# M2BIT: A Multi-Modal Bilingual Instruction Tuning Dataset for Visual Language Models

Anonymous ACL submission

#### Abstract

In this paper, we introduce the Multi-Modal 001 Bilingual Instruction Tuning dataset (M2BIT), specifically designed to enhance the perfor-004 mance of vision language models (VLMs). Our M2BIT dataset is one of the largest multi-006 modal instruction tuning datasets available, covering 40 diverse vision-language tasks in both English and Chinese. It comprises 2 million instances, each accompanied by 400 manually written task instructions. With a care-011 fully curated annotation process, we strive to elevate the quality of response, thereby enrich-012 ing the user experience while minimizing the generation of potential hallucinations. To val-014 idate the efficacy of M2BIT, we train a VLM known as Ying-VLM using this dataset, delving into the impact of instruction tuning across di-017 verse languages and modalities. Upon comparing it with strong VLM baselines, Ying-VLM demonstrates superior performance on complex knowledge vision question answering tasks. Moreover, it exhibits a lower propensity for hal-023 lucination, displays greater generalization capabilities to previously unseen video tasks, and better comprehends novel instructions in Chinese. We will open-source the M2BIT dataset and trained models to facilitate future research. 027

### 1 Introduction

028

034

040

Following the substantial success of Chat-GPT (OpenAI, 2022), the interest in designing a versatile intelligent assistant that can understand and interact with the multi-modal world has surged. The potential of transforming Large Language Models (LLMs) into powerful Vision Language Models (VLMs) has been demonstrated by further training on image-text pairs or implementing specialized visual instruction tuning (Zhu et al., 2023; Liu et al., 2023b). This enhancement allows LLMs to *see* the world, offering promising capabilities that could significantly assist individuals with disabilities (OpenAI, 2023).

Dataset	Tasks	LAN	Samples	Ins. / Task
MiniGPT4	N / A	En	5K	N / A
LLaVA	3	En	158K	N / A
MultiModalGPT	3	En	6K	5
MultiInstruct	26	En	235K	5
InstructBLIP	28	En	1.6M	9.7
M2BIT (Ours)	40	En, Zh	2M	10

Table 1: Comparison of different multi-modal instruction tuning datasets. Ins. denotes Task Instruction and N / A means instructions are artificially generated. Our M2BIT is one of the largest datasets, covering 40 tasks in English and Chinese.

042

043

044

047

050

051

054

056

060

061

062

063

064

065

066

067

068

069

The quality of the visual instruction tuning dataset is pivotal in the development of VLMs, as indicated by findings in LLMs (Zhou et al., 2023). Recent research efforts in this area can be broadly grouped into two categories. The first stream of research strives to augment existing academic vision-language (V+L) datasets with manually written task instructions (Xu et al., 2022; Dai et al., 2023). Although VLMs trained on these datasets exhibit notable performance on academic benchmarks, they often generate responses that are excessively terse and blunt (Chen et al., 2023a). This brevity, a typical characteristic of academic datasets, compromises the user experience during interactions. The second stream of research employs image annotation tools to generate textual descriptions of original images, subsequently leveraging models like ChatGPT/GPT-4 to create dialog-style instruction tuning datasets (Zhu et al., 2023; Liu et al., 2023b; Zhao et al., 2023a). However, a pitfall of training VLMs with these pseudogrounded dialogs is the risk of exacerbating LLMs' hallucination problem, possibly resulting in inconsistent image descriptions featuring non-existent objects. Besides, current studies mainly focus on English tasks, which limits the investigation of the cross-lingual effects of instruction tuning.

In this paper, we introduce M2BIT, a Multi-

# List[String]: the base64 string representation of a profile photo of F. Scott Fitzgerald
Images: ["IVBORw0KGg5ErkJggg=="]
# String: task instruction
Instruction: "Analyze the image and provide an
appropriate response to the question. "
# String: task-specific inputs, e.g., a question related
to the image.
Inputs: "On which book by this man, Baz luhrmann's
planned a film?"
# String: task outputs, e.g., the correct answer for the question.
Outputs: "Baz Luhrmann has planned a film adaptation
of the book The Great Gatsby."
# Dictionary: a meta information dictionary contains
original data.
Meta Data: {"original_answer" : "The Great
Gatsby", ,"wikipedia_id": "152171"}

Figure 1: The unified data format schema of our M2BIT.

Modal Bilingual Instruction Tuning dataset, which leverages valuable academic benchmarks and the 071 capabilities of ChatGPT for dialog style enhancement. M2BIT is meticulously constructed in three stages: (1) Task Selection and Manual Instruction Writing: Our dataset consists of diverse tasks, including traditional image-text tasks like visual question answering and image captioning, as well 077 as video-related tasks such as video question answering. Annotators are instructed to review the dataset paper thoroughly and craft 10 unique instructions for each task. We further incorporate Chinese vision-language datasets with corresponding Chinese instructions, resulting in a comprehensive compilation of 40 diverse tasks and 400 instructions. (2) Data Format Unification: We ensure that all tasks within our dataset adhere to a unified vision-to-text format. This format comprises four fields: images, instruction, inputs and outputs. Additional information, such as bounding 089 box details, is embedded within the images, and short answers are rephrased using ChatGPT while incorporating contextual information, where available. (3) Quality Check: For quality control, we assign an extra annotator to each task to review 20 examples from each split of every dataset. A task is considered complete only after the annotator has verified the accuracy and consistency of the images, instructions, inputs, and outputs for each instance. As demonstrated in Table 1, M2BIT is one of the largest multi-modal instruction tuning datasets regarding the number of instructions and samples, 101 covering diverse tasks in English and Chinese. 102

To substantiate the effectiveness of the M2BIT dataset, we develop Ying-VLM, merging the ca-

103

104

pabilities of a vision encoder, BLIP-2 (Li et al., 2023a), with Ziya-13B (Zhang et al., 2022), which is a derivative of LLaMA (Touvron et al., 2023a). We leverage the proven methodology of incorporating visual tokens as prefix prompts in LLMs and utilize a two-stage training regime. The initial stage aligns vision features with text embeddings via image captioning on LAION-400M (Schuhmann et al., 2021), while the second stage enhances the model by performing instruction tuning on M2BIT.

105

106

107

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

Our evaluation of the instruction tuning effect with M2BIT is threefold. Firstly, we evaluate Ying-VLM on knowledgeable VQA (KVQA) tasks, including OK-VQA, A-OKVQA, and a held-out dataset, ViQuAE. These tasks, which demand VLM's comprehension of image context and reasoning with the knowledge acquired by LLMs, have gained wide recognition as benchmarks for evaluating VLMs (Dai et al., 2023; Bai et al., 2023b). Secondly, we conduct a hallucination evaluation on image captioning, following the methodology by Li et al. (2023b). Lastly, we perform zero-shot transfer evaluations on Chinese V+L tasks and video-language tasks to scrutinize the cross-language/modal effect of instruction tuning. The experimental results highlight that Ying-VLM surpasses potent baseline VLMs in KVQA tasks, is less susceptible to hallucination than models trained with generated pseudo-grounded dialogs, and demonstrates enhanced generalization capabilities when confronted with unseen video and cross-lingual tasks. These results underscore the potential of our proposed M2BIT dataset in constructing robust VLMs and investigating instruction tuning effects across languages and modalities.

# 2 M2BIT: A Multi-Modal Bilingual Instruction Tuning Dataset

In this section, we introduce the M2BIT dataset by first elaborating the task coverage (§ 2.1), followed by the annotation process details (§ 2.2). In § 2.3, we present the dataset format and the statistics of the crafted dataset instructions.

#### 2.1 Task Coverage

Our dataset makes a wide coverage of the current existing visual-language and video-language benchmarks, from simple image captioning to complicated reasoning based on the image even beyond the visual content:

Captioning This task aims to produce descrip-

tions of the given images according to different 154 needs. We include MS COCO (Lin et al., 2014) 155 (the Karpathy split) for generic image descriptions. 156 TextCaps (Sidorov et al., 2020) requires models 157 to capture the text presented in the image and 158 generate captions accordingly. Image-Paragraph-159 Captioning (Krause et al., 2017) focuses on gener-160 ating detailed descriptions for images. 161

162

163

164

165

166

167

168

170

171

- **Reasoning** This task evaluates specific reasoning capabilities. We incorporate CLEVR (Johnson et al., 2017) and NLVR (Suhr et al., 2017) for spatial reasoning, Visual Commonsense Reasoning (VCR) (Zellers et al., 2019) for commonsense reasoning, Visual MRC (Tanaka et al., 2021) for reading comprehensive over images, and Winoground (Thrush et al., 2022) for fine-grained semantics reasoning over text descriptions and image contents.
- Visual Question Answering (VQA) This is the 172 most widely studied multi-modal task, which re-173 quires the model to answer a given question based 174 on the image correctly. Tasks include VQA 175 v2 (Goyal et al., 2017b), Shapes VQA (Andreas 176 et al., 2016), DocVQA (Mathew et al., 2021), 177 OCR-VQA (Mishra et al., 2019), ST-VQA (Biten 178 et al., 2019), Text-VQA (Singh et al., 2019), and 179 GQA (Hudson and Manning, 2019).
- **Knowledgeable Visual Question Answering** 181 (KVQA) Unlike traditional VQA tasks focusing on 182 the question relevant to the content image, KVQA requires the model to draw upon outside knowledge to answer questions. We incorporate two outside 185 knowledge VQA datasets: OK-VQA (Marino et al., 186 2019) and A-OK-VQA (Schwenk et al., 2022), Sci-187 enceQA (Lu et al., 2022) which contains multimodal science questions, and ViQuAE (Lerner et al., 2022) focusing on knowledge facts of named entities in images. 191
- Classification This task involves classifying an 192 image based on a given set of candidate labels. 193 ImageNet (Russakovsky et al., 2015), Grounded 194 Object Identification (COCO-GOI) (Lin et al., 195 2014), COCO-Text (Veit et al., 2016), Image Text 197 Matching (COCO-ITM) (Lin et al., 2014), e-SNLI-VE (Kayser et al., 2021), Multi-modal Fact Check-198 ing (Mocheg) (Yao et al., 2022), and IQA (Duanmu 199 et al., 2021) are included. Due to language model input length constraints, we reduce the number of 201 202 options in some datasets with extensive candidate labels, such as ImageNet.
- **Generation** Visual conditional generation requires

Number of different instructions - Image Captioning	400
- Classification	113
- Visual Question Answering	95
- Knowledgeable Visual QA	40
- Generation	40
Tokens per instruction	$24.4\pm9.6$
Instruction edit distance among the same task	$76.6\pm37.2$
Instruction edit distance across tasks	$106.6\pm39.5$

Table 2: The statistics of our instructions.

models to understand the visual content and make a composition meeting the task demand. We have Visual Storytelling (VIST) (Huang et al., 2016), Visual Dialog (VisDial) (Das et al., 2017), and multimodal machine translation Multi30k (Elliott et al., 2016) in this category. 205

206

207

208

209

210

211

212

213

214

215

216

217

218

219

220

221

222

223

224

225

226

227

228

229

230

231

232

233

234

235

236

237

238

239

240

241

242

243

**Chinese Vision-Language Tasks** To examine the effect of instruction tuning on different languages, we incorporate several Chinese vision-language tasks including FM-IQA (Gao et al., 2015) for VQA, COCO-CN (Li et al., 2019) and Flickr8k-CN (Li et al., 2016) for captioning, Chinese Food Net (Chen et al., 2017) for classification, and MM-Chat (Zheng et al., 2022) for generation.

Video-Language Tasks Beyond the static images, we are interested in whether instruction tuning can be applied to video-language tasks. We include the classic MSR-VTT datasets (Xu et al., 2016) for video captioning, MSRVTT-QA (Xu et al., 2017), ActivityNet-QA (Yu et al., 2019), iVQA (Yang et al., 2021) and MSVD-QA (Xu et al., 2017) for video question answering, Something-Something (Goyal et al., 2017a) for video action classification.

# 2.2 Annotation Process

To build high-quality multi-modal instruction datasets, we rewrite various datasets into a visionto-text format. The annotation process includes three steps: (1) writing instructions for each task, (2) structuring images and texts into a unified schema, and (3) checking the overall dataset quality. Eight authors of this work are employed as human annotators, each of whom is a graduate student familiar with relevant literature.

**Stage I: Instruction Writing** To build high-quality instructions, we first ask annotators to carefully read the dataset paper and check the original dataset with some instances to get a clear understanding of the task. After that, they are required to write 10 diverse task instructions manually, covering the



Figure 2: (Left) On region-based tasks, bounding boxes are added to serve as a visual referring prompt. (Right) Short answer rephrasing to improve the response quality, e.g., incorporating rationales into answers.

key characteristics of the task. Table 2 shows the statistics of the written instructions for each task and Figure 7 in the Appendix visualizes the instruction verb distribution. In total, we annotate 400 instructions for all tasks. The average length per instruction is 24.4. To evaluate the diversity of annotated instructions, we employ the average edit distance to measure the similarity between two strings. The average edit distance within the same task is 76.6, indicating a good range of instruction diversity.

245

247

248

249

254

Stage II: Data Format Unification After the in-256 struction has been written according to the task characteristics, we further process the images and corresponding text for a unified instance schema. For most datasets, we keep the original images and 260 text, where images are converted into correspond-261 ing base64 encoded strings for easy data loading. We perform two modifications on potential examples: (1) Adding Bounding Box to Images. For tasks designed for specific regions in the image, a straightforward solution is to provide the bounding box information in natural language for informing the language models of the regions in interest. However, the image pre-processing techniques adopted by different vision encoders may resize the original image, and the original bounding box annotation thus needs further adjustments. Inspired by the recent observation that vision encoders such 273 as CLIP (Radford et al., 2021) are sensitive to the visual prompt (Shtedritski et al., 2023), we directly tag the bounding box as a red rectangle to the image, serving as a visual referring prompting (OpenAI, 2023) for VLMs to focus on the target region. 278 (2) Short Answer Rephrasing. As recent studies 279 have shown that the original short and brief answers

Teels	Description		Total #samples			
Task			Val	Test		
CAP	Given an image, write a description for the image.	679,087	41,462	27,499		
CLS	Given an image, classify the image into pre-defined categories.	238,303	100,069	21,206		
VQA	Given an image, answer a question relevant to the image.	177,633	46,314	10,828		
KVQA	Given an image, answer the question requires outside knowledge.	39,981	11,682	5,477		
REA	Given an image, conduct reasoning over the images.	99,372	11,500	10,000		
GEN	Given an image, make compositions with certain requirements.	145,000	11,315	17,350		
Chinese	CAP, CLS, VQA, and GEN tasks in Chinese.	192,076	77,306	4,100		
Video	CAP, CLS, and VQA tasks on video-language datasets.	20,868	7,542	9,294		
Total	2,005,264					

Table 3: M2BIT task descriptions and statistics. We aggregate instance counts for training, validation, and test sets across all tasks, totaling 2M instances.

in the common VQA dataset could negatively influence the model generation performance (Dai et al., 2023; Chen et al., 2023b), we propose to utilize the ChatGPT (gpt-3.5-turbo-0301) (OpenAI, 2022) model for rephrasing the original answers, by providing origin question and answer with potential extra contextual information. Contextual information includes the caption of the original images, rationales for specific VQA tasks and OCR tokens for the scene-related question, which make the rephrased answers more engaging and informative. The original short answers are also kept in our dataset for potential usage. Figure 2 shows these two modifications.

**Stage III: Quality Check** In this stage, we assign a different annotator to each task to review 20 examples from each split of every dataset. During this stage, we identify minor format inconsistencies between tasks and address them by standardizing the task formats. We observe that answer rephrasing greatly improves the response quality, e.g., more than 95% instances we checked are perceived as better than original concise answers, while Chat-GPT refused to rephrase a few answers (less than 3% of examined instances) due to insufficient image information. We employ simple heuristics to



Figure 3: The overview of Ying-VLM model architecture and training stages.

filter these failed answers and use a basic template to convert the original answer into a sentence, with manual validation to ensure the quality. Finally, the task dataset is deemed complete once the annotator can successfully load it and re-examine the accuracy of the instructions, inputs, and outputs for each instance examined.

#### **Dataset Format and Statistics** 2.3

307

310

311

312

314

317

321

323

324

329

331

333

339

340

341

The instance in our dataset consists of five fields: (1) Images: we represent the images with the po-316 tentially added bounding box by a base64 string. (2) Instruction: we randomly select an instruction from the task instruction pool for each instance. 319 (3) Inputs: we allocate this field for providing task-specific inputs to the model, e.g., the question in the VQA tasks. For tasks such as captioning, there is no extra input so the corresponding field is left as an empty string. (4) Outputs: the required output to the specific tasks, such as the 325 description of the image for captioning tasks and the answer to the image-related question. (5) Meta Data: we provide this field to preserve important 328 information such as image ID for referencing the original dataset. With the clear distinction of these fields, the user of our benchmark can construct the training instances needed flexibly and evaluate the models conveniently. Table 3 gives the statistics aggregated by tasks and Figure 1 illustrates a case in M2BIT. For the detailed statistics, and the license of each dataset in M2BIT, we refer readers to Appendix **B**.

#### **Experiments** 3

We first introduce the experimental setups ( $\S$  3.1), then report and discuss the evaluation results (§ 3.2). Lastly, we analyze instruction robustness and provide a case study ( $\S$  4).

#### 3.1 **Experimental Settings**

#### **Implementation Details** 3.1.1

Model Architecture: Inspired by the recent success of BLIP-2 (Li et al., 2023a), we adopt the vision encoder and the Q-former architecture in the BLIP2-OPT-2.7B (Li et al., 2023a) model to extract relevant visual features from images. For the large language models, we utilize Ziya-13B (Zhang et al., 2022) derived from LLaMA (Touvron et al., 2023a) with bilingual (English and Chinese) ability.

343

345

346

347

350

351

353

354

355

357

358

359

360

361

362

363

364

365

366

367

368

369

370

371

372

373

374

375

376

377

378

379

381

382

383

385

388

389

390

391

392

Training Datasets: We employ a two-staged training consisting of a visual-text alignment stage (Stage I) and a multi-modal instruction tuning stage (Stage II). Stage I aims to align the visual and textual feature space, where we utilize the instructions in the coco captioning and perform an initial alignment training on LAION 400M (Schuhmann et al., 2021). In Stage II, we adopt our M2BIT dataset for training the model, excluding ViQuAE for KVQA evaluation, Winogound as it only has a test set, and all Chinese V+L and video-language datasets for cross-lingual/modality effect investigation. Details of Stage II training datasets are in Figure 8 and Figure 3 provides an overview of the model and training stages.

Hyper-parameter Settings In Stage I, we only train the Q-Former and the language projection for modality alignment, resulting in a total 130M parameters to optimize with AdamW (Loshchilov and Hutter, 2019). The batch size is set to 256 to maximize the utilization of GPU and the model is trained with 300k steps. The learning rate linearly increases to a peak value of 5e-5 in the first 2000 steps and follows a cosine decay scheduler. The weight decay is set to 0.05. We train the model after Stage I for 3 epochs and with a lower learning rate of 1e-5 and a warmup stage of 1000 steps. Inspired by LoRa tuning (Hu et al., 2022), the weights for mapping query and value vectors in the attention layer of LLMs are learnable in this stage to better adapt to the instruction tuning dataset. Other training parameters are consistent with Stage I. All experiments are conducted with 8 NVIDIA 80GB A100 GPUs. It took about 10 days for Stage I and Stage II can be finished in a day.

#### **3.1.2** Evaluation Setup

Evaluated Tasks We conduct three evaluations to understand the instruction tuning effect comprehensively with our M2BIT: (1) Evaluation of KVQA tasks, which includes OK-VQA, A-OKVQA and a

Model	OK-VQA	A-OKVQA	ViQuAE
BLIP2-Flan-T5-XXL	9.1	15.6	9.7
MiniGPT4	23.3	21.8	24.4
InstructBLIP	7.1	5.9	7.3
Ying-VLM <sub>LLaVA</sub>	26.4	22.5	24.3
Ying-VLM	27.5	24.5	29.6

Table 4: ROUGE-L evaluation results of KVQA tasks.

held-out dataset ViQuAE. These tasks are widely adopted in evaluation for VLMs (Dai et al., 2023; Bai et al., 2023b) as they pose a great challenge for VLMs to understand the image context and perform reasoning with the knowledge acquired by LLMs. (2) Evaluation of object hallucination, which refers to a phenomenon that the model produces image descriptions that contain objects that are not anchored 400 with or even absent from the target image. We fol-401 low the exact setup in Li et al. (2023b) to perform 402 hallucination analysis on 2,000 images randomly 403 404 sampled from the MSCOCO dataset (Lin et al., 2014). (3) Evaluation of cross-language/modality 405 transferability. We hold out all Chinese V+L tasks 406 and video-language tasks during the instruction 407 tuning stage, then perform a zero-shot transfer to 408 409 investigate whether instruction tuning is generalizable across languages, i.e., English to Chinese, and 410 modalities, i.e., images to videos. For a compre-411 hensive evaluation, we also present the results on 412 the generation performance on Visual-Dialog and 413 image classification on ImageNet in Appendix D. 414 In all experiments, we use greedy decoding in in-415 ference for deterministic results. 416

Metrics We adopt ROUGE-L (Lin, 2004) as an 417 automatic metric to assess the consistency between 418 predictions and ground-truth answers, focusing on 419 evaluating the model's conversational abilities. As 420 the automatic metric may not fully capture the nu-421 ances of conversational quality, we further intro-422 423 duce GPT-4 as a proxy of human evaluators (§ 4). For the object hallucination, we follow Li et al. 424 (2023b) to adopt  $CHAIR_I$  and  $CHAIR_S$  (Rohrbach 425 et al., 2018). CHAIR $_I$  denotes the proportion of 426 hallucinated ones in all generated objects, while 427 CHAIR<sub>S</sub> describes the hallucination at the sen-428 tence level, i.e., the proportion of generated cap-429 tions that contain hallucinated objects. Appendix H 430 provides a detailed definition for these two metrics. 431 Baselines We employ open-sourced VLMs as the 432 off-the-shelf models for comparison, using the 433 original model weights and greedy decoding: (1) 434 BLIP-2-Flan-T5-XXL (Li et al., 2023a) where an 435

Model	Len	$\operatorname{CHAIR}_{I}(\downarrow)$	$\mathrm{CHAIR}_S(\downarrow)$	Avg. $(\downarrow)$
mPLUG-Owl*	98.5	30.2	76.8	53.5
LLaVA*	90.7	18.8	62.7	40.8
MiniGPT-4*	116.2	9.2	31.5	20.4
InstructBLIP*	7.5	<b>2.5</b>	<b>3.4</b>	<b>3.0</b>
Ying-VLM <sub>LLaVA</sub>	62.7	<b>11.0</b>	36.0	23.5
Ying-VLM	34.2	12.6	<b>16.8</b>	14.7

Table 5: Object hallucination evaluation with instruction "*Provide a brief description of the given image*". Len denotes the average length of generated captions. \* denotes results collected from Li et al. (2023b).

Model	Flickr-8k-CN	FM-IQA	Chinese-FoodNet
MiniGPT4	9.6	20.1	5.0
InstructBLIP	5.2	2.3	1.0
Ying-VLM	20.5	33.3	49.8
Ying-VLM <sub>zh</sub>	20.0	39.8	0.1

Table 6: Results Chinese vision-language tasks. Our model generalizes well on unseen Chinese captioning, VQA and classification tasks.

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

instruction-tuned Flan-T5 (Wei et al., 2022) is connected with a powerful vision encoder to perform a series of multi-modal tasks;<sup>1</sup> (2) MiniGPT-4 which aligns a CLIP visual encoder with a frozen Vicuna-13B (Chiang et al., 2023) with artificially collected dialog dataset;<sup>2</sup> (3) InstructBLIP, an instruction tuning enhanced VLM with Vicuna-13B with converted multi-model datasets and the LLaVA (Liu et al., 2023b) dataset generated by text-only GPT-4.<sup>3</sup> We refer readers to Appendix E for a detailed introduction for these methods and related work. To isolate the effect of the base LLM, we train a variant of our model, Ying-VLM<sub>LLaVA</sub>, by only replacing the Stage II training datasets with the LLaVA dataset (Liu et al., 2023b), using the same training setup with Ying-VLM.

#### 3.2 Results

**Evaluation of Knowledgeable VQA** The results on the KVQA benchmarks are shown in Table 4. In comparison to the strongest baseline, our model achieves an improvement of 1.1 and 2.0 ROUGE-L points for OK-VQA and A-OKVQA, respectively. Additionally, Ying-VLM delivers the best performance on the held-out ViQuAE dataset. These findings indicate that instruction tuning on M2BIT

https://huggingface.co/Salesforce/ blip2-flan-t5-xxl

<sup>&</sup>lt;sup>2</sup>https://github.com/Vision-CAIR/ MiniGPT-4

<sup>&</sup>lt;sup>3</sup>https://huggingface.co/Salesforce/ instructblip-vicuna-13b

461 effectively harnesses knowledge from LLMs and462 elevates response quality.

Evaluation of Object Hallucination As shown 463 in Table 5, the VLMs trained on instruction tun-464 ing datasets generated by ChatGPT/GPT-4 exhibit 465 serious hallucination problems, as indicated by 466 the relatively high average CHAIR of LLaVA and 467 MiniGPT-4. InstructBLIP performs the best on this evaluation. However, it should be noted that the 469 too-short answers provided by InstructBLIP may 470 result in a lack of politeness in responses, which 471 can harm the user experience. This was validated in 472 the later evaluation with GPT-4. When comparing 473 Ying-VLM<sub>LLaVA</sub> to Ying-VLM, it can be observed 474 that Ying-VLM achieves a much lower CHAIR<sub>S</sub> 475 score (16.8 v.s. 36.0), demonstrating that it suffers 476 significantly less from the sentence-level hallucina-477 tion issue. These results suggest that M2BIT could 478 help VLMs to achieve a better balance between the 479 hallucination problem and response quality. 480

Cross-Language Transferability We assess mod-481 els on three unseen Chinese vision-language tasks 482 to investigate the cross-language generalization ef-483 fect of instruction tuning. BLIP-2 and Flan-T5 are 484 not considered here as they do not support Chinese 485 outputs.<sup>4</sup> As illustrated in Table 6, our model per-486 forms well on all evaluated tasks compared with 487 MiniGPT4 and InstructBLIP. While the gain can be attributed to the Chinese ability of the underly-489 ing Ziya-13B LLM, it promisingly indicates that 490 instruction tuning with English datasets can effec-491 tively generalize to different languages. Further, we 492 perform continual training on Chinese V+L tasks, 493 resulting in YingVLM<sub>zh</sub>. As shown in the last row 494 of Table 6, the scores of FM-IQA can be further 495 enhanced. However, the model suffers from catas-496 trophic forgetting as there are no classification tasks 497 in the continual training, resulting in poor perfor-498 mance on the Chinese-FoodNet. An in-depth inves-499 tigation for this problem could be promising (Zhai 500 et al., 2023). 501

> **Cross-Modality Transferability** To evaluate performance on video-language tasks, we uniformly sample 8 frames from each video. MiniGPT4 is excluded as it does not support video inputs. Following InstructBLIP (Dai et al., 2023), we concatenate the visual embedding extracted from the Q-former of each frame as a prefix embedding to the language model. As demonstrated in Table 7, our model

502

503

506

509



Figure 4: Evaluation results using GPT-4 as an evaluator. Our model outperforms MiniGPT-4 and InstructBLIP with a winning rate at 55.6% and 65.5%, respectively.

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

excels in these challenging settings, significantly surpassing the BLIP-series baselines. It is worth noting that the training dataset does not include any videos inputs, implying that our instruction tuning effectively aids the model in generalizing to inputs with a temporal dimension. Furthermore, we train YingVLM on the video-language datasets in M2BIT. The resulting YingVLM<sub>video</sub> achieves the best performance on all tasks, suggesting the effectiveness of video instruction tuning in boosting the video understanding abilities of VLMs.

### 4 Analysis

Evaluation with GPT-4 To further validate the quality of the generated response, we propose to utilize GPT-4 as a proxy of human evaluators (Peng et al., 2023; Gilardi et al., 2023). Specifically, following Vicuna (Chiang et al., 2023), we query GPT-4 to rate the performance of different models against our Ying-VLM. For each sample, we construct a prompt consisting of the original question, its corresponding reference answer, the response generated by our Ying-VLM, and a baseline system output. GPT-4 (gpt-4-0613) is asked to rate both responses on a scale of 10 based on the given question and its reference answer. The ratings are primarily based on the accuracy, relevance, and naturalness of the response to meet the requirements when humans are interacting with multi-modal agents (evaluation template is provided in Appendix G). We swap the order of candidate responses to mitigate potential evaluation biases (Wang et al., 2023). Considering the API cost of GPT-4, 300 examples are randomly sampled from OK-VQA, A-OKVQA and ViQuAE datasets as a subset for evaluation. Figure 4 shows that our Ying-VLM outperforms baseline models in most samples. For example, the GPT-4 evaluator favors Ying-VLM over MiniGPT4 on 167 over 300 tested

<sup>&</sup>lt;sup>4</sup>For all models, we introduce a prompt to promote Chinese outputs. See Appendix F for details.

Model	Video Captioning	y Video Question Answering				
	MSRVTT	iVQA	ActivityNet-QA	MSRVTT-QA	MSVD-QA	
BLIP-2-Flan-T5-XXL	8.8	11.1	8.9	10.3	13.2	
InstructBLIP	14.3	6.3	9.3	4.0	7.0	
Ying-VLM	14.2	23.5	21.9	18.3	21.4	
Ying-VLM <sub>video</sub>	16.1	41.3	44.8	32.0	40.5	

Table 7: ROUGE-L results of video-language tasks. YingVLM shows great generalization from images to videos and training on video datasets (YingVLM<sub>video</sub>) further boosts the performance.

Model	iVQA	ActivityNet-QA	MSRVTT-QA	MSVD-QA
BLIP2-Flan-T5-XXL InstructBLIP	25.0	30.0 38.0	19.0 23.0	32.0
Ying-VLM <sub>Video</sub>	40.0	40.0	40.0	58.0





Human:请给图中的食 Classify the food in the Human: Who plays at this Hu man: On which book by this n, Baz luhrmann's planned a (A)薯条 (French fri (B)凉拌西红柿 (Te film MiniGPT4: The stadium is 副大虾 (13 ed Sh 南贝 (Scallop in Shell) alled Wembley Stadium, ome to the English nation MiniGPT4: Th (E)集低 (Oysters) (F)螃蟹 (Crab) (G)扬州炒饭 (Yang (H)披萨 (Pizza) InstructBLIP: Arsenal F.C. uctBLIP: Great (I)蛋挞 (Egg Tart) (J)肉酱意大利面 (Spaghetti with m Ying-VLM: The stadium you are looking at is called Emirates Stadium and it is the Ying-VLM: The book on which Baz Luhrmann's planned film is MiniGPT4: This is a bowl of fried rice with Emirates Stadium and it is home of Arsenal Football Club. Arsenal is a professional football club based in London, England and they play in the Premi based is called **'The Great Gat** and was written Fitzgerald. imn, and ego en by F. Scot ructBLIP: (K) 🐺 🌲 Ying-VLM: (G) 扬州妙体

Figure 6: Case study of model outputs: Correct answers in green, wrong answers in red, and non-compliant answers in grey. Ying-VLM accurately responds to entitycentric questions and generalizes to unseen Chinese food classification tasks.

572

573

574

575

576

577

578

579

580

581

583

584

585

586

587

589

591

592

593

594

595

Figure 5: Ying-VLM performs stably on unseen instructions, with an average 0.23 ROUGE-L drop.

samples, indicating that VLMs tuned on M2BIT can produce more accurate and engaging responses on the challenging KVQA tasks. Similarly, we supplement the video-language evaluation using GPT-4 as a judge to obtain the question-answering accuracy, following the recent practice in VLM evaluation (Yu et al., 2023). Table 8 shows the results on a 100-sample subset of each video-language task. We find that Ying-VLM<sub>video</sub> achieves the best performance, confirming the effectiveness of multimodal instruction tuning.

548

549

550

551

553

554

556

557

Robustness on Unseen Instructions We further explore the robustness of models on unseen instruc-560 tions, where models are trained on 4 randomly se-561 lected instructions and evaluated with the left 6 562 instructions on each task. As shown in Figure 5, the model performs stably on the unseen instruc-564 tions with a moderate 0.23 ROUGE-L score drop, 565 indicating that it generalizes well on the unseen 566 instructions. Further analysis in Appendix I shows 568 that the instruction diversity also affects the performance, and Ying-VLM shows promising scalability with the increased task number.

571 Qualitative Results We conduct a case study

to provide a more intuitive understanding of instruction-tuned models. The cases are chosen from the held-out ViQuAE and ChineseFoodNet datasets. As shown in Figure 6, Ying-VLM trained with M2BIT provides correct and informative responses to entity-centric questions requiring external world knowledge. In contrast, MiniGPT4 produces an incorrect answer for the stadium question on the left, and InstructBLIP consistently provides concise but less engaging answers. Ying-VLM also generalizes to Chinese inputs, accurately classifying the food image. These cases underscore the importance of instruction tuning dataset quality and validate the value of our M2BIT.

#### 5 Conclusion

In this paper, we introduce M2BIT, a multi-modal bilingual instruction tuning dataset consisting of 2 million instances and 400 task instructions across 40 tasks. We develop Ying-VLM as a proof-ofconcept model to demonstrate the effectiveness of our dataset. Extensive results confirm that Ying-VLM excels in KVQA tasks, exhibits reduced hallucination, and demonstrates superior generalization in unseen video and Chinese tasks.

## Limitations

596

Limitations of Dataset Collection The number of Chinese tasks in our M2BIT is limited, as most high-quality multi-modal resources are English only. In the future, we look forward to incorporating more Chinese V+L tasks in our dataset and exploring machine translation techniques to improve the V+L task coverage in Chinese. Besides, our M2BIT focuses on image-to-text and video-to-text tasks, while more modalities, such as audio (Kim et al., 2019; You et al., 2022; Mei et al., 2023), can be considered further.

**Limitations of Experimental Exploration** In this paper, we curate the M2BIT dataset to provide a resource for developing powerful VLMs and explore 610 the cross-lingual/modality effect of multi-modal 611 instruction tuning. However, there are still under-612 explored setups worth investigating. Promising av-613 enues include exploring improved methodologies 614 for instruction and task selection by taking the in-615 terdependence of different tasks into consideration 616 and exploring the effects of generalization across different languages and modalities. Furthermore, 618 we only adopted the Ziya-13B LLM in our experiments due to its promising bilingual ability. Re-620 cently, many powerful foundation LLMs have been released, such as LLaMA-2 (Touvron et al., 2023a) and Baichuan-2 (Baichuan, 2023). It would also 623 be interesting to perform a comprehensive analysis regarding different model families and scales.

### Ethic Considerations

628

631

633

635

639

642

In line with established practices in language instruction tuning (Mishra et al., 2022; Longpre et al., 2023), our M2BIT dataset has been carefully curated by gathering and unifying NLP datasets from various sources, including academic papers and projects, making them suitable for research purposes. The licenses for the included tasks can be found in Appendix B. However, it should be noted that there are certain tasks for which the license information is not publicly available. We only release the datasets with clear license information to avoid potential problems.

## References

Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katie Millican, Malcolm Reynolds, Roman Ring, Eliza Rutherford, Serkan Cabi, Tengda Han, Zhitao Gong, Sina Samangooei, Marianne Monteiro, Jacob Menick, Sebastian Borgeaud, Andy Brock, Aida Nematzadeh, Sahand Sharifzadeh, Mikolaj Binkowski, Ricardo Barreira, Oriol Vinyals, Andrew Zisserman, and Karen Simonyan. 2022. Flamingo: a visual language model for few-shot learning. *ArXiv*, abs/2204.14198. 644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

- Jacob Andreas, Marcus Rohrbach, Trevor Darrell, and Dan Klein. 2016. Neural module networks. In 2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016, pages 39–48.
- Anas Awadalla, Irena Gao, Josh Gardner, Jack Hessel, Yusuf Hanafy, Wanrong Zhu, Kalyani Marathe, Yonatan Bitton, Samir Gadre, Shiori Sagawa, Jenia Jitsev, Simon Kornblith, Pang Wei Koh, Gabriel Ilharco, Mitchell Wortsman, and Ludwig Schmidt. 2023. Openflamingo: An open-source framework for training large autoregressive vision-language models. *arXiv preprint arXiv:2308.01390*.
- Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, Binyuan Hui, Luo Ji, Mei Li, Junyang Lin, Runji Lin, Dayiheng Liu, Gao Liu, Chengqiang Lu, Keming Lu, Jianxin Ma, Rui Men, Xingzhang Ren, Xuancheng Ren, Chuanqi Tan, Sinan Tan, Jianhong Tu, Peng Wang, Shijie Wang, Wei Wang, Shengguang Wu, Benfeng Xu, Jin Xu, An Yang, Hao Yang, Jian Yang, Shusheng Yang, Yang Yao, Bowen Yu, Hongyi Yuan, Zheng Yuan, Jianwei Zhang, Xingxuan Zhang, Yichang Zhang, Zhenru Zhang, Chang Zhou, Jingren Zhou, Xiaohuan Zhou, and Tianhang Zhu. 2023a. Qwen technical report. *arXiv preprint arXiv:2309.16609*.
- Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. 2023b. Qwen-vl: A frontier large vision-language model with versatile abilities. *arXiv* preprint arXiv:2308.12966.
- Baichuan. 2023. Baichuan 2: Open large-scale language models. arXiv preprint arXiv:2309.10305.
- Rohan Bavishi, Erich Elsen, Curtis Hawthorne, Maxwell Nye, Augustus Odena, Arushi Somani, and Sağnak Taşırlar. 2023. Introducing our multimodal models.
- Ali Furkan Biten, Rubèn Tito, Andrés Mafla, Lluís Gómez i Bigorda, Marçal Rusiñol, C. V. Jawahar, Ernest Valveny, and Dimosthenis Karatzas. 2019. Scene text visual question answering. In 2019 IEEE/CVF International Conference on Computer Vision, ICCV 2019, Seoul, Korea (South), October 27 - November 2, 2019, pages 4290–4300.
- Delong Chen, Jianfeng Liu, Wenliang Dai, and Baoyuan Wang. 2023a. Visual instruction tuning with polite flamingo. *arXiv preprint arXiv:2307.01003*.

756

758

773 774 775

776

780 781

782 783 785

788

789

790

791

792

794

797

798

799

800

801

802

803

804

805

806

807

808

810

811

812

Delong Chen, Jianfeng Liu, Wenliang Dai, and Baoyuan Wang. 2023b. Visual instruction tuning with polite flamingo. arXiv preprint arXiv:2307.01003.

703

704

710

711

713

715

718

719

721

722

725

726

727

729

731

733

736

737

738

739

740

741

742

743

744

745

746

747

748

749

751

- Xi Chen, Josip Djolonga, Piotr Padlewski, Basil Mustafa, Soravit Changpinyo, Jialin Wu, Carlos Riquelme Ruiz, Sebastian Goodman, Xiao Wang, Yi Tay, Siamak Shakeri, Mostafa Dehghani, Daniel M. Salz, Mario Lucic, Michael Tschannen, Arsha Nagrani, Hexiang Hu, Mandar Joshi, Bo Pang, Ceslee Montgomery, Paulina Pietrzyk, Marvin Ritter, A. J. Piergiovanni, Matthias Minderer, Filip Pavetic, Austin Waters, Gang Li, Ibrahim M. Alabdulmohsin, Lucas Beyer, Julien Amelot, Kenton Lee, Andreas Steiner, Yang Li, Daniel Keysers, Anurag Arnab, Yuanzhong Xu, Keran Rong, Alexander Kolesnikov, Mojtaba Seyedhosseini, Anelia Angelova, Xiaohua Zhai, Neil Houlsby, and Radu Soricut. 2023c. Palix: On scaling up a multilingual vision and language model. ArXiv, abs/2305.18565.
  - Xin Chen, Yu Zhu, Hua Zhou, Liang Diao, and Dongyan Wang. 2017. Chinesefoodnet: A large-scale image dataset for chinese food recognition. ArXiv preprint, abs/1705.02743.
  - Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. 2023. Vicuna: An opensource chatbot impressing gpt-4 with 90%\* chatgpt quality.
  - Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale Fung, and Steven Hoi. 2023. Instructblip: Towards general-purpose vision-language models with instruction tuning. ArXiv preprint, abs/2305.06500.
  - Abhishek Das, Satwik Kottur, Khushi Gupta, Avi Singh, Deshraj Yadav, José M. F. Moura, Devi Parikh, and Dhruv Batra. 2017. Visual dialog. In 2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017, pages 1080-1089.
  - Zhengfang Duanmu, Wentao Liu, Zhongling Wang, and Zhou Wang. 2021. Quantifying visual image quality: A bayesian view. Annual Review of Vision Science, 7:437-464.
  - Desmond Elliott, Stella Frank, Khalil Sima'an, and Lucia Specia. 2016. Multi30K: Multilingual English-German image descriptions. In Proceedings of the 5th Workshop on Vision and Language, pages 70-74.
- Haoyuan Gao, Junhua Mao, Jie Zhou, Zhiheng Huang, Lei Wang, and Wei Xu. 2015. Are you talking to a machine? dataset and methods for multilingual image question. In Advances in Neural Information Processing Systems 28: Annual Conference on Neural Information Processing Systems 2015, December 7-12, 2015, Montreal, Quebec, Canada, pages 2296-2304.

- Fabrizio Gilardi, Meysam Alizadeh, and Maël Kubli. 2023. Chatgpt outperforms crowd-workers for textannotation tasks. ArXiv preprint, abs/2303.15056.
- Raghav Goyal, Samira Ebrahimi Kahou, Vincent Michalski, Joanna Materzynska, Susanne Westphal, Heuna Kim, Valentin Haenel, Ingo Fründ, Peter Yianilos, Moritz Mueller-Freitag, Florian Hoppe, Christian Thurau, Ingo Bax, and Roland Memisevic. 2017a. The "something something" video database for learning and evaluating visual common sense. In *IEEE* International Conference on Computer Vision, ICCV 2017, Venice, Italy, October 22-29, 2017, pages 5843-5851.
- Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. 2017b. Making the V in VQA matter: Elevating the role of image understanding in visual question answering. In 2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017, pages 6325-6334.
- Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. Lora: Low-rank adaptation of large language models. In The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022.
- Ting-Hao Kenneth Huang, Francis Ferraro, Nasrin Mostafazadeh, Ishan Misra, Aishwarya Agrawal, Jacob Devlin, Ross Girshick, Xiaodong He, Pushmeet Kohli, Dhruv Batra, C. Lawrence Zitnick, Devi Parikh, Lucy Vanderwende, Michel Galley, and Margaret Mitchell. 2016. Visual storytelling. In Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pages 1233-1239.
- Drew A. Hudson and Christopher D. Manning. 2019. GQA: A new dataset for real-world visual reasoning and compositional question answering. In IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019, pages 6700-6709.
- Justin Johnson, Bharath Hariharan, Laurens van der Maaten, Li Fei-Fei, C. Lawrence Zitnick, and Ross B. Girshick. 2017. CLEVR: A diagnostic dataset for compositional language and elementary visual reasoning. In 2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017, pages 1988–1997.
- Maxime Kayser, Oana-Maria Camburu, Leonard Salewski, Cornelius Emde, Virginie Do, Zeynep Akata, and Thomas Lukasiewicz. 2021. e-vil: A dataset and benchmark for natural language explanations in vision-language tasks. In 2021 IEEE/CVF International Conference on Computer Vision, ICCV 2021, Montreal, QC, Canada, October 10-17, 2021, pages 1224–1234.

813

- 821 822 823 824 825 826 827 828 829 830 831 832 833 832
- 835 836 837 838 839 840 841
- 841 842 843 844 845
- 846 847
- 850 851
- 853 854
- 8
- 8

862

8 8

865

866

- Chris Dongjoo Kim, Byeongchang Kim, Hyunmin Lee, and Gunhee Kim. 2019. AudioCaps: Generating captions for audios in the wild. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), pages 119–132, Minneapolis, Minnesota. Association for Computational Linguistics.
- Jonathan Krause, Justin Johnson, Ranjay Krishna, and Li Fei-Fei. 2017. A hierarchical approach for generating descriptive image paragraphs. In 2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017, pages 3337–3345.
- Paul Lerner, Olivier Ferret, Camille Guinaudeau, Hervé Le Borgne, Romaric Besançon, José G Moreno, and Jesús Lovón Melgarejo. 2022. Viquae, a dataset for knowledge-based visual question answering about named entities. In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 3108– 3120.
- Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. 2023a. Blip-2: Bootstrapping language-image pretraining with frozen image encoders and large language models. *ArXiv preprint*, abs/2301.12597.
- Xirong Li, Weiyu Lan, Jianfeng Dong, and Hailong Liu. 2016. Adding chinese captions to images. In Proceedings of the 2016 ACM on international conference on multimedia retrieval, pages 271–275.
- Xirong Li, Chaoxi Xu, Xiaoxu Wang, Weiyu Lan, Zhengxiong Jia, Gang Yang, and Jieping Xu. 2019. Coco-cn for cross-lingual image tagging, captioning, and retrieval. *IEEE Transactions on Multimedia*, 21(9):2347–2360.
- Yifan Li, Yifan Du, Kun Zhou, Jinpeng Wang, Wayne Xin Zhao, and Ji-Rong Wen. 2023b. Evaluating object hallucination in large vision-language models. *arXiv preprint arXiv:2305.10355*.
- Chin-Yew Lin. 2004. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, pages 74–81.
- Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. 2014. Microsoft coco: Common objects in context. In Computer Vision– ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V 13, pages 740–755. Springer.
- Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. 2023a. Improved baselines with visual instruction tuning.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023b. Visual instruction tuning. *ArXiv preprint*, abs/2304.08485.

Shayne Longpre, Le Hou, Tu Vu, Albert Webson, Hyung Won Chung, Yi Tay, Denny Zhou, Quoc V Le, Barret Zoph, Jason Wei, et al. 2023. The flan collection: Designing data and methods for effective instruction tuning. *ArXiv preprint*, abs/2301.13688. 868

869

870

871

872

873

874

875

876

877

878

879

880

881

882

883

884

885

886

887

888

889

890

891

892

893

894

895

896

897

898

899

900

901

902

903

904

905

906

907

908

909

910

911

912

913

914

915

916

917

918

919

920

921

922

- Ilya Loshchilov and Frank Hutter. 2019. Decoupled weight decay regularization. In 7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019.
- Jiasen Lu, Jianwei Yang, Dhruv Batra, and Devi Parikh. 2018. Neural baby talk. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7219–7228.
- Pan Lu, Swaroop Mishra, Tony Xia, Liang Qiu, Kai-Wei Chang, Song-Chun Zhu, Oyvind Tafjord, Peter Clark, and Ashwin Kalyan. 2022. Learn to explain: Multimodal reasoning via thought chains for science question answering. In *The 36th Conference on Neural Information Processing Systems (NeurIPS)*.
- Kenneth Marino, Mohammad Rastegari, Ali Farhadi, and Roozbeh Mottaghi. 2019. OK-VQA: A visual question answering benchmark requiring external knowledge. In *IEEE Conference on Computer Vision* and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019, pages 3195–3204.
- Minesh Mathew, Dimosthenis Karatzas, and CV Jawahar. 2021. Docvqa: A dataset for vqa on document images. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, pages 2200–2209.
- Xinhao Mei, Chutong Meng, Haohe Liu, Qiuqiang Kong, Tom Ko, Chengqi Zhao, Mark D Plumbley, Yuexian Zou, and Wenwu Wang. 2023. Wavcaps: A chatgpt-assisted weakly-labelled audio captioning dataset for audio-language multimodal research. *arXiv preprint arXiv:2303.17395*.
- Anand Mishra, Shashank Shekhar, Ajeet Kumar Singh, and Anirban Chakraborty. 2019. Ocr-vqa: Visual question answering by reading text in images. In 2019 international conference on document analysis and recognition (ICDAR), pages 947–952. IEEE.
- Swaroop Mishra, Daniel Khashabi, Chitta Baral, and Hannaneh Hajishirzi. 2022. Cross-task generalization via natural language crowdsourcing instructions. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3470–3487.

OpenAI. 2022. Introducing chatgpt.

OpenAI. 2023. Gpt-4v(ision) system card.

Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:27730–27744. Baolin Peng, Chunyuan Li, Pengcheng He, Michel Galley, and Jianfeng Gao. 2023. Instruction tuning with gpt-4. *ArXiv preprint*, abs/2304.03277.

923

924

926

930

931

933

934

935

938

941

942

943

944

945

946

947

949

951

952

955

958

962

963

964

965

966

967

969

970

971

972

973 974

975

976

977

979

- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. 2021. Learning transferable visual models from natural language supervision. In *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24* July 2021, Virtual Event, volume 139 of Proceedings of Machine Learning Research, pages 8748–8763.
- Anna Rohrbach, Lisa Anne Hendricks, Kaylee Burns, Trevor Darrell, and Kate Saenko. 2018. Object hallucination in image captioning. In *Proceedings of the* 2018 Conference on Empirical Methods in Natural Language Processing, pages 4035–4045.
  - Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. 2015. Imagenet large scale visual recognition challenge. *International journal of computer vision*, 115:211–252.
- Christoph Schuhmann, Richard Vencu, Romain Beaumont, Robert Kaczmarczyk, Clayton Mullis, Aarush Katta, Theo Coombes, Jenia Jitsev, and Aran Komatsuzaki. 2021. Laion-400m: Open dataset of clipfiltered 400 million image-text pairs. ArXiv preprint, abs/2111.02114.
- Dustin Schwenk, Apoorv Khandelwal, Christopher Clark, Kenneth Marino, and Roozbeh Mottaghi. 2022.
  A-okvqa: A benchmark for visual question answering using world knowledge. In *Computer Vision–ECCV* 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part VIII, pages 146–162. Springer.
- Aleksandar Shtedritski, Christian Rupprecht, and Andrea Vedaldi. 2023. What does clip know about a red circle? visual prompt engineering for vlms. *ArXiv preprint*, abs/2304.06712.
- Oleksii Sidorov, Ronghang Hu, Marcus Rohrbach, and Amanpreet Singh. 2020. Textcaps: a dataset for image captioning with reading comprehension. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part II 16*, pages 742–758. Springer.
- Amanpreet Singh, Vivek Natarajan, Meet Shah, Yu Jiang, Xinlei Chen, Dhruv Batra, Devi Parikh, and Marcus Rohrbach. 2019. Towards VQA models that can read. In *IEEE Conference on Computer Vision* and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019, pages 8317–8326.
- Alane Suhr, Mike Lewis, James Yeh, and Yoav Artzi. 2017. A corpus of natural language for visual reasoning. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume* 2: Short Papers), pages 217–223.

Ryota Tanaka, Kyosuke Nishida, and Sen Yoshida. 2021. Visualmrc: Machine reading comprehension on document images. In *Thirty-Fifth AAAI Conference* on Artificial Intelligence, AAAI 2021, Thirty-Third Conference on Innovative Applications of Artificial Intelligence, IAAI 2021, The Eleventh Symposium on Educational Advances in Artificial Intelligence, EAAI 2021, Virtual Event, February 2-9, 2021, pages 13878–13888. 980

981

982

983

984

985

986

987

988

989

990

991

992

993

994

995

996

997

998

999

1000

1001

1004

1005

1006

1007

1008

1009

1010

1011

1012

1013

1014

1015

1016

1017

1018

1019

1020

1021

1023

1024

1025

1026

1027

1028

1029

1030

1031

1032

1033

1036

- Tristan Thrush, Ryan Jiang, Max Bartolo, Amanpreet Singh, Adina Williams, Douwe Kiela, and Candace Ross. 2022. Winoground: Probing vision and language models for visio-linguistic compositionality. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5238– 5248.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023a. Llama: Open and efficient foundation language models. *ArXiv preprint*, abs/2302.13971.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023b. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.
- Andreas Veit, Tomas Matera, Lukas Neumann, Jiri Matas, and Serge Belongie. 2016. Coco-text: Dataset and benchmark for text detection and recognition in natural images. *ArXiv preprint*, abs/1601.07140.
- Peiyi Wang, Lei Li, Liang Chen, Dawei Zhu, Binghuai Lin, Yunbo Cao, Qi Liu, Tianyu Liu, and Zhifang Sui. 2023. Large language models are not fair evaluators. *ArXiv preprint*, abs/2305.17926.
- Yizhong Wang, Swaroop Mishra, Pegah Alipoormolabashi, Yeganeh Kordi, Amirreza Mirzaei, Atharva Naik, Arjun Ashok, Arut Selvan Dhanasekaran, Anjana Arunkumar, David Stap, Eshaan Pathak, Giannis Karamanolakis, Haizhi Lai, Ishan Purohit, Ishani Mondal, Jacob Anderson, Kirby Kuznia, Krima Doshi, Kuntal Kumar Pal, Maitreya Patel, Mehrad Moradshahi, Mihir Parmar, Mirali Purohit, Neeraj Varshney, Phani Rohitha Kaza, Pulkit Verma, Ravsehaj Singh Puri, Rushang Karia, Savan Doshi, Shailaja Keyur Sampat, Siddhartha Mishra, Sujan Reddy A, Sumanta Patro, Tanay Dixit, and Xudong Shen. 2022. Super-NaturalInstructions: Generalization via declarative instructions on 1600+ NLP tasks. In Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, pages 5085-5109.
- Zirui Wang, Jiahui Yu, Adams Wei Yu, Zihang Dai, Yulia Tsvetkov, and Yuan Cao. 2021. Simvlm: Simple visual language model pretraining with weak supervision. In *International Conference on Learning Representations*.

Jason Wei, Maarten Bosma, Vincent Y. Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M. Dai, and Quoc V. Le. 2022. Finetuned language models are zero-shot learners. In *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022.* 

1038

1039

1040

1042

1044

1045

1046

1047

1048

1049

1050

1051

1052

1053

1054

1055

1056

1057

1058

1059

1060

1061

1062

1063

1064

1065

1066

1067

1068

1069

1070

1071

1072

1075

1078

1081

1082

1087

1088

1089

1092

- Dejing Xu, Zhou Zhao, Jun Xiao, Fei Wu, Hanwang Zhang, Xiangnan He, and Yueting Zhuang. 2017.
  Video question answering via gradually refined attention over appearance and motion. In *Proceedings of the 2017 ACM on Multimedia Conference, MM 2017, Mountain View, CA, USA, October 23-27, 2017*, pages 1645–1653.
- Jun Xu, Tao Mei, Ting Yao, and Yong Rui. 2016. MSR-VTT: A large video description dataset for bridging video and language. In 2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016, pages 5288–5296.
- Zhiyang Xu, Ying Shen, and Lifu Huang. 2022. Multiinstruct: Improving multi-modal zero-shot learning via instruction tuning. *ArXiv preprint*, abs/2212.10773.
- Antoine Yang, Antoine Miech, Josef Sivic, Ivan Laptev, and Cordelia Schmid. 2021. Just ask: Learning to answer questions from millions of narrated videos. In 2021 IEEE/CVF International Conference on Computer Vision, ICCV 2021, Montreal, QC, Canada, October 10-17, 2021, pages 1666–1677.
- Barry Menglong Yao, Aditya Shah, Lichao Sun, Jin-Hee Cho, and Lifu Huang. 2022. End-to-end multimodal fact-checking and explanation generation: A challenging dataset and models. *ArXiv preprint*, abs/2205.12487.
- Chenyu You, Nuo Chen, Fenglin Liu, Shen Ge, Xian Wu, and Yuexian Zou. 2022. End-to-end spoken conversational question answering: Task, dataset and model. In *Findings of the Association for Computational Linguistics: NAACL 2022*, pages 1219–1232, Seattle, United States. Association for Computational Linguistics.
- Weihao Yu, Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Zicheng Liu, Xinchao Wang, and Lijuan Wang. 2023. MM-Vet: Evaluating large multimodal models for integrated capabilities. arXiv preprint arXiv:2308.02490.
- Zhou Yu, Dejing Xu, Jun Yu, Ting Yu, Zhou Zhao, Yueting Zhuang, and Dacheng Tao. 2019. Activitynet-qa: A dataset for understanding complex web videos via question answering. In *The Thirty-Third AAAI Conference on Artificial Intelligence, AAAI 2019, The Thirty-First Innovative Applications of Artificial Intelligence Conference, IAAI 2019, The Ninth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2019, Honolulu, Hawaii, USA, January 27 - February 1, 2019,* pages 9127–9134.

Rowan Zellers, Yonatan Bisk, Ali Farhadi, and Yejin1093Choi. 2019. From recognition to cognition: Visual<br/>commonsense reasoning. In IEEE Conference on<br/>Computer Vision and Pattern Recognition, CVPR10942019, Long Beach, CA, USA, June 16-20, 2019, pages<br/>6720–6731.1097

1099

1100

1101

1102

1103

1104

1105

1106

1107

1108

1109

1110

1111

1112

1113

1114

1115

1116

1117

1118

1119

1120

1121

1122

1123

1124

1125

1126

1127

1128

1129

1130

1131

1132

1133

1134

1135

1136

1137

1138

1139

1140

1141

1142

- Yuexiang Zhai, Shengbang Tong, Xiao Li, Mu Cai, Qing Qu, Yong Jae Lee, and Yi Ma. 2023. Investigating the catastrophic forgetting in multimodal large language models. *arXiv preprint arXiv:2309.10313*.
- Jiaxing Zhang, Ruyi Gan, Junjie Wang, Yuxiang Zhang, Lin Zhang, Ping Yang, Xinyu Gao, Ziwei Wu, Xiaoqun Dong, Junqing He, Jianheng Zhuo, Qi Yang, Yongfeng Huang, Xiayu Li, Yanghan Wu, Junyu Lu, Xinyu Zhu, Weifeng Chen, Ting Han, Kunhao Pan, Rui Wang, Hao Wang, Xiaojun Wu, Zhongshen Zeng, and Chongpei Chen. 2022. Fengshenbang 1.0: Being the foundation of chinese cognitive intelligence. *ArXiv preprint*, abs/2209.02970.
- Yanzhe Zhang, Ruiyi Zhang, Jiuxiang Gu, Yufan Zhou, Nedim Lipka, Diyi Yang, and Tong Sun. 2023. Llavar: Enhanced visual instruction tuning for textrich image understanding.
- Bo Zhao, Boya Wu, and Tiejun Huang. 2023a. Svit: Scaling up visual instruction tuning. *arXiv preprint arXiv:2307.04087*.
- Haozhe Zhao, Zefan Cai, Shuzheng Si, Xiaojian Ma, Kaikai An, Liang Chen, Zixuan Liu, Sheng Wang, Wenjuan Han, and Baobao Chang. 2023b. Mmicl: Empowering vision-language model with multi-modal in-context learning. *arXiv preprint arXiv:2309.07915*.
- Yinhe Zheng, Guanyi Chen, Xin Liu, and Jian Sun. 2022. MMChat: Multi-modal chat dataset on social media. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 5778–5786.
- Chunting Zhou, Pengfei Liu, Puxin Xu, Srini Iyer, Jiao Sun, Yuning Mao, Xuezhe Ma, Avia Efrat, Ping Yu, Lili Yu, et al. 2023. Lima: Less is more for alignment. *arXiv preprint arXiv:2305.11206*.
- Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. 2023. Minigpt-4: Enhancing vision-language understanding with advanced large language models. *ArXiv preprint*, abs/2304.10592.

### A Datasheet for M2BIT

### A.1 Motivation

**For what purpose was the dataset created?** M2BIT is created to facilitate multi-modal multilingual instruction tuning for large language models.

#### A.2 License

All annotated instructions are licensed under the1143CC-BY 4.0 license. For the licenses of original1144datasets, we refer users to Table 9 for more details.1145

#### A.3 **Maintenance Plan**

1146

We commit to continually updating the dataset and 1147 rectifying any potential errors. Previous versions 1148 of the dataset can still be accessed in the Git history. 1149 Users can submit their questions and suggestions in 1150 the dataset hub, and we will promptly address their 1151 inquiries. We also encourage community contribu-1152 tions to expand the range of datasets by submitting 1153 pull requests to the dataset repository. 1154

#### A.4 Composition 1155

What do the instances that comprise the dataset 1156 represent? (e.g., documents, photos, people, 1157 countries) Our data is provided in JSON format. 1158 Each data instance consists of (1) an instruction 1159 prompt, (2) a list of base64 strings representing 1160 images (3) a task-specific input, such as the ques-1161 tion of the image, (4) a desired output, such as 1162 the answer for the image-related question, and (5)1163 a metadata dictionary for referencing the original 1164 dataset. 1165

How many instances are there in total (of 1166 each type, if appropriate)? The statistics of our 1167 dataset can be found in Table 9. 1168

Does the dataset contain all possible instances 1169 or is it a sample (not necessarily random) of in-1170 stances from a larger set? We tried to transform 1171 the original whole dataset into a unified schema. 1172 However, due to the disk limitation and the cost of 1173 paraphrasing short answers, we chose a randomly 1174 sampled subset from the original dataset to perform 1175 the transformation. 1176

Is there a label or target associated with each in-1177 1178 stance? Yes, the outputs field serves as the label.

Is any information missing from individual in-1179 stances? No. 1180

Are relationships between individual instances 1181 made explicit (e.g., users' movie ratings, social 1182 **network links**)? N/A. 1183

Are there recommended data splits (e.g., train-1184 ing, development/validation, testing)? Yes. We 1185 made the transformation based on the original 1186 dataset split. 1187

Are there any errors, sources of noise, or redun-1188 dancies in the dataset? No. 1189

Is the dataset self-contained, or does it link to or	1190
otherwise rely on external resources (e.g., web-	1191
sites, tweets, other datasets)? Yes.	1192
Does the dataset contain data that might be con-	1193
sidered confidential? No.	1194
Does the dataset contain data that, if viewed di-	1195
rectly, might be offensive, insulting, threatening,	1196
or might otherwise cause anxiety? No.	1197
A.5 Uses	1198
Has the dataset been used for any tasks al-	1199
ready? Yes. We have used the M2BIT dataset to	1200
train a vision-language model, which demonstrates	1201
promising results on knowledgeable VQA tasks	1202
and generalizes well to video-language tasks and	1203
Chinese vision-language tasks. Please see Section	1204
4 of the main paper for details.	1205
What (other) tasks could the dataset be used	1206
for? M2BIT is a useful resource for instruction	1207
tuning studies in the multi-modal field. Future stud-	1208
ies can utilize M2BIT to investigate the influence	1209
of instruction tuning and improve the general per-	1210
formance of vision-language models.	1211
Is there a repository that links to any or all pa-	1212
pers or systems that use the dataset? No.	1213
Is there anything about the composition of the	1214
dataset or the way it was collected and prepro-	1215
cessed/cleaned/labeled that might impact future	1216
<b>uses?</b> The bounding boxes are added to the im-	1217
age with a red rectangle box to inform the model	1218
of regions in interest. For those models with a vi-	1219
sion encoder that is not sensitive to these visual	1220
prompts, the effect of this operation can be lim-	1221
ited. Besides, short answers in some VQA tasks	1222
are paraphrased by ChatGPT, which is designed to	1223
improve the response quality of the model while	1224
potentially impacting the language diversity of the	1225
model.	1226
<b>B</b> Dataset Statistics	1227
Table 9 lists the detailed statistics in our benchmark	1228
and Figure 1 illustrates the unified schema adopted	1229
in our dataset. We collect the dataset license from	1230
PaperWithCode. <sup>5</sup> For datasets under Unknown and	1231

Custom licenses, we suggest the users check the

project page or contact the dataset owner before

1232

1233

1234

usage.

<sup>&</sup>lt;sup>5</sup>https://paperswithcode.com/

Task	Dataset	Used	Train	#samples Val	Test	License
Captioning	MS COCO (Lin et al., 2014) TextCaps (Sidorov et al., 2020) Image-Paragraph-Captioning (Krause et al., 2017)	Yes Yes Yes	566,747 97,765 14,575	25,010 13,965 2,487	25,010 0 2,489	Custom Unknown Custom
Classification	COCO-GOI (Lin et al., 2014) COCO-Text (Veit et al., 2016) ImageNet (Russakovsky et al., 2015) COCO-ITM (Lin et al., 2014) e-SNLI-VE (Kayser et al., 2021) Mocheg (Yao et al., 2022) IQA (Duanmu et al., 2021)	Yes Yes Yes Yes Yes Yes Yes	30,000 118,312 30,000 30,000 20,000 4,991 5,000	2,000 27,550 50,000 5,000 14,339 180 1,000	$0 \\ 0 \\ 5,000 \\ 14,740 \\ 466 \\ 1,000$	Custom Custom Non-commercial Custom Unknown CC BY 4.0 Custom
VQA	VQA v2 (Goyal et al., 2017b) Shapes VQA (Andreas et al., 2016) DocVQA (Mathew et al., 2021) OCR-VQA (Mishra et al., 2019) ST-VQA (Biten et al., 2019) Text-VQA (Singh et al., 2019) GQA (Hudson and Manning, 2019)	Yes Yes Yes Yes Yes Yes Yes	30,000 13,568 39,463 11,414 26,074 27,113 30,001	$\begin{array}{c} 30,000 \\ 1,024 \\ 5,349 \\ 4,940 \\ 0 \\ 0 \\ 5,001 \end{array}$	$0\\1,024\\0\\4,070\\5,734\\0$	CC-BY 4.0 Unknown Unknown Unknown CC BY 4.0 CC BY 4.0
KVQA	OK-VQA (Marino et al., 2019) A-OK-VQA (Schwenk et al., 2022) ScienceQA (Lu et al., 2022) ViQuAE (Lerner et al., 2022)	Yes Yes Yes No	9,009 17,056 12,726 1,190	5,046 1,145 4,241 1,250	0 0 4,241 1,236	Unknown Unknown CC BY-NC-SA CC By 4.0
Reasoning	CLEVR (Johnson et al., 2017) NLVR (Suhr et al., 2017) VCR (Zellers et al., 2019) VisualMRC (Tanaka et al., 2021) Winoground (Thrush et al., 2022)	Yes Yes Yes Yes No	30,000 29,372 25,000 15,000 0	2,000 2,000 5,000 2,500 0	0 0 5,000 5,000 800	CC BY 4.0 Unknown Custom Unknown Unknown
Generation	Visual Storytelling (Huang et al., 2016) Visual Dialog (Das et al., 2017) Multi30k (Elliott et al., 2016)	Yes Yes Yes	5,000 50,000 90,000	4,315 1,000 6,000	4,350 1,000 12,000	Unknown CC By 4.0 Non-commercial
Chinese	FM-IQA (Gao et al., 2015) COCO-Caption CN (Li et al., 2019) Flickr-8k-Caption CN (Li et al., 2016) Chinese Food Classification (Chen et al., 2017) Multimodal Chat (Zheng et al., 2022)	No No No No	164,735 18,341 6,000 0 3,000	75,206 1,000 1,000 0 1,000	0 1,000 1,000 1,100 1,000	Unknown Non-commercial CC By 3.0 Unknown Unknown
Video	Action-Classification (Goyal et al., 2017a) iVQA (Yang et al., 2021) MSVD QA (Xu et al., 2017) ActivityNet QA (Yu et al., 2019) MSRVTT QA (Xu et al., 2017) MSRVTT Captioning (Xu et al., 2016)	No No No No No	2,000 5,994 1,161 3,200 6,513 2,000	2,000 2,000 245 1,800 497 1,000	2,000 2,000 504 800 2,990 1,000	Custom Unknown Unknown Unknown Unknown Unknown

Table 9: Detailed task descriptions and statistics of our instruction tuning tasks, including all datasets in all types of tasks. The column "Used" indicates whether we use this dataset in the instruction tuning stage.

1235

1236

1237

1238

1239

1240

1241

1242

1243

1244

1245

1246

1247

# C Template for Answer Rephrasing

We provide the paraphrase template in Table 10 for querying the ChatGPT to re-write the original short answers, where {Q} and {A} is filled with the question and the answer need to be paraphrased, respectively. We incorporate an example to better inform the model of the paraphrasing tasks. For VQAv2 tasks, we add an extra {Caption} field in the template filled with corresponding captions from the COCO dataset to provide extra context information to help to rephrase. For A-OKVQA tasks, the rationale of each sample is adopted to enrich the final answer.

## D Results on Generation and Classification

For a comprehensive evaluation, we incorporate the results on an image-to-text generation task, i.e., Visual-Dialog and an image classification task ImageNet in Table 11. The results indicate that our M2BIT can provide more balanced coverage of different tasks, helping Ying-VLM achieve better overall performance across various types of tasks.

## **E** Related Work

Our work is motivated by the recent progress in de-1258veloping powerful VLMs based on large language1259

1248

1249

- 1250 1251
- 1252

1253

1254

1255

1256



Figure 7: Top 20 verbs distribution for each task type.

You are an AI visual assistant. Now you are given a question related to an image and a short ground-truth answer. Your task is to transform the ground-truth answer into a natural and convincing response. Make sure the response is accurate, highly relevant to the question, and consistent with the original answer.

Question: Which NASA space probe was launched to this planet in 1989? Answer: Magellan Transformed Answer: NASA sent the Magellan spacecraft to Venus in 1989, which was the first planetary spacecraft launched from a space shuttle. Question: {Q}

Answer: {A} Transformed Answer:

Table 10: Template used to query ChatGPT for answer paraphrasing.

1260models. The efforts can be divided into the ex-1261ploration of the model architecture and the cura-1262tion of high-quality multi-modal instruction tuning1263datasets.

Model	Visual-Dialog	ImageNet
BLIP2-Flan-T5-XXL	6.4	96.0
InstructBLIP	7.0	74.5
MiniGPT4	20.3	41.7
Ying-VLM	14.7	<u>91.9</u>

Table 11: ROUGE-L evaluation on VisualDialog and accuracy score on ImageNet dataset.

Model Architecture of VLMs The architecture 1264 of VLMs usually consists of three key modules, 1265 (i) a vision encoder, (ii) a modality alignment 1266 module, and (iii) an LLM backbone. The vi-1267 sion encoder is responsible for converting raw im-1268 ages/videos into vector representations, which are 1269 then mapped into the embedding space of the LLM 1270 backbone by the alignment module. Finally, the 1271 LLMs decode the response given the multimodal 1272 context. Previous studies mainly employ CLIP 1273 models as the vision encoder due to the transfer-1274 ability of the learned visual representations (Rad-1275 ford et al., 2021). For the LLM backbone, perfor-1276 mant open-sourced LLMs such as LLaMA (Tou-1277 vron et al., 2023b), Vicuna (Chiang et al., 2023) 1278 and Qwen-LM (Bai et al., 2023a) are usually 1279 adopted. Chen et al. (2023c) delves deeper into 1280 the scaling effects of the vision encoder and the 1281 language model, showcasing the great potential 1282 advantages of scaling both modules. There are 1283 many research efforts paid to the design of a better modality alignment module. The pilot studies 1285 (Open)Flamingo (Alayrac et al., 2022; Awadalla 1286 et al., 2023) design a gated cross-attention layer, 1287 where the keys and values in these layers are ob-1288 tained from the vision features while the queries are derived from the language inputs. This design 1290 enables the LLM to produce visual-grounded text 1291 given interleaved context. Another common prac-1292 tice is to treat the visual features as a prefix (Wang 1293 et al., 2021), and train the alignment module via the 1294 construction of the captions. In this way, the align-1295 ment module can become a simple linear projection 1296 layer (Liu et al., 2023b; Zhu et al., 2023) and a two-1297 layer MLP for better capacity (Liu et al., 2023a). 1298 The Q-Former module introduced by BLIP-2 (Li 1299 et al., 2023a) achieves the alignment goal by intro-1300 ducing a suite of query tokens in the transformer 1301 to attend to the visual features. InstructBLIP (Dai 1302 et al., 2023) and MM-ICL (Zhao et al., 2023b) fur-1303 ther incorporate the instruction information during 1304 the alignment process. Fuyu (Bavishi et al., 2023) 1305 presents a new architecture that directly projects 1306



Figure 8: Tasks coverage of M2BIT. We hold out ViQuAE for evaluating the KVQA ability of VLMs and Winogound, as it only provides a small test set. All Chinese V+L and Video-language datasets are excluded during training for cross-lingual/modality effect investigation.

the raw pixels of images into the embedding space of the LLM to perform the following training and decoding.

1309

Multi-modal Instruction Tuning Datasets Cu-1310 rating high-quality multi-modal instruction tuning 1311 datasets has also played a vital role in empowering 1312 VLMs, as evidenced in the progress of develop-1313 ing LLMs (Ouyang et al., 2022; Longpre et al., 1314 2023; Wang et al., 2022). The efforts can be di-1315 vided into two streams. The first uses existing 1316 vision-text benchmarks to create an instruction-1317 tuning dataset (Xu et al., 2022; Dai et al., 2023). 1318 Traditional multi-modal tasks such as VQA and 1319 image captioning are collected, and supplemented 1320 with task instructions annotated by humans or semi-1321 automatically generated, to serve as an instruction-1322 tuning dataset. The second line uses image an-1323 notation tools for providing ChatGPT/GPT-4 with 1324 the textual description of images and then gener-1325 ates a dialog-style dataset (Zhu et al., 2023; Liu 1326 et al., 2023b). For example, the LLaVA dataset (Liu 1327 et al., 2023b) leverages image captions and bound-1328 ing boxes to serve as the textual context of the 1329 original images. The context along with a designed 1330 template is used to prompt the GPT-4 model for 1331 generating instruction tuning samples. Following 1332 this paradigm, Zhao et al. (2023a) further scales 1333 this framework to create 4.2 million samples, and 1334 LLaVAR (Zhang et al., 2023) extends the idea with 1335

OCR results to create text-rich image understanding instructions.

1336

1337

1338

1339

1340

1341

1342

1343

1344

1345

1346

1347

1348

1349

In this paper, we present M2BIT, a multi-modal multilingual instruction tuning dataset with 2 million samples and 400 tasks, providing more comprehensive coverage of multi-modal tasks across different languages and modalities. Besides, we rephrase the short answers of original datasets into natural expressions to improve user experience. Built on the success recipe of VLM architecture, we develop Ying-VLM, which validates the effectiveness of our dataset.

## F Prompt for Zero-Shot Chinese Vision-Language Tasks

In our experiments, all VLMs are fine-tuned ex-1350 clusively using English data. In our preliminary study, we observe that these models tend to gen-1352 erate English responses, even when the input and 1353 instructions are written in Chinese. We introduce 1354 a simple Chinese dialogue context during the zero-1355 shot Chinese Vision-Language Task evaluation for all models, as illustrated in Table 12, Interestingly, 1357 this minor adjustment can encourage models to 1358 produce reasonable Chinese output. We leave the 1359 analysis of instruction-tuned VLM models' multi-1360 lingual capabilities for future research. 1361 <human>:
 请根据我的指示,以及所给的图片,做
 出相应的回答。
 <bot>:
 好的。
 <human>:
 {Instruction}
 {Input}
 <bot>:
 好的。

Table 12: Prompt for promoting Chinese outputs.

### G Template for GPT-4 Evaluation

1362

1363

1364

1365

1369

1370

1372

1373

1374

1375

1376

1377

1378

1380

1381

1382

1384

1385

1386

1388

1389

1390

1391

1392

1393

1394

We adopt the template in Table 13 to query GPT-4 and obtain the evaluation results with FairEval<sup>6</sup> to obtain more stable results. Specifically, each tested instance is a quaternion: (question, reference, response1, response2), where response1 and response2 are two responses from our Ying-VLM and the baseline model, respectively. For each instance, we query GPT-4 to judge which response is of better quality regarding accuracy, relevance and naturalness. We populate the quaternion into the evaluation template to form two query prompts: T(Q=question, R=reference, R1=response1, R2=response2)

and T(Q=question, R=reference, R1=response2, R2=response1). We set the temperature of GPT-4 to 1 and sample three completions for each query prompt. Therefore, each response will receive 6 scores, and we use the average score as the final score for each response. The response with the higher final score is considered the better response. The GPT-4 evaluation incurred a cost of \$20.45 for InstructBlip and \$20.90 for MiniGPT-4.

### H Object Hallucination Metrics

For object hallucination evaluation, we adopt Caption Hallucination Assessment with Image Relevance (CHAIR) proposed by Rohrbach et al. (2018), a metric for evaluating object hallucination in image captioning tasks. Specifically, given the existing objects in the image, CHAIR calculates the proportion of objects that appear in the caption but not the image. CHAIR has two variants, i.e., CHAIR $_I$  and CHAIR $_S$ , which evaluate the halluci-1396nation degree at the object instance level and the1397sentence level, respectively. Formally, these two1398metrics are defined as:1399

$$CHAIR_{I} = \frac{|\{ \text{hallucinated objects } \}|}{|\{ \text{ all mentioned objects } \}|}$$

$$CHAIR_{S} = \frac{|\{ \text{ captions w/ hallucinated objects } \}|}{|\{ \text{ all captions } \}|}.$$
1400

1401

1402

1403

1404

1405

1406

1407

1408

1409

1410

1411

1412

1413

1414

1415

1416

1417

1418

1419

1420

1421

1422

1423

Intuitively, CHAIR<sub>I</sub> denotes the proportion of hallucinated ones in all generated objects, while CHAIR<sub>S</sub> describes the hallucination at the sentence level, i.e., the proportion of generated captions that contain hallucinated objects. We follow the settings adopted in Rohrbach et al. (2018), which only consider 80 objects in the MSCOCO segmentation challenge. Following (Li et al., 2023b), a synonym list (Lu et al., 2018) is used for synonymous word unification in the generated captions to avoid misjudging hallucinated objects.

## I Effect of Instruction Diversity and Robustness

To investigate the influence of instruction diversity, we randomly select 1, 2, 4, and 8 instructions from each dataset, resulting in varied instruction diversity. The other training parameters are consistent with those used in previous experiments on task number investigation. Figure 9 shows that the performance varies with the level of diversity. Specifically, our results suggest that using four instructions per task is sufficient for achieving decent performance.

We investigate the influence of task numbers by 1424 randomly shuffling our tasks and then selecting 1425 a subset to train the model during the instruction 1426 tuning stage. Due to the computational resource 1427 limitation, we set up a maximum of 5k examples 1428 for each task and train all the models for 5k steps 1429 with a batch size of 64. We select 0, 4, 8, 16 and 1430 all 27 tasks for training, and report the individual 1431 ROUGE-L score and the average score. As illus-1432 trated in Figure 10, increasing the number of tasks 1433 greatly improves the results of the generalization 1434 performance. Besides, the performance gain is not 1435 diminished as the task number increases. This is 1436 promising as it indicates that we can continually 1437 improve performance by introducing more tasks 1438 into the training. 1439

<sup>&</sup>lt;sup>6</sup>https://github.com/i-Eval/FairEval

[Question] {Q} [The Start of Reference Answer] {R} [The End of Reference Answer] [The Start of Assistant 1's Answer] {R1} [The End of Assistant 1's Answer] [The Start of Assistant 2's Answer] {R2} [The End of Assistant 2's Answer]

## [System]

We would like to request your feedback on the performance of two AI assistants in response to the user's multimodal question displayed above. We provided no multimodal inputs other than question text, but we provided a reference answer for this question. You need to evaluate the quality of the two responses based on the question and the reference answer.

Please rate the on the follow aspects:

1. Accuracy: whether the candidate's response is consistent with the original answer, this is important as we do not want a misleading result;

2. Relevance: whether the candidate's response is highly relevant to the question and image content;

3. Naturalness: whether the candidate's response is engaging, providing a great communication experience for the user when interacting with the AI visual assistant.

of the two Assistants' responses.

Each assistant receives an overall score on a scale of 1 to 10, where a higher score indicates better overall performance.

Please first provide a comprehensive explanation of your evaluation, avoiding any potential bias and ensuring that the order in which the responses were presented does not affect your judgment. Then, output two lines indicating the scores for Assistant 1 and 2, respectively.

Output with the following format: Evaluation evidence: <evaluation explanation here> The score of Assistant 1: <score> The score of Assistant 2: <score>

Table 13: Template used to query GPT-4 for evaluating the response quality of different models.





Figure 9: ROUGE-L Score changes with the varied number of instructions used for training.

Figure 10: ROUGE-L score increases when models are trained with more instruction tuning datasets.