer questions with multimodal outputs.
- We construct a large-scale dataset PM209 to support MPMQA. It contains not only semantic labels for manual contents, but also multimodal answers for questions.
- For the MPMQA task, we design a unified model named URA that can both retrieve relevant pages and generate multimodal answers. It achieves comparable results with multiple task-specific models.

Related Works

Product Manuals Question Answering

To build an automatic question-answering system, existing works explore constructing datasets based on product manuals. TechQA (Castelli et al. 2020) collects 1400 user questions from the online forums and annotates the corresponding answers from IBM technical documents. For each question, TechQA annotates a single text span answer in the documents, similar to the strategy in SQuAD (Rajpurkar et al. 2016; Rajpurkar, Jia, and Liang 2018). Nandy et al. (Nandy et al. 2021) propose S10 QA and Smart TV/Remote QA datasets. They extract multiple text spans from the two Samsung device manuals to answer each question. These works leverage textual contents in manuals to build automatic QA systems, but ignore crucial vision information. In this work, we propose the MPMQA task, which requires models to understand both text and vision information to generate multimodal answers. Besides, our dataset PM209 is much bigger than the aforementioned datasets in terms of the number of products and the number of question-answering pairs.

Multimodal Question Answering

Many efforts have been made to answer questions from a multimodal context. TextVQA (Singh et al. 2019), ST-VQA (Biten et al. 2019), and EST-VQA (Wang et al. 2020) explore question answering on the image with scene

texts. They typically require the model to extract correct scene text according to the question. ManyModalQA (Hannan, Jain, and Bansal 2020) and MultiModalQA (Talmor et al. 2020) reason across text, tables and images from Wikipedia. DocVQA (Mathew, Karatzas, and Jawahar 2021) performs question answering on industry documents. VisualMRC (Tanaka, Nishida, and Yoshida 2021), WebSRC (Chen et al. 2021), WebQA (Chang et al. 2022) and DuReader_{vis} (Qi et al. 2022) require comprehension on web pages. InfographicVQA (Mathew et al. 2022) focuses on arithmetic reasoning over infographics. Different from previous multimodal inputs, the product manual is a specific domain in terms of the question type and the content. Since product manuals usually contain detailed operation instructions for a specific device, the questions beginning with 'How to' are very common (Nandy et al. 2021), while this type of contents and questions rarely occur in general domain datasets. Moreover, the answers in the above-mentioned works are all in text format, including text span, multi-choice, and generative sentences. Multimodal answers are less studied in the existing literature. MIMOQA (Singh et al. 2021) explores incorporating a Wikipedia-sourced image as a part of the answer. Apart from the domain difference, the setting in MIMOQA is rather ideal, as it assumes all text answers associated to at least one complementary image. This assumption does not hold in product manuals. The visual-part answer in MPMQA is very diverse, not restricted to images. It can also be regions like titles and tables. Moreover, most aforementioned works search for answers within a single document or web page. However, in the real scenario of PMQA, the target pages are not given in advance, and models have to locate relevant regions by themselves from an entire manual. To better fit the real application scenarios, our MPMQA task is designed to answer a question according to a complete manual rather than a single page, which is much more challenging than previous works.

MPMQA Task and PM209 Dataset

This section first presents a formal definition of the MPMQA task, and then describes the detailed process of constructing the PM209 dataset.

MPMQA Task Definition

TASK (MPMQA). Given a question Q and an n -page product manual $M = \{P_i\}_1^n$, where $P_i = \{r_{i1}, \dots, r_{ik}\}$ refers to a page in M and r_{ij} represents a semantic region in P_i , the model produces a multimodal answer $A = (T, R)$ containing two parts, with T as the textual-part answer in natural language sentences and $R = \{r_i\}_1^m$ as the visual-part answer consisting of multiple semantic regions.

Since almost all questions are relevant to a very small number of pages in a manual, the MPMQA task can be naturally split into the following two subtasks:

SUBTASK I (Page Retrieval). Given a question Q and an n -page product manual $M = \{P_i\}_1^n$, the model finds the smallest subset $\{P_{(i)}\}_1^k$ that contains the answer of Q .

Dataset	# Manuals	# Brands	# QA pairs	Multimodal content	Answer type
S10 QA	1	1	904	✗	Extractive
Smart TV/Remote QA	1	1	950	✗	Extractive
PM209 (ours)	209	27	22,021	✓	Multimodal

Table 1: Comparison between other question answering datasets on product manuals.

SUBTASK II (Multimodal QA). Given a question Q and k relevant pages $\{P_{(i)}\}_1^k$, the model generates a multimodal answer $A = (T, R)$ as defined in the TASK MPMQA.

PM209 Dataset Construction

We construct the PM209 dataset to support the MPMQA task. We first collect a set of product manuals. Crowd workers from the Maadaa Platform¹ then annotate the semantic regions r for each page P in the manuals. After that, the OCR words $W = \{w_i, b_i\}_1^n$ inside each semantic region r are automatically extracted. Finally, crowd workers create (question, multimodal answer) pairs (Q, A) based on the content of each manual. All crowd workers who participated in this project are proficient English speakers.

Product Manual Collection. We collect 209 English product manuals in total from well-known consumer electronic brands. These manuals cover 27 brands and 90 categories.

To ensure the diversity, we only keep the longest manual for the products in the same series. All manuals are born-in-digital PDF files and we render each page into image. We manually remove pages that are not suitable for posing questions, such as empty pages and cover pages, and ensure that all manuals in PM209 contain not less than 10 valid pages.

Semantic Region Annotation. Thirteen crowd workers are recruited to annotate the semantic regions r_i of each page in the product manuals. Two crowd workers then further validate the annotations. A semantic region consists of a bounding box b_i and a semantic label c_i . We define six semantic regions as follows.

- **Text.** The body paragraphs that convey major textual information in the product manual.
- **Title.** The words summarize or indicate the section of the whole page or nearby paragraph. Titles typically consist of a few words and have different fonts than the words in the paragraph (e.g. larger size, in bold or different color).
- **Product image.** Product relevant images in the manual, including the picture of product, operating interface, and components of the product etc. Product irrelevant images such as decorative drawings are not included.
- **Illustration.** Visually rich regions to describe a particular function, operation, and purpose of the product. They usually but not always consist of a combination of a product image and a surrounding text notes.
- **Table.** Regions that convey the information of text in a row-column format.
- **Graphic.** Visually rich regions indicating the name and position of a product component. It typically consists of

a product image, some surrounding texts, and indicators (lines, arrows, and serial numbers) that align the names in the text regions with positions in the product image.

To reduce the burden of human annotation, we leverage PyMuPDF (McKie and Liu 2016) to automatically extract bounding boxes of paragraphs and images in each page. We attach the 'Text' and 'Product image' labels to the paragraph and image bounding boxes produced by PyMuPDF respectively. The crowd workers then modify these initial bounding boxes and generate the above-mentioned semantic regions. The modification options include moving, resizing, relabeling categories, creating, and deleting.

Word Extraction. Since the product manuals are born-in-digital, we automatically extract OCR words $\{w_i, b_i\}_1^n$ in each region through PyMuPDF (McKie and Liu 2016).

QA Annotation. Twenty crowd workers are recruited to create (question, multimodal answer) pairs (Q, A) for each product manual. Considering the large cognitive load for reading the entire manual, and the fact that a question is usually only relevant to a few pages, we divide the entire product manuals into groups, and each group contains consecutive 5 pages. Crowd workers focus on each group and pose two questions for each page. For each question, they create a multimodal answer, which consists of two parts: the textual part T that is written to describe the answer in natural language sentences, and the visual part $R = \{r_i\}_1^m$ that is selected from the semantic regions. To simulate the real user scenario, the annotators are encouraged to write the question in the first person, and provide textual part answer T in the second person.

Statistics and Analysis

This section presents the statistics and analysis of the proposed PM209 dataset.

Manuals. PM209 consists of 209 product manuals in 27 well-known consumer electronic brands and 90 product categories. Figure 2 shows the top 10 products and brands. Note that the top 10 products cover less than 50% of all manuals, which reveals that the manuals in PM209 are highly diverse.

We also analyze the distribution over the number of pages in Figure 3a. It shows that PM209 are also diverse in lengths, ranging from 10 pages to 500 pages. The average length of the manuals is 50.76 pages.

Semantic regions. Figure 3b presents the statistics of the semantic regions. We observe that product manuals indeed contain rich layout information. Specifically, 65.1% of pages contain visually-rich regions such as product images, illustrations, tables and graphics. And 22.1% of these regions occur in the visual-part answer.

¹<https://maadaa.ai>

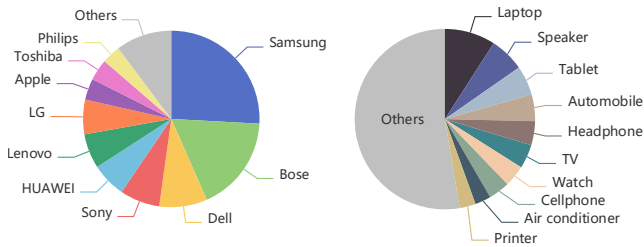
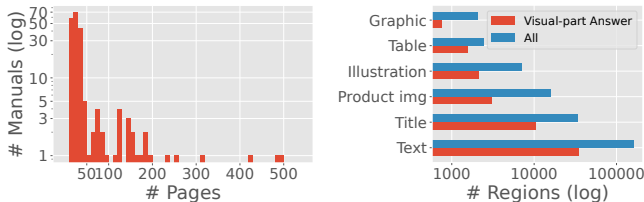


Figure 2: Top 10 brands (left) and products (right) in PM209.



(a) Count of manuals with a particular number of pages. (b) The number of semantic regions.

Figure 3: Statistics over pages and regions.

Questions and Textual-part Answers. The comparison between PM209 and other Multimodal Question Answering datasets is shown in Table 2. PM209 has a higher percentage of unique questions (98.46%) and unique answers (98.35%). It further indicates the high diversity of the PM209 dataset, since we avoid the appearance of similar product manuals, and both the questions and the answers in PM209 are specifically designed for each product. In addition, PM209 has the longest answer compared to other datasets, since the instruction and procedural answers can be long in product manuals.

Figure 4 shows the word cloud of the questions and textual-part answers. We find that questions in PM209 contain both factual words such as 'function' and 'information', and procedural words including 'begin', 'step', and 'after'. Apart from guidance-related questions such as 'what' and 'how', the frequency of pronoun 'I' has a high frequency in the questions. Correspondingly, the word 'you' appear frequently in the answers. This is as expected since we simulate the real-world scenarios where users pose questions in the first person, while the QA system answers the questions in the second person. Figure 5 shows the first 4-grams of questions and answers. Most questions begin with the word 'what' (43.91%) and 'how' (25.72%). Questions with 'how' tend to ask about the procedural process of an operation. Questions with 'what' are typically about factual information about the product usage, except in the case of 'What should I do ...', which are also procedural questions. Besides, there are 7.71% of questions starting with the word 'can'. These questions are usually confirming something uncertain about the product, e.g. 'Can I use this device underwater?'. Their answers usually begin with 'yes' or 'no'.

Visual-part Answers. Apart from the textual-part answers, each question in PM209 is also paired with a set of regions in the product manual. These regions can be seen as comple-

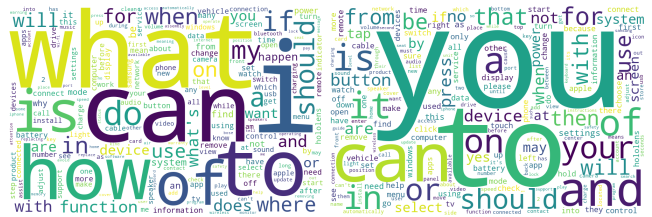


Figure 4: Word clouds for questions (left) and textual-part answers (right) in PM209.

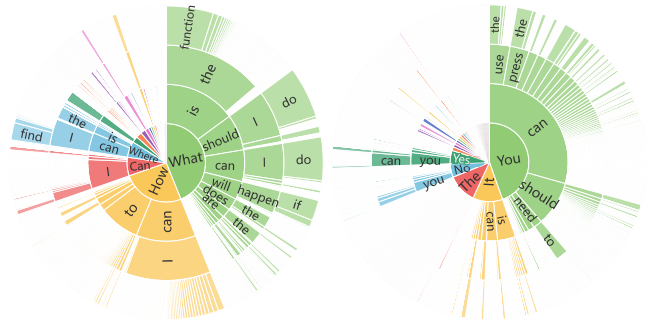


Figure 5: First 4-grams of questions (left) and answers (right) in PM209.

Dataset	Question		Answers		Page Length
	%Uniq.	Length	%Uniq.	Length	
ST-VQA	84.84	8.80	65.63	1.56	7.52
TextVQA	80.36	8.12	51.74	1.51	12.17
DocVQA	72.34	9.49	64.29	2.43	182.75
VisualMRC	96.26	10.55	91.82	9.55	151.46
InfographicVQA	99.11	11.54	48.84	1.60	217.89
PM209	98.46	9.77	98.35	15.74	231.36

Table 2: Comparison of Multimodal Question Answering datasets w.r.t. uniqueness rate and the average length of questions and answers, and the average length per page.

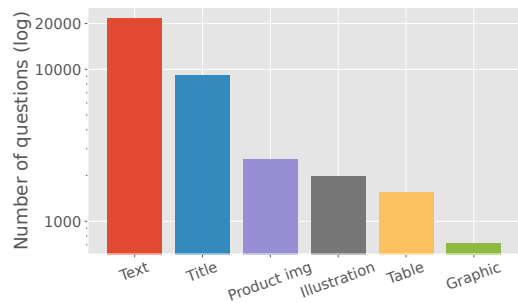


Figure 6: Visual-part answers break down by semantic labels

mentary to understanding the text answers. Figure 6 shows the number of visual-part answers broken down by semantic labels. A significant portion (21.8%) of questions include visually-rich regions (product images, illustrations, tables, and graphics) in their visual-part answers. This portion is

	Train	Val	Test
# Manuals	146	21	42
# Pages	7004	1011	2003
# QAs	15839	2257	3925

Table 3: Number of samples in each data split.

higher than VisualMRC, in which 9.1% of questions are relevant with visually-rich regions (picture and data). It indicates that visual components can be more important in understanding product manuals than open-domain web pages.

Data splits. We divide the manuals in the PM209 dataset into Train/Val/Test as shown in Table 3.

Proposed Model

We propose a **Unified Retrieval and Question Answering (URA)** model for the new MPMQA task, which can perform page retrieval and multimodal QA all together. As shown in Figure 7, the model consists of three key components: a URA Encoder, a URA Decoder, and a Region Selector. For the page retrieval task, URA encodes the questions and the pages separately, and calculates their relevant scores with token-level interaction. For the multimodal question answering, URA encodes questions and pages jointly, and produces the textual part and visual part of the multimodal answer through the Decoder and Region Selector.

Input Embeddings

URA embeds questions and pages similar to LayoutT5 (Tanaka, Nishida, and Yoshida 2021).

Question tokens. Question Q is tokenized into subword units with SentencePiece (Kudo and Richardson 2018). The special token $\langle /s \rangle$ denotes the end of the question.

$$x_Q^{\text{token}} = [q_1, q_2, \dots, q_m, \langle /s \rangle] \quad (1)$$

Region tokens. The region tokens consist of a special token $\langle c_i \rangle$ followed by the OCR words in this region. $\langle c_i \rangle$ denotes the semantic label of r_i .

$$x_{r_i}^{\text{token}} = [\langle c_i \rangle, w_{i1}, \dots, w_{ik}] \quad (2)$$

Page tokens. The sequence of page tokens is the concatenation of all region tokens in the page:

$$x_P^{\text{token}} = [x_{r_1}^{\text{token}}, \dots, x_{r_n}^{\text{token}}] \quad (3)$$

Special embeddings. Apart from the token embeddings, we add segment embedding z^{seg} to distinguish question/page tokens. To incorporate visual and layout information, we add 2D positional embeddings z^{pos} (Xu et al. 2020) and ROI embeddings z^{roi} (Anderson et al. 2018) to each page token.

$$z_Q = z_Q^{\text{token}} + z_Q^{\text{seg}} \quad (4)$$

$$z_P = z_P^{\text{token}} + z_P^{\text{seg}} + z^{\text{pos}} + z^{\text{roi}} \quad (5)$$

Page Retrieval

Page Retrieval aims to find the relevant pages for a question, which requires producing relevant scores between the question and pages. Our URA encoder f processes Q and P separately.

$$h_Q = f(z_Q; \theta_f) \quad (6)$$

$$h_P = f(z_P; \theta_f) \quad (7)$$

Since the clues to answer a question usually only appear in a small part of the page, considering the large content of the page, it is difficult for a single global feature to retain detailed clues. Thus, different from general retrieval methods that calculate the cosine similarity between global features, we perform token-level interaction (Yao et al. 2021) between Q and P as shown in Figure 7(a). Specifically, We calculate the token-level relevant scores s^{ij} between each token in h_Q^i and h_P^j , and aggregate them into two global relevant scores: question-to-page relevant score $S_{Q \rightarrow P}$ and page-to-question relevant score $S_{P \rightarrow Q}$:

$$s^{ij} = \|h_Q^i\|^T \|h_P^j\| \quad (8)$$

$$S_{Q \rightarrow P} = \frac{1}{N} \sum_i \max_j (s^{ij}) \quad (9)$$

$$S_{P \rightarrow Q} = \frac{1}{M} \sum_j \max_i (s^{ij}) \quad (10)$$

We optimize the model by minimizing the NCE loss (Gutmann and Hyvärinen 2010) on both the $Q \rightarrow P$ and $P \rightarrow Q$ directions. The loss function for Page Retrieval is written as:

$$\mathcal{L}'(\cdot) = \frac{1}{B} \sum_i \log \frac{\exp(S_{(\cdot)}^i / \tau)}{\sum_j \exp(S_{(\cdot)}^{ij} / \tau)} \quad (11)$$

$$\mathcal{L}_{\text{PR}} = \frac{1}{2} (\mathcal{L}'_{Q \rightarrow P} + \mathcal{L}'_{P \rightarrow Q}) \quad (12)$$

Where $\tau = 0.01$ denotes the temperature parameter of NCE, and B denotes the batch size. Note that since we focus on retrieving pages relevant to a given question during inference, we use the score $S_{Q \rightarrow P}$ to rank the candidate pages.

Multimodal Question Answering

Compared to finding relevant pages for a question, answering a question requires a stronger understanding of both the question and the multimodal contents of the page. Thus, different from Page Retrieval, URA encodes question Q and page P jointly to perform early interaction for Multimodal QA. We get the joint hidden state H as follows:

$$H = f([z_Q, z_P]; \theta_f) \quad (13)$$

Textual-part Answer. As shown in Figure 7(b), the URA decoder receives H and generates the textual-part of the multimodal answer auto-regressively. We train the model in a teacher-forcing manner by minimizing negative log-likelihood loss as below:

$$\mathcal{L}_{\text{TA}} = -\frac{1}{N} \log p(Y|Y_{<}, H) \quad (14)$$

Where $Y = [y_1, \dots, y_N]$ are the ground truth tokens.

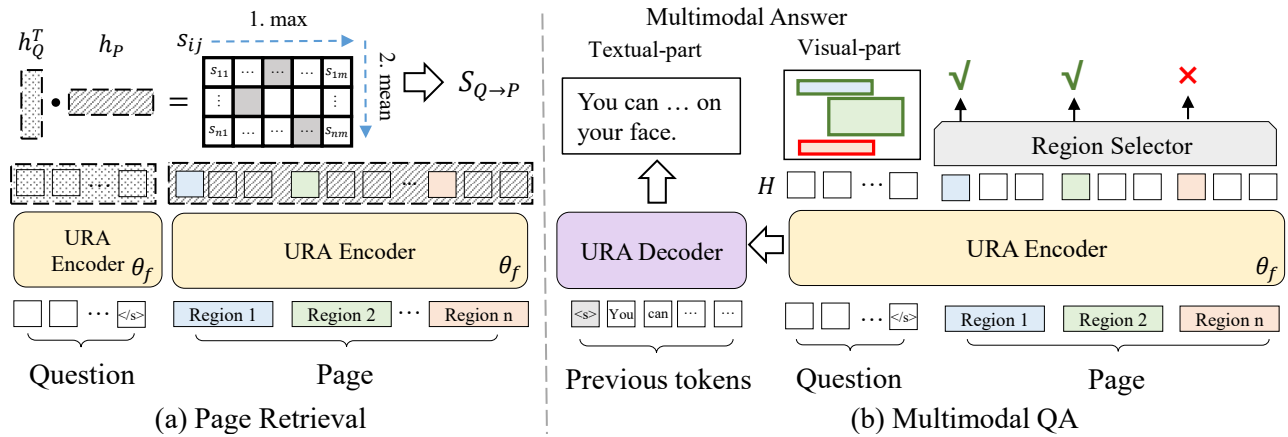


Figure 7: Overview of the Unified Retrieval and Question Answering (URA) model.

Visual-part Answer. The Region Selector RS selects a set of regions to compose the visual-part of the multimodal answer. RS is implemented as a linear projection layer followed by a sigmoid activation. The encoder hidden states corresponding to the $\langle c_i \rangle$ token is chosen to decide whether r_i is relevant to the question. We minimize the BCE loss to train the model as follows:

$$p_i = \text{RS}(H_{\langle c_i \rangle}; \theta_{\text{RS}}) \quad (15)$$

$$\mathcal{L}_{\text{VA}} = -\frac{1}{N} \sum_i y_i \log(p_i) + (1 - y_i) \log(1 - p_i) \quad (16)$$

Where $y_i = \{0, 1\}$ denotes whether the region r_i belongs to the ground truth vision-part answer.

Multitask Learning

Finally, URA is optimized in a multitask learning manner, where the final loss function is calculated as follows:

$$\mathcal{L}_{\text{URA}} = \mathcal{L}_{\text{PR}} + \mathcal{L}_{\text{TA}} + \mathcal{L}_{\text{VA}} \quad (17)$$

Experiments

We conduct experiments to validate our URA model on the proposed PM209 dataset.

Evaluation Setup

Evaluation settings. As mentioned before, MPMQA can be naturally split into two subtasks. Thus, we design two evaluation settings for subtask II, Multimodal QA: 1) *separate setting*: evaluating QA given the ground-truth pages; 2) *cascade setting*: evaluating QA given retrieved pages. We adopt the *separate setting* by default if not specified.

Evaluation metrics. For Page Retrieval, we calculate Recall@{1,3,5} in each manual, and weigh the scores across manuals by the number of pages. For Textual-part Answer, we report sequence generation metrics BLEU4 (B4) (Papineni et al. 2002), METEOR (M) (Banerjee and Lavie 2005), ROUGE-L (R-L) (Lin 2004) and CIDEr (Vedantam, Lawrence Zitnick, and Parikh 2015). For Visual-part Answer, we report the average Precision (P), Recall (R), and F1 scores on the whole dataset.

Baselines

We compare our URA model with the following baselines:

- PR: the Page Retrieval task-specific model. It can conduct the Page Retrieval task only.
- PR_g: the Page Retrieval task-specific model that uses global features to measure the relevancy between questions and pages.
- PR+TA: the multi-task model that is jointly optimized with Page Retrieval and Textual-part Answer tasks.
- PR_g+TA: the multi-task model that is jointly optimized with the Page Retrieval (global features) and Textual-part Answer tasks.
- 3 Single: 3 separate task-specific models for Page Retrieval, Textual-part Answer, and Visual-part Answer.

Implementation Details

We implement the above-mentioned models based on Pytorch (Paszke et al. 2019) and Huggingface Transformers (Wolf et al. 2020). The encoder and decoder of the models are in standard transformer architecture (Vaswani et al. 2017) with T5 (Raffel et al. 2020) initialization. The models adopt the T5_{BASE} structure that consists of 12 transformer layers with 768-d hidden states. We train the models for 20 epochs with a batch size of 8 and a learning rate of 3e-5. It takes about 20 hours to converge on 1 NVIDIA RTX A6000 GPU. We choose the model that performs best on the validation set, and report its performance on the test set. We consider the most relevant page for the Multimodal QA task.

Results and Analysis

Comparing URA with several baselines. Table 4 shows the comparison between URA and the baselines described above. Comparing row 1 and 2, we observe that retrieving with global features performs much worse than with the token-level interaction method described in the previous section, which indicates that the Page Retrieval task requires fine-grained interaction between questions and pages, since question-related clues usually occur in local area of the page. Additionally, jointly optimizing TA with PR_g (row 3) hurts

	Model	Page Retrieval			Textual-part Answer				Visual-part Answer		
		R@1	R@3	R@5	B4	M	R-L	C	P	R	F1
1	PR _g	39.0	61.2	71.3	-	-	-	-	-	-	-
2	PR	80.3	93.5	95.8	-	-	-	-	-	-	-
3	PR _g +TA	38.3	60.8	70.4	41.5	31.8	57.4	345.3	-	-	-
4	PR+TA	80.7	93.0	95.6	42.4	32.4	58.5	355.3	-	-	-
5	URA (PR+TA+VA)	81.8	94.4	96.4	42.9	33.0	59.5	361.6	81.1	56.6	66.7
6	3 Single	80.3	93.5	95.8	42.4	32.4	59.6	367.7	75.7	60.1	67.0

Table 4: Comparing the URA model with several baselines on Page Retrieval and Multimodal QA.

Model	Textual-part Answer				Visual-part Answer		
	B4	M	R-L	C	P	R	F1
3 Single	38.0	29.6	55.1	323.9	76.9	50.1	60.7
URA	38.9	30.3	55.5	324.6	82.5	48.3	61.0

Table 5: Evaluate Multimodal QA under the cascade setting.

both TA (row 6) and PR_g (row 1). In contrast, jointly training TA and PR (row 4) affects each other less. It may be because that question answering task requires fine-grained understanding, and it does not conflict with the token-level interaction in PR. Finally, jointly training with VA also helps both the Page Retrieval and TA task (row 5). Compared to the 3 task-specific models, URA achieves even better performance over Page Retrieval and Multimodal QA.

Multimodal QA under cascade setting. Table 5 shows that URA also outperforms multiple task-specific models under the cascade setting. However, we observe a large performance gap between the separate setting and the cascade setting. Considering that the cascade setting is closer to real applications, the Page Retrieval task could be the bottleneck of the MPMQA. Thus, investigating more powerful retrieval models, or models that can directly answer questions from the whole manual will benefit the MPMQA task.

Multimodal QA broken down by Semantic Regions. Table 6 shows the URA performance breaking down in semantic categories. URA performs well on text regions, but worse on visually-rich regions such as product images, illustrations, tables, and graphics, which indicates that a more powerful multimodal understanding ability is required to achieve better performance in the MPMQA task.

Human evaluation. We conduct a human evaluation to verify whether the multimodal answer is helpful for user understanding. We sample 50 question-answer pairs from the test set of PM209. We inference the TA task-specific model and URA on this subset to get the Text-only Answer and Multimodal Answer respectively. Considering the annotators may easily distinguish the two models according to whether there are visual-part outputs, we attach visual-part outputs to the text-only answers with two baseline approaches: 1) random region: randomly selecting two regions from the page; 2) nearest region: selecting the neighbor region that shares the most OCR words with the textual answer. We provide the question and four answers simultaneously to 20 human evaluators, and ask them to rate each

Region	Textual-part Ans.				Visual-part Ans.		
	B4	M	R-L	C	P	R	F1
Text	43.1	33.2	59.8	364.9	84.6	72.6	78.2
Title	38.8	30.6	57.6	366.6	68.4	44.8	54.1
Product image	36.5	29.2	55.2	305.9	58.1	4.5	8.3
Illustration	37.7	30.1	58.9	332.6	32.6	3.0	5.4
Table	36.8	29.1	50.4	271.1	57.4	15.8	24.8
Graphic	32.9	27.1	51.3	271.2	60.8	16.2	25.5

Table 6: Multimodal QA results on each semantic region.

Model	MOS
TA+random region	2.11
TA+nearest region	2.67
URA	3.52
Human	4.70

Table 7: Mean Opinion Score of the human evaluation.

answer by 1-5 points according to whether the answer is helpful to address the given question. The four answers include: two text-only answers attached with visual-part outputs, the multimodal answer produced by our model, and the ground truth multimodal answer by humans. The Mean Opinion Score (MOS) of the 4 answers is shown in Table 7. It shows that the multimodal answer produced by URA are more helpful than text-only answers.

Conclusion

In this paper, we propose the Multimodal Product Manual Question Answering (MPMQA) task, which requires the model to comprehend multimodal content in an entire product manual and answer questions with multimodal outputs. To support the MPMQA task, we construct the large-scale dataset PM209 with human annotations. It contains 22,021 multimodal question-answering pairs on 209 product manuals across 27 well-known consumer brands. The multimodal answer to each question consists of a textual-part in natural language sentences, and a visual-part consisting of regions from the manual. For the MPMQA task, we further propose a unified model that retrieves relevant pages and generates multimodal answers based on multitask learning. It achieves competitive results compared to multiple task-specific models. We release the dataset, code, and model at <https://github.com/AIM3-RUC/MPMQA>.

Acknowledgements

This work was partially supported by the National Natural Science Foundation of China (No. 62072462) and the National Key R&D Program of China (No. 2020AAA0108600).

References

- Anderson, P.; He, X.; Buehler, C.; Teney, D.; Johnson, M.; Gould, S.; and Zhang, L. 2018. Bottom-up and top-down attention for image captioning and visual question answering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 6077–6086.
- Banerjee, S.; and Lavie, A. 2005. METEOR: An automatic metric for MT evaluation with improved correlation with human judgments. In *Proceedings of the acl workshop on intrinsic and extrinsic evaluation measures for machine translation and/or summarization*, 65–72.
- Biten, A. F.; Tito, R.; Maffa, A.; Gomez, L.; Rusinol, M.; Valveny, E.; Jawahar, C.; and Karatzas, D. 2019. Scene text visual question answering. In *Proceedings of the IEEE/CVF international conference on computer vision*, 4291–4301.
- Castelli, V.; Chakravarti, R.; Dana, S.; Ferritto, A.; Florian, R.; Franz, M.; Garg, D.; Khandelwal, D.; McCarley, S.; McCawley, M.; Nasr, M.; Pan, L.; Pendus, C.; Pitrelli, J.; Pujar, S.; Roukos, S.; Sakrajda, A.; Sil, A.; Uceda-Sosa, R.; Ward, T.; and Zhang, R. 2020. The TechQA Dataset. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, 1269–1278. Online: Association for Computational Linguistics.
- Chang, Y.; Narang, M.; Suzuki, H.; Cao, G.; Gao, J.; and Bisk, Y. 2022. Webqa: Multihop and multimodal qa. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 16495–16504.
- Chen, X.; Zhao, Z.; Chen, L.; Ji, J.; Zhang, D.; Luo, A.; Xiong, Y.; and Yu, K. 2021. WebSRC: A Dataset for Web-Based Structural Reading Comprehension. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, 4173–4185. Online and Punta Cana, Dominican Republic: Association for Computational Linguistics.
- Gutmann, M.; and Hyvärinen, A. 2010. Noise-contrastive estimation: A new estimation principle for unnormalized statistical models. In *Proceedings of the thirteenth international conference on artificial intelligence and statistics*, 297–304. JMLR Workshop and Conference Proceedings.
- Hannan, D.; Jain, A.; and Bansal, M. 2020. Manymodalqa: Modality disambiguation and qa over diverse inputs. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 05, 7879–7886.
- Kudo, T.; and Richardson, J. 2018. SentencePiece: A simple and language independent subword tokenizer and detokenizer for Neural Text Processing. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, 66–71. Brussels, Belgium: Association for Computational Linguistics.
- Lin, C.-Y. 2004. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*, 74–81.
- Mathew, M.; Bagal, V.; Tito, R.; Karatzas, D.; Valveny, E.; and Jawahar, C. 2022. InfographicVQA. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 1697–1706.
- Mathew, M.; Karatzas, D.; and Jawahar, C. 2021. Docvqa: A dataset for vqa on document images. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, 2200–2209.
- McKie, J. X.; and Liu, R. 2016. PyMuPDF. <https://github.com/pymupdf/PyMuPDF>. Accessed: 2022-05-01.
- Nandy, A.; Sharma, S.; Maddhashiya, S.; Sachdeva, K.; Goyal, P.; and Ganguly, N. 2021. Question Answering over Electronic Devices: A New Benchmark Dataset and a Multi-Task Learning based QA Framework. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, 4600–4609.
- Papineni, K.; Roukos, S.; Ward, T.; and Zhu, W.-J. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, 311–318.
- Paszke, A.; Gross, S.; Massa, F.; Lerer, A.; Bradbury, J.; Chanan, G.; Killeen, T.; Lin, Z.; Gimelshein, N.; Antiga, L.; et al. 2019. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32.
- Qi, L.; Lv, S.; Li, H.; Liu, J.; Zhang, Y.; She, Q.; Wu, H.; Wang, H.; and Liu, T. 2022. DuReadervis: A: A Chinese Dataset for Open-domain Document Visual Question Answering. In *Findings of the Association for Computational Linguistics: ACL 2022*, 1338–1351.
- Raffel, C.; Shazeer, N.; Roberts, A.; Lee, K.; Narang, S.; Matena, M.; Zhou, Y.; Li, W.; Liu, P. J.; et al. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *J. Mach. Learn. Res.*, 21(140): 1–67.
- Rajpurkar, P.; Jia, R.; and Liang, P. 2018. Know What You Don’t Know: Unanswerable Questions for SQuAD. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, 784–789.
- Rajpurkar, P.; Zhang, J.; Lopyrev, K.; and Liang, P. 2016. SQuAD: 100,000+ Questions for Machine Comprehension of Text. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, 2383–2392.
- Singh, A.; Natarjan, V.; Shah, M.; Jiang, Y.; Chen, X.; Parikh, D.; and Rohrbach, M. 2019. Towards VQA Models That Can Read. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 8317–8326.
- Singh, H.; Nasery, A.; Mehta, D.; Agarwal, A.; Lamba, J.; and Srinivasan, B. V. 2021. Mimoqa: Multimodal input multimodal output question answering. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 5317–5332.
- Talmor, A.; Yoran, O.; Catav, A.; Lahav, D.; Wang, Y.; Asai, A.; Ilharco, G.; Hajishirzi, H.; and Berant, J. 2020. Multi-ModalQA: complex question answering over text, tables and

images. In *International Conference on Learning Representations*.

Tanaka, R.; Nishida, K.; and Yoshida, S. 2021. Visualmrc: Machine reading comprehension on document images. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 15, 13878–13888.

Vaswani, A.; Shazeer, N.; Parmar, N.; Uszkoreit, J.; Jones, L.; Gomez, A. N.; Kaiser, Ł.; and Polosukhin, I. 2017. Attention is all you need. *Advances in neural information processing systems*, 30.

Vedantam, R.; Lawrence Zitnick, C.; and Parikh, D. 2015. Cider: Consensus-based image description evaluation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 4566–4575.

Wang, X.; Liu, Y.; Shen, C.; Ng, C. C.; Luo, C.; Jin, L.; Chan, C. S.; Hengel, A. v. d.; and Wang, L. 2020. On the general value of evidence, and bilingual scene-text visual question answering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 10126–10135.

Wolf, T.; Debut, L.; Sanh, V.; Chaumond, J.; Delangue, C.; Moi, A.; Cistac, P.; Rault, T.; Louf, R.; Funtowicz, M.; et al. 2020. Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 conference on empirical methods in natural language processing: system demonstrations*, 38–45.

Xu, Y.; Li, M.; Cui, L.; Huang, S.; Wei, F.; and Zhou, M. 2020. Layoutlm: Pre-training of text and layout for document image understanding. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 1192–1200.

Yao, L.; Huang, R.; Hou, L.; Lu, G.; Niu, M.; Xu, H.; Liang, X.; Li, Z.; Jiang, X.; and Xu, C. 2021. FILIP: Fine-grained Interactive Language-Image Pre-Training. In *International Conference on Learning Representations*.