

# DATASHEET FOR MDCD-VQA

## **Anonymous authors**

Paper under double-blind review

## 1 DATASHEET

Because of the double-blind mechanism, we have deleted some of the questions, which will be added to the subsequent paper as they become publicly available.

### 1.1 MOTIVATION

#### **For what purpose was the dataset created?**

At a time when large models are rapidly evolving, research efforts are starting to focus more and more on this end-to-end document understanding capability, and a lot of work has been done to validate their performance on DocVQA tasks. Currently, however, almost all of the work has only demonstrated the effect on English data, which limits the application of related models to other languages, such as Chinese. Moreover, these datasets are usually derived from a single domain, which does not provide a good proof of the generalization ability of the models. We propose this dataset, which on the one hand can help to have a good evaluation of the existing models on the Chinese document understanding task, and on the other hand, the training data and labels it provides can be treated as high-quality data to be added to the training of the large models.

### 1.2 COMPOSITION

#### **What do the instances that comprise the dataset represent (e.g., documents, photos, people, countries)?**

The MDCD-VQA dataset comprises document images, accompanied by their respective annotation files. These images represent scanned, photographed or Born-Digital-file transferred documents and are stored in the Joint Photographic Experts Group (JPEG) format.

#### **How many instances are there in total (of each type, if appropriate)?**

The MDCD-VQA contains 5,071 images in total, and a detailed category distribution can be found in paper appendix.

#### **Does the dataset contain all possible instances or is it a sample (not necessarily random) of instances from a larger set?**

Part of the images of MDCD-VQA are randomly sampled from several public datasets, including EPHOIE, EATEN, SCID, CER-VIR, ComFinTab, CDLA, DI, XFUNSD, HUST-CELL, and Baidu-FEST. For other self-collected samples, we collect them from some public newspaper resources database and searching engine (e.g., Google). The number of samples is based on the diversity of samples in the original dataset. We have annotated these images with VQA labels and category labels.

#### **What data does each instance consist of?**

Each instance in the dataset consists of an image and associated annotations. These annotations, stored in JSON format, include a list of question-answer pairs, the OCR annotations (bounding boxes and transcriptions for each line of text), and the document category label.

#### **Is there a label or target associated with each instance?**

Yes.

#### **Is any information missing from individual instances?**

No.

#### **Are relationships between individual instances made explicit (e.g., users' movie ratings, social network links)?**

054 MDCD-VQA desensitizes all information that may contain personally identifiable information (such  
055 as a person’s name in a ticket).

056 **Are there recommended data splits (e.g., training, development/validation, testing)?**

057 The dataset provides a division of the training, validation and test sets, each containing 3,561/ 762/  
058 755 of images, respectively.

059 **Are there any errors, sources of noise, or redundancies in the dataset? If so, please provide a  
060 description.** There could have some potential noise of QA annotation and OCR annotation.

061 **Is the dataset self-contained, or does it link to or otherwise rely on external resources (e.g.,  
062 websites, tweets, other datasets)?**

063 Some of the data in MDCD-VQA are sampled from other public datasets, and for some of the data  
064 where redistribution is restricted, we provide the link to download the original data collection and the  
065 indexes of the corresponding sampled data in the data description. All data labels are self-contained.  
066

067 **Does the dataset contain data that might be considered confidential (e.g., data that is protected  
068 by legal privilege or by doctor–patient confidentiality, data that includes the content of indi-  
069 viduals’ non-public communications)?**

070 No. All samples in MDCD-VQA are public available.

071 **Does the dataset contain data that, if viewed directly, might be offensive, insulting, threatening,  
072 or might otherwise cause anxiety? If so, please describe why.**

073 No.

074 **Does the dataset relate to people? If not, you may skip the remaining questions in this section.**

075 The vast majority of the data in MDCD-VQA comes from other publicly available datasets, which  
076 we have verified have been desensitized to remove personal information when they are made publicly  
077 available. Other parts of our own collection do not contain private information directly related to  
078 people.  
079

### 080 1.3 COLLECTION PROCESS

081 **How was the data associated with each instance acquired?**

082 The collection process is described in Section 3.1 of the main paper. The data is directly observable.  
083

084 **What mechanisms or procedures were used to collect the data (e.g., hardware apparatuses or  
085 sensors, manual human curation, software programs, software APIs)?**

086 The annotation process is mainly described in Section 3.1 of the main paper. And we specifically  
087 developed a web-based tool for data annotation, along with the detailed instructions on the web  
088 page, as shown in Figure 1. With this tool, annotators are able to add, modify, verify, or delete the  
089 QA pairs for an image.

090 **If the dataset is a sample from a larger set, what was the sampling strategy (e.g., deterministic,  
091 probabilistic with specific sampling probabilities)?**

092 The sampling process is a manual selection process in which each image is carefully examined by  
093 the dataset constructors. The main guidelines for selection are representativeness and diversity, i.e.  
094 based on the diversity of data in the original dataset, for data with a certain layout (e.g. train tickets)  
095 the sample size is no more than 50 images because the QA questions that can be constructed based  
096 on this data are basically similar. In contrast, for data with a more diverse and rich layout, as many  
097 different data types as possible are covered.

098 **Who was involved in the data collection process (e.g., students, crowdworkers, contractors)  
099 and how were they compensated (e.g., how much were crowdworkers paid)?**

100 For the annotators, they are all Chinese native speakers from our lab or partner labs. According to the  
101 total working hours and their average salary, the expenditure for the whole annotation is estimated  
102 to be \$5000.

103 **Over what timeframe was the data collected?**

104 Collection of the dataset began in July 2023 and it took an estimated 4 months to complete data  
105 collection and labeling.

106 **Were any ethical review processes conducted (e.g., by an institutional review board)?**

107 We have conduct an internal ethical review process by the company’s legal compliance department.

108  
109  
110  
111  
112  
113  
114  
115  
116  
117  
118  
119  
120  
121  
122  
123  
124  
125  
126  
127  
128  
129  
130  
131  
132  
133  
134  
135  
136  
137  
138  
139  
140  
141  
142  
143  
144  
145  
146  
147  
148  
149  
150  
151  
152  
153  
154  
155  
156  
157  
158  
159  
160  
161



(a) Screenshot of the web-based annotation tool.

(b) Illustration of the instruction for annotators.

Figure 1: We developed a web-based annotation tool, as shown in (a). The annotation tool will randomly push the images to the annotator, who can see the history of the annotated content. In the table on the right side, annotators can directly edit and delete the history of the annotation content, and can also click on the verification button to review a annotation; the annotator can also submit new annotation information in the form below. Below the tool will display the number of images that have been added, modified and reviewed by the current annotator in real time. The annotation instructions is shown in (b).

**Does the dataset relate to people? If not, you may skip the remainder of the questions in this section**

The annotators is only asked to design the QA data based on the images provided, and there is no collection of individual-related information involved.

1.4 PREPROCESSING/CLEANING/LABELING

**Was any preprocessing/cleaning/labeling of the data done (e.g., discretization or bucketing, tokenization, part-of-speech tagging, SIFT feature extraction, removal of instances, processing of missing values)?**

The annotations of the MDCD-VQA dataset come from two parts, with each part accounting for about 50%. The first part is based on the original data labels provided by the original public dataset (e.g., the labels of the KIE task), and we randomly selected some of them and transformed them into the form of QA by designing various templates. And the question form is enriched by natural language extensions (e.g., mutual translation). The other part is designed by the annotators directly based on the provided images, and this part of QA will avoid the existing QA style as much as possible to achieve a more liberalized question design.

For each image, annotators can see what has been annotated so far and can choose to modify or delete the current QA content. Each question (including both semi-automated and manually annotated questions) must be reviewed by at least one person (by clicking the Reviewed button). When the number of questions in an image reaches the threshold and all questions have been reviewed, the annotator is no longer pushed. The annotator will continue to annotate or review until no new data is pushed.

**Was the “raw” data saved in addition to the preprocessed/cleaned/labeled data (e.g., to support unanticipated future uses)?**

No. The images contained in the MDCD-VQA dataset have not been additionally processed, and the referenced original data and data labels can be accessed via the provided connection.

**Is the software that was used to preprocess/clean/label the data available?**

Currently not. We will consider open-sourcing this web-based tool in the future.

162 1.5 USES

163

164 **Has the dataset been used for any tasks already?**

165 No.

166 **Is there a repository that links to any or all papers or systems that use the dataset?**

167 It is a new dataset that haven't been used by current works. We run existing state-of-the-art models  
168 and release the code.

169 **What (other) tasks could the dataset be used for?** In addition to being used directly for tasks  
170 related to large model evaluation or document VQA, this dataset can be used for various OCR or  
171 document understanding related tasks, e.g., information extraction tasks can be constructed using  
172 only the extractive questions contained in the dataset; document classification tasks can be con-  
173 structed using the document category labels provided by the dataset; the OCR labels provided in  
174 this dataset can be used to improve the generalization ability of the OCR model in various document  
175 scenarios, etc.

176 **Is there anything about the composition of the dataset or the way it was collected and prepro-  
177 cessed/cleaned/labeled that might impact future uses?**

178 Since the current dataset is new, none of the existing large models have used it in their own training  
179 tasks. However, if the dataset is made public, there is a high probability that it will be used by future  
180 new big models in their training process (the training process of existing large models generally  
181 collects as much high-quality data as possible to improve the generalization ability of the models),  
182 which will somewhat affect the fairness of evaluation among models on this dataset.

183 **Are there tasks for which the dataset should not be used?**

184 The dataset should not be used for commercial usage.

185

186 1.6 DISTRIBUTION

187

188 **Will the dataset be distributed to third parties outside of the entity (e.g., company, institution,  
189 organization) on behalf of which the dataset was created?**

190 No.

191 **How will the dataset will be distributed (e.g., tarball on website, API, GitHub)?**

192 All code and dataset will be publicly distributed on GitHub.

193 **When will the dataset be distributed?**

194 It will be published along with the paper.

195 **Will the dataset be distributed under a copyright or other intellectual property (IP) license,  
196 and/or under applicable terms of use (ToU)?**

197 The dataset is under Creative Attribution-NonCommercial-ShareAlike 4.0 International (CC BY-  
198 NC-SA 4.0) License.

199 **Have any third parties imposed IP-based or other restrictions on the data associated with the  
200 instances?**

201 No.

202 **Do any export controls or other regulatory restrictions apply to the dataset or to individual  
203 instances?**

204 All authors bear all responsibility for the dataset in case of violation of rights, etc.

205

206 1.7 MAINTENANCE

207

208 **How can the owner/curator/manager of the dataset be contacted (e.g., email address)?**

209 The contact email will be updated in real time in the introduction page of the dataset.

210 **Is there an erratum?**

211 No.

212 **Will the dataset be updated (e.g., to correct labeling errors, add new instances, delete in-  
213 stances)?**

214

216 If significant errors are reported by dataset users, we will consider updating the dataset accordingly.  
217 The update information will also be posted on the dataset page.

218  
219 **If the dataset relates to people, are there applicable limits on the retention of the data asso-**  
220 **ciated with the instances (e.g., were the individuals in question told that their data would be**  
221 **retained for a fixed period of time and then deleted)?**

222 No.

223 **Will older versions of the dataset continue to be supported/hosted/maintained?**

224 Yes. If we plan to update the data, we will keep the original version available and then release the  
225 follow-up version.

226 **If others want to extend/augment/build on/contribute to the dataset, is there a mechanism for**  
227 **them to do so?**

228 Yes, others who are providing some issue fixes can just make a pull request on Github or contact us  
229 privately. If it's to expand new work based on our data, it just needs to follow the license.

230

231

232

233

234

235

236

237

238

239

240

241

242

243

244

245

246

247

248

249

250

251

252

253

254

255

256

257

258

259

260

261

262

263

264

265

266

267

268

269