

---

# KnowBLIP: Large Vision-Language Models meet Interleaved World Knowledge

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Currently, Large Visual Language Models (LVLMs) have integrated various tasks,  
2 and some “any2any” models have emerged, enabling support for diverse modalities  
3 for input and various output formats, thereby expanding the impact of LVLMs.  
4 However, we note that existing models remain incomplete and distant from real-  
5 world applicability, primarily due to two key factors: partial functionality and lack  
6 of connection with world knowledge. Firstly, we stress the importance of a com-  
7 prehensive LVLM that accommodates various input and output paradigms, which  
8 should be interleaved, multi-turn, and multi-lingual. However, there is no multi-  
9 modal framework capable of fully encompassing all functionalities. Moreover,  
10 textual knowledge often lacks substantial connections with the visible objective  
11 world, resulting in visual and language alignment limited to the semantic level,  
12 detached from objective facts. This limitation leads to models generated primarily  
13 through imagination rather than a grounded understanding of the world. To bridge  
14 the gap between LVLMs and the visible real world, we propose a novel multi-modal  
15 framework, KnowBLIP, designed to support a comprehensive, world-grounded,  
16 “any2any” paradigm. Specifically, we have collected a dataset, IWK-500k, of multi-  
17 modal objective knowledge derived from reasoning on multi-modal knowledge  
18 graphs, characterized by interleaved, multi-turn, and multi-lingual attributes. Fur-  
19 thermore, we have developed the KnowBLIP, capable of accommodating various  
20 input and output paradigms based on IWK-500k dataset. Extensive experiments  
21 have demonstrated the effectiveness of our approach, and both the dataset and code  
22 will be made publicly available<sup>1</sup>.

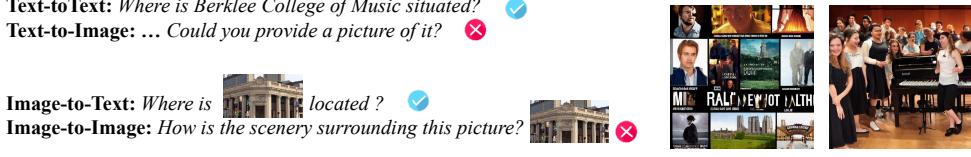
## 1 Introduction

23 With the development of the internet and social media, data across various modalities has seen explo-  
24 sive growth. Harnessing this vast amount of data, *unimodal* models have achieved groundbreaking  
25 successes in their respective domains, such as the text-oriented ChatGPT<sup>2</sup> and the vision-oriented  
26 ViT[1]. However, their lack of exposure to data from other modalities renders them fall short in  
27 multi-modal tasks. To address this *multimodal* challenge, Large-scale Vision-Language Models  
28 (**LVLMS**) like GPT4-V and LLaVA[2] have emerged. By understanding, integrating, and aligning  
29 data from different modalities, LVLMS have demonstrated remarkable performance across numerous  
30 multimodal tasks.

32 Along this line, we revisit the aforementioned pioneer LVLMS: MiniGPT-4[3] and LLaVA[2] (*Visual*  
33 *Question Answering*), Shikra[4] and LiSA[5] (*Grounding and Segmentation*), MiniGPT-5[6] and

<sup>1</sup>[https://drive.google.com/drive/folders/1ZgfApnIhzXg8n3R\\_kwkiix5C0K1v6UP8](https://drive.google.com/drive/folders/1ZgfApnIhzXg8n3R_kwkiix5C0K1v6UP8)

<sup>2</sup><https://openai.com/blog/chatgpt>



(a) Interleaved Question revolves around World Knowledge. (b) Generated images by GILL.

Figure 1: The four ideal paradigms exist within LVLMs. However, most baselines are unable to generate images. Even models that can generate images, like GILL, struggle to obtain objective results, as demonstrated in the case of “Boston”.

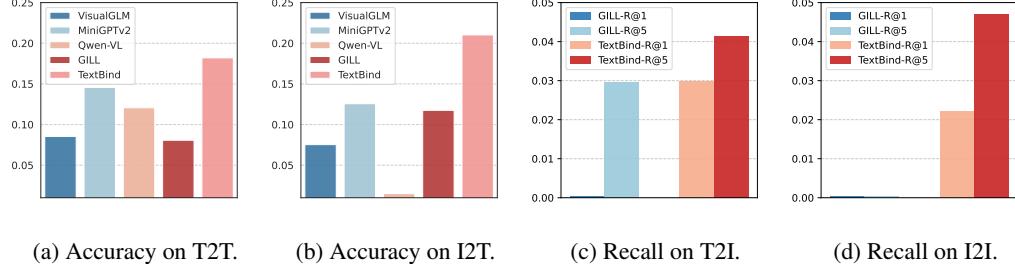


Figure 2: The zero-shot testing of world knowledge answering on existing models. The accuracy of knowledge and the similarity between generated images and real images are insufficient.

34 DreamLLM[7] (*Text-to-Image Generation*). Despite these LVLMs advancing the field of multimodal  
35 tasks, it’s regrettable that they are “**Non-any2any**” **multimodal models** (i.e., *not supporting arbitrary*  
36 *input and output modalities*) but only a subset of “any2any” modes (e.g., T2I, I2T) as shown in  
37 Table1. Considering that humans communicate in an “any2any” mode in real life, the “any2any”  
38 paradigm is essential path toward multimodal general models. In light of this, researchers have  
39 focused and proposed several “any2any” models, such as TextBind[8], NExT-GPT[9], to aid LVLMs  
40 in their real-world practices. Attracted by the powerful capabilities of existing “any2any” models,  
41 we employed these models to accomplish multimodal tasks, but discover a common issue — *the*  
42 *generated content do not align with the objective facts of the real world* — we refer to as “**A gap**  
43 **with world knowledge**”. For example, as shown in Figure1, when addressing the question “Where  
44 is the Berklee College of Music located?”, the text generated by VisualGLM, the image produced  
45 by GILL[10], or the text-image content generated by TextBind[8], are all incorrect. Further, to  
46 verify the “A gap with world knowledge” issue indeed exists in current “any2any” baselines, we  
47 meticulously conducted four multimodal world knowledge QA tasks. As illustrated in Figure2, *the*  
48 *poor performance* of text generation ( $\leq 25\%$ ) and image generation ( $\leq 5\%$ ) tasks validate that  
49 *existing “any2any” models rely on imagining rather than understanding world to generate content*.  
50 Therefore, we underscore the importance of integrating world knowledge into “any2any” multimodal  
51 models, thereby facilitating the generation of content that aligns with objective facts.

52 Facing the challenges of “*Non-any2any*” and “*A gap with world knowledge*”, in this paper, we  
53 develop a comprehensive multimodal benchmark to foster a close connection between “any2any”  
54 LVLMs and the real world: 1)-we collected a multimodal objective knowledge dataset **IWK-500**, via  
55 multimodal knowledge graph (MMKG) reasoning, as depicted in radar chart 5. Notably, to bring  
56 this dataset closely mimic real-world QA scenarios, we added three key attributes: *i*)-*Interleaved*:  
57 both queries and answers contain text and images; *ii*)-*Multi-turn*: the query-answer pairs support the  
58 form of multi-turn dialogues; *iii*)-*Multilingual*: the query-answer pairs support multiple languages.  
59 2)-we developed an “any2any” LVLM **KnowBLIP** driven by autoregressive text generation and  
60 contrastive-based image retrieval tasks, as illustrated in the framework diagram 6.

61 Our contributions are summarized as follows:

- 62 We pioneered a comprehensive multimodal benchmark with the creation of the IWK-500  
63 dataset, a **first-of-its-kind** that supports “*any2any*”, *interleaved*, *multi-turn*, *multi-lingual*, and  
64 *incorporates world-knowledge aspects*. Leveraging this, we introduced “any2any” LVLMs,  
65 KnowBLIP, designed for *practical use and aligned with human needs*.
- 66 Over **six months**, we meticulously prepared prompts and engaged **10 annotators** to compile  
67 the IWK-500 dataset. KnowBLIP then outperformed **8 leading baselines**, such as GPT-4V

Table 1: Comparision between our collected datasets and other existing datasets. Among them, the “T2I” means Text-to-Image, aka image generation.

Dataset	Type	Multi-Trun	Interleaved	World Knowledge	Multi-lingual	Image volume	Text volume
ShareGPT52K	T2T	✓	✗	✗	✓	-	90K
Alpaca-GPT4[11]	T2T	✗	✗	✗	✓	-	52K
Open Assistant	I2T	✗	✗	✗	✓	-	84K
COCO Caption[12]	I2T	✗	✗	✗	✓	164K	1M
CC12M[13]	I2T	✗	✗	✗	✗	12M	12M
LAION-COCO[14]	I2T	✗	✗	✗	✗	600M	300M
VQAv2[15]	I2T & T2T	✗	✗	✗	✗	265K	1.4M
OK-VQA[16]	T2I & T2T	✗	✗	✗	✗	14K	14K
GQA[17]	T2I & T2T	✗	✗	✗	✗	113K	22M
Visual Genome[18]	T2I & T2T	✗	✗	✗	✗	108K	1.77M
ScienceQA[19]	T2I & T2T	✗	✗	✗	✗	10k	21K
TextVQA[20]	T2I & T2T	✗	✗	✗	✗	28k	45K
LLaVA-instruct-150K[2]	T2I & T2T	✓	✗	✗	✗	82K	150K
MMDialog[21]	T2I & I2T & T2T & I2I	✓	✓	✗	✗	1.53M	1.08M
Ours	T2I & I2T & T2T & I2I	✓	✓	✓	✓	650K	1.2M

and Qwen-VL, on the xxxx task, setting a **new record for state-of-the-art performance** in downstream multimodal tasks.

## 2 Related Work

### 2.1 Large Vision-Language Models

Recent progresses of computational resources has greatly facilitated the research into large-scale foundational models incorporated with multi-modal learning [22]. Powered by open sourcing large language models such as LLaMA[23] and Vicuna[24], LVLMs understand and generate diverse content in a more comprehensive way by integrating information from different modalities, such as text, images, and audio. Qwen-vl[25], VisualGLM[26] and MiniGPT-4[3] take a step forward in this field, allowing users to interact with these intelligence with images and texts as prompts. All of them share the same two training phases, i.e., pre-trained feature alignment and instruction fine-tuning[27], to help the model to comprehend the format of instruction input. With the assist of Diffusion model[28], LVLMs are extended to image generation task. GILL[10], TextBind[8] translates hidden representations of text into the embedding space of the visual models, enabling LVLMs to leverage the strong text representations of the LLM for visual outputs. However, all of aforementioned LVLMs suffer from lacking connection with the visible world. Consequently, we mainly conduct the experiments on these models to solve this problem in our paper.

### 2.2 Multi-modal Knowledge Graph

In this paper, we will utilize a multi-modal knowledge graph (MMKG) to establish connections between text and the physical world[29]. A knowledge graph (KG) provides a structured representation of knowledge, typically in the form of triples consisting of (head, predicate, tail). In this schema, the head and tail denote entities, while the predicate indicates the type of relationship between them. In the context of MMKG, entities can also encompass images and other modalities, for instance (France, hasImage, <IMG>), where <IMG> represents the flag of France. This framework facilitates the seamless linkage of abstract textual concepts to tangible real-world entities, enabling the model to transcend reliance on abstract notions and establish concrete associations with the physical realm.

## 3 Proposed Method

To facilitate the LVLM in acquiring objective knowledge of the visible real world, we collect a comprehensive dataset by reasoning on MMKG. Moreover, in Section 3.2, we will elaborate on how our framework achieves comprehensive and world-grounded performance.

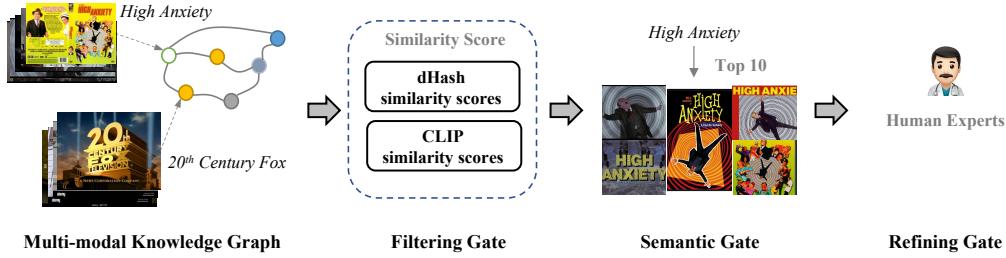


Figure 3: The pipeline for selecting an image from Multi-modal Knowledge Graph.

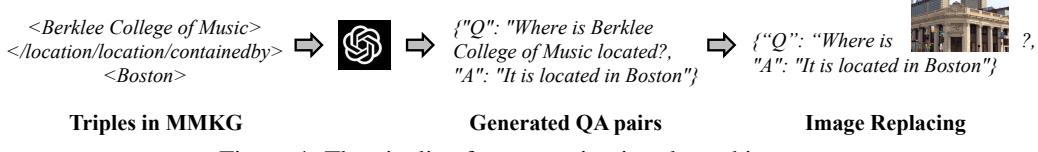


Figure 4: The pipeline for generating interleaved images.

### 98 3.1 IWK-300K Collection

99 We generated a comprehensive dataset of factual knowledge using the multi-modal knowledge graph  
100 (MMKG), which is built upon FB15K[29] to create interleaved data in diverse forms, including  
101 Text-to-Text (T2T), Image-to-Text (I2T), Text-to-Image (T2I), and Image-to-Image (I2I). In FB15K,  
102 knowledge is structured into triples (head entity, relation, tail entity), where each entity has multiple  
103 images. To select the most representative image, as depicted in Figure 3, three gates are designed: the  
104 filter gate, semantic gate, and refining gate. The filter gate utilizes dHash and CLIP[30] to compute  
105 image similarity, thereby retaining images with minimal repetition. Subsequently, the semantic gate  
106 selects the most relevant images to the textual entity. Finally, expert selection is employed to choose  
107 an image to represent the entity. Following rigorous endeavors, we successfully obtained a dataset of  
108 298,957 triplets, including 14,505 entries.

109 As shown in Figure 4, in order to produce QA pairs suitable for LVLMs training, we employed  
110 GPT-4[31] to fashion a natural language description of triples, as well as polished QA pairs derived  
111 from the subsequent prompt.

112 **Input:** Please generate QA pair for the description "Berklee College of Music is located in Boston"  
113 and return them in JSON format. Such as ...  
**Output:** {"Q": "Where is Berklee College of Music located?", "A": "It is located in Boston" } ...

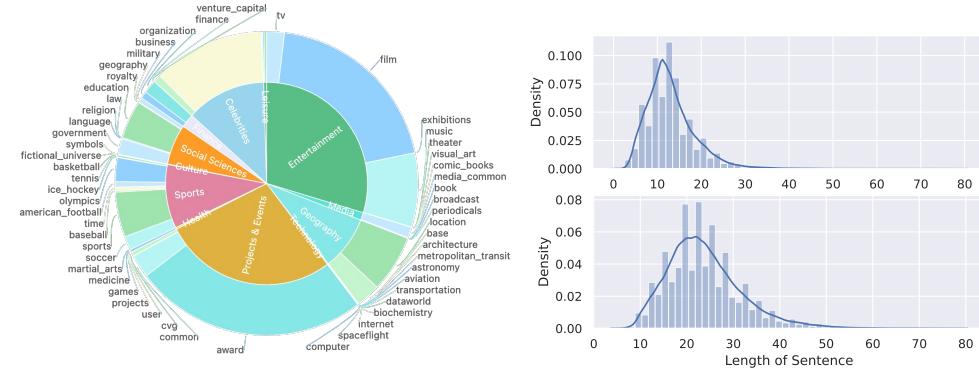
114 Thereafter, entity concepts are randomly replaced by corresponding images with a 50% probability,  
115 forming pairs like (<IMG> of head entity, relation, tail entity).

116 Furthermore, recognizing the significance of multi-turn conversations in enhancing the few-shot[32]  
117 capability of LVLMs, we have curated multi-turn data by consolidating identical entities. Notably,  
118 every output produced by GPT-4 will undergo meticulous scrutiny and refinement by 10 human  
119 experts to guarantee the precision of knowledge. Ultimately, we acquired a multi-turn, interleaved  
120 dataset comprising 650K images and 1.2M text volumes, named IWK-300k.

121 **Statistics** We have compiled fundamental attributes of IWK-300k, summarized in Table  
122 2. The dataset comprises 298,957 rounds of  
123 conversation, averaging 3.49 rounds per complete conversation. Additionally, we conducted  
124 an analysis of the occurrence of "wh" words in QA pairs to ensure question diversity. Notably,  
125 owing to the rich location attributes in MMKGs, "who" questions are predominant, totaling 136,862 occurrences. From a visual  
126 standpoint, the average number of images per

Table 2: Statistics of IWK-300k dataset.

Total datas	298,957
Total entitys	14,505
Avg. #data/entity	20.61
Avg. #image/entity	61.23
Avg. #turns in conversations	3.49
Total number of QA pairs	
"What"	124,111
"When"	9,271
"Which"	5,669
"Who"	136,862
"Where"	21,948
"Is/Are/Did"	1,096



(a) The distribution of topics.

(b) The distribution of sentence length.

Figure 5: The statistics in our IWK-300k dataset.

entity is 61.23, with each entity featuring in an average of 20.61 questions. Lastly, we compared the performance of GPT-4 before and after refinement and enhancement, illustrated in Figure 5. The length of questions per round increased from an average of 10 words to over 18 words. More details are disclosed in the Appendix A.

**Diversity** To evaluate sentence diversity within the IWK-300k datasets, we perform a topic frequency analysis of world knowledge present in the IWK-300k corpus. As illustrated in Figure 5, we initially quantify the total occurrence of each topic among entities. Among them, there are 63 distinct topic types in total. Through human classification and summarization, these topics are categorized into 12 areas, encompassing nearly comprehensive and diverse real-world domains.

Formally, we define this dataset constructed as  $D=\{(X_u^1, Y_m^1), \dots, (X_u^n, Y_m^n)\}$ , where the  $X_u^i$  represent the i-th instruction from user while the  $Y_m^i$  represent the i-th answer from model. Meanwhile, a interleaved input  $X_u^i$  is defined as  $X_u^i=(x_t^1, x_t^2, x_i^1, \dots, x_t^n, x_i^k)$ , where  $x_t^*$  represents the text token of multi-turn conversations ,  $x_i^*$  represents the visual part interleaving in the conversations. The n denotes the total number of text tokens, while the k denotes the total number of images in conversations.

### 3.2 KnowBLIP Framework

The overall training framework is depicted with an illustration in Figure 6. In line with most LVLMs, our framework employs ViT-L/14 from BLIP-2[33] as the visual encoder, a pretrained Q-Former as the visual adapter, and LLaMA-2-7B as the LLM. During training, the adapter undergoes tuning[34] while the ViT and LLM remain frozen.

On the training function, we extended the autoregressive loss  $\mathcal{L}_{AR}$  of the language model. Formally, given a question in a single turn conversation  $X_u^i=(x_t^1, x_t^2, x_i^1, \dots, x_t^n, x_i^k)$ , the model expects to output the answer  $Y_m^i=(y_t^1, y_t^2, \dots, y_t^k, y_t^n)$ .

$$\mathcal{L}_{AR} = -\log p(y_t^k | X_u^i, y_t^j, j < k). \quad (1)$$

To bolster the model’s capability in generating objective knowledge, particularly in visual generation, we define the image generation as a retrieval task. This necessitates the LVLM to not only output textual knowledge (i.e. Boston), but also to “generate” Boston images via retrieval process. Specifically, we introduce two special tokens: the knowledge token and the retrieval token, represented as  $\langle KNO \rangle$  and  $\langle RET \rangle$ , for knowledge generation and image retrieval, respectively.

For knowledge generation task, we use  $\langle KNO \rangle$  token as an anchor to mark the world knowledge entity and compute a knowledge loss  $\mathcal{L}_{KNO}$ . Specifically, a question and answer in a single turn conversation are represented as  $X_u^i=(x_t^1, x_t^2, x_i^1, \dots, x_t^n, x_i^k)$  and  $Y_m^i=(y_t^1, y_t^2, \dots, y_{kno}, y_t^n)$ , where  $y_{kno}$  denotes the marked world knowledge entity  $\langle KNO \rangle$ . This  $\langle KNO \rangle$  token will be computed loss with the correctly world knowledge entity:

$$\mathcal{L}_{KNO} = -\log p(y_{kno} | X_u^i, y_t^j, j < k), \quad (2)$$

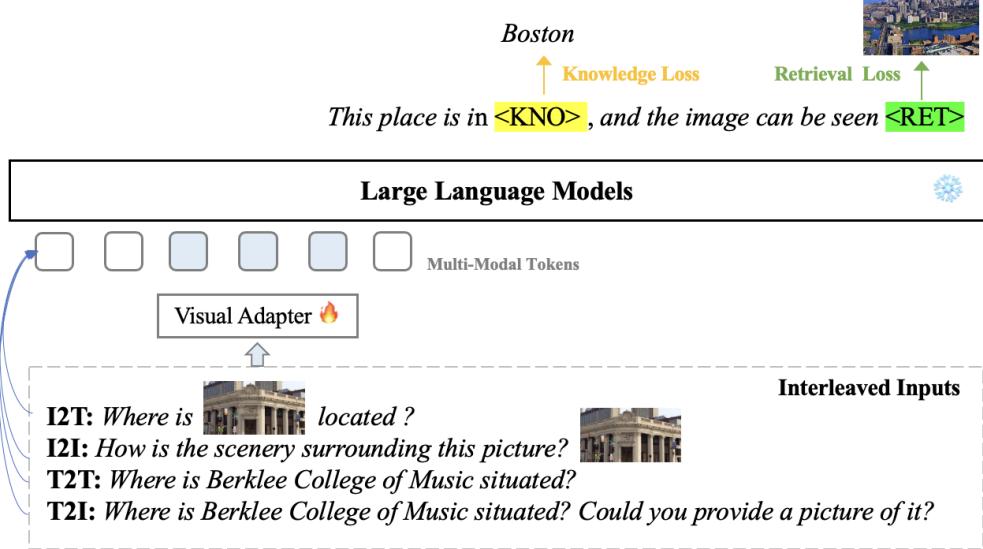


Figure 6: The overview of KnowBLIP for interleaved inputs, which consists of three losses: regression loss  $\mathcal{L}_{AR}$ , knowledge loss  $\mathcal{L}_{KNO}$ , and retrieval loss  $\mathcal{L}_{RET}$ .

165 Through this loss, the model will be guided to focus on generating correct knowledge.  
 166 For image retrieval task, we introduce a specific token <RET>, which serves to “generate” objective  
 167 images. In the training phase, this token is affixed to the textual context in order to tag a objective  
 168 image. Correspondingly, during the inference stage, if the model generates a <RET> token, our  
 169 system proceeds with the retrieval process. The computation of the retrieval loss  $\mathcal{L}_{RET}$ , hinges on  
 170 the ability of the representation of the <RET> token to retrieve images. Similar to CLIP, we employ  
 171 contrastive loss for this task. Specifically, the embeddings of all potential images are acquired via  
 172 ViT-L/14 and then mapped onto a unified space through the projection matrix,  $W_{img} \in \mathbb{R}^{768 \times p}$ ,  
 173 where  $p$  is the linear dimension of the <RET> vector. Subsequently, the image feature  $E_{img}^i$  and the  
 174 token feature  $E_{ret}^i$  are used to calculate contrastive loss[35],

$$\mathcal{L}_{RET} = -\log \frac{\exp(\text{sim}(E_{ret}, E_{img}^i)/\tau)}{\sum_{j=1}^n \exp(\text{sim}(E_{ret}, E_{img}^j)/\tau)}, \quad (3)$$

$$\text{sim}(E_{ret}, E_{img}^i) = \frac{(E_{ret})(E_{img}^i)^T}{\|E_{ret}\| \|E_{img}^i\|}, \quad (4)$$

175 where  $i$  represent the matching image and  $j$  denotes negative images that is randomly sampled,  $\tau$  is  
 176 the temperature coefficient.

177 Finally, our training function  $\mathcal{L}$  consists of three parts, namely regression loss  $\mathcal{L}_{AR}$ , knowledge loss  
 178  $\mathcal{L}_{KNO}$ , and retrieval loss  $\mathcal{L}_{RET}$ ,

$$\mathcal{L} = \mathcal{L}_{AR} + \lambda_1 \mathcal{L}_{KNO} + \lambda_2 \mathcal{L}_{RET}, \quad (5)$$

180 where  $\lambda_1$  and  $\lambda_2$  are balance parameters.

## 181 4 Experiments

182 Our KnowBLIP is a LVLM designed to incorporate world knowledge, and thus our experiments  
 183 primarily center on assessing its capacity to produce world knowledge images and texts. In order  
 184 to compare the effectiveness of our model with existing baselines in terms of world knowledge, we  
 185 carried out experiments on four subsets of IWK-300k (e.g., Text2Image).

### 186 4.1 Datasets

187 Our datasets are split into 3 non-overlapping subsets, where 0.6, 0.2 and 0.2 are used for training,  
 188 validation and testing. Specifically, for one triplet that has four formats in different formats of

Table 3: Zero-shot testing of baselines on IWK-300k. Among these, “–” indicates not applicable, as the model is unable to execute the relevant task.

	Text2Text	Image2Text	Text2Image R@1	Text2Image R@5	Image2Image R@1	Image2Image R@5
VisualGLM[26]	8.54%	7.53%	–	–	–	–
Qwen-VL[25]	14.57%	12.56%	–	–	–	–
MiniGPTv2[37]	12.06%	1.50%	–	–	–	–
GILL[10]	8.07%	11.75%	0.00%	2.95%	–	–
TextBind[8]	18.20%	21.02%	2.98%	4.13%	2.21%	4.70%
LLaVA-1.5-7B[38]	6.69%	10.20%	–	–	–	–
CogVLM-chat-v1.1[39]	14.45%	17.41%	–	–	–	–
GPT-4V	61.53%	55.88%	–	–	–	–
Gemini-1.5-pro[40]	56.41%	41.53%	–	–	–	–

Table 4: Performance comparison of baselines after training on IWK-300k. Among these, “–” indicates not applicable, as the model is unable to execute the relevant task. **Bold** represents optimal performance, while underlined represents suboptimal.

	Text2Text	Image2Text	Text2Image R@1	Text2Image R@5	Image2Image R@1	Image2Image R@5
VisualGLM[26]	50.72%	22.50%	–	–	–	–
Qwen-VL[25]	51.69 %	48.74%	–	–	–	–
MiniGPTv2[37]	40.62%	30.57%	–	–	–	–
GILL[10]	47.38%	70.03%	6.79%	24.30%	–	–
TextBind[8]	59.38%	64.29%	5.66%	6.30%	3.72%	6.79%
<b>KnowBLIP (Ours)</b>	<b>66.73%</b>	<b>70.60%</b>	<b>29.71%</b>	<b>78.24%</b>	<b>17.47%</b>	<b>64.38%</b>

189 question-answer pair, we would split its different formats into training, validation and testing sets,  
190 which ensures the world knowledge of this triplet having trained.

## 191 4.2 Evaluation Metrics

192 The evaluation is carried out in two aspects: the accuracy of the generated knowledge and the quality  
193 of the generated images. These metrics are calculated using <KNO> token and <RET> token. In our  
194 framework, the <RET> token will be utilized to retrieve images, while in the baselines, the embedding  
195 of the final generated image will be used for retrieval. For knowledge accuracy, the accuracy will be  
196 employed for T2T and I2T tasks to ascertain whether the response aligns with the world knowledge  
197 entity. Regarding image, recall@1 and recall@5 are utilized to assess whether the images obtained  
198 from I2I and T2I tasks are sufficiently objective.

## 199 4.3 Implementation Details

200 In the retrieve process, feature linear mapping  $W_{RET}$  and linear mapping  $W_{img}$  are 768-dimensional  
201 vectors. We use the Adam optimizer with  $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$  for all gradient-based methods, where  
202 the mini-batch size and the learning rate were searched and are set up to 16 and  $10^{-4}$ , respectively.  
203 We iteratively trained 20K steps on a server equipped with 4 A100 GPUs[36].

## 204 4.4 Overall Performance Comparison

205 In this subsection, we begin by conducting robustness experiments (i.e. zero-shot testing) on the  
206 baselines, followed by an overall comparison of the trained baselines.

207 During zero-shot testing, we evaluate a wide range of models, and metrics are not calculated if the  
208 corresponding subtasks cannot be executed. Notably, in GILL, we use the <IMG> token for retrieval,  
209 while in TextBind[8], the images generated by stable diffusion are utilized. As depicted in Table  
210 3, the following observations are made: 1) Virtually no model inherently supports all four tasks,  
211 underscoring the value of developing a comprehensive LVLM. 2) The majority of models exhibit  
212 limited knowledge abilities, indicating that baselines often produce outputs without comprehensive  
213 world knowledge.

Table 5: Ablation analysis on different subtasks. Among these, “–” indicates not applicable, as the model is unable to execute the relevant task. **Bold** represents optimal performance, while underline represents suboptimal.

	Text2Text	Image2Text	Text2Image R@1	Text2Image R@5	Image2Image R@1	Image2Image R@5
VisualGLM[26]	56.00%	15.46%	–	–	–	–
Qwen-VL[25]	74.19%	29.47%	–	–	–	–
MiniGPTv2[37]	80.99%	17.51%	–	–	–	–
GILL[10]	41.61%	61.12%	<u>39.77%</u>	<u>71.92%</u>	–	–
TextBind[8]	<u>91.00%</u>	64.16%	6.86%	9.97%	3.88%	7.89%
<b>KnowBLIP (Ours)</b>	<b>91.29%</b>	<b>71.32%</b>	<b>25.13%</b>	<b>72.54%</b>	<b>24.67%</b>	<b>70.88%</b>

214 In overall comparison, we selected a diverse range of state-of-the-art models. These encompass three  
 215 baselines that are proficient in text generation, Qwen-VL[25], VisualGLM[26], MiniGPTv2[37]; two  
 216 models for image generation: GILL, TextBind[8]. Specifically, for text generation baselines, we only  
 217 conduct experiments on T2T and I2T tasks. For Image generation frameworks, we train TextBind[8]  
 218 within the origin code and retrieve images via the similarity of images generated by stable diffusion.

219 The experimental results are presented in Table 4, and the following observations can be made:

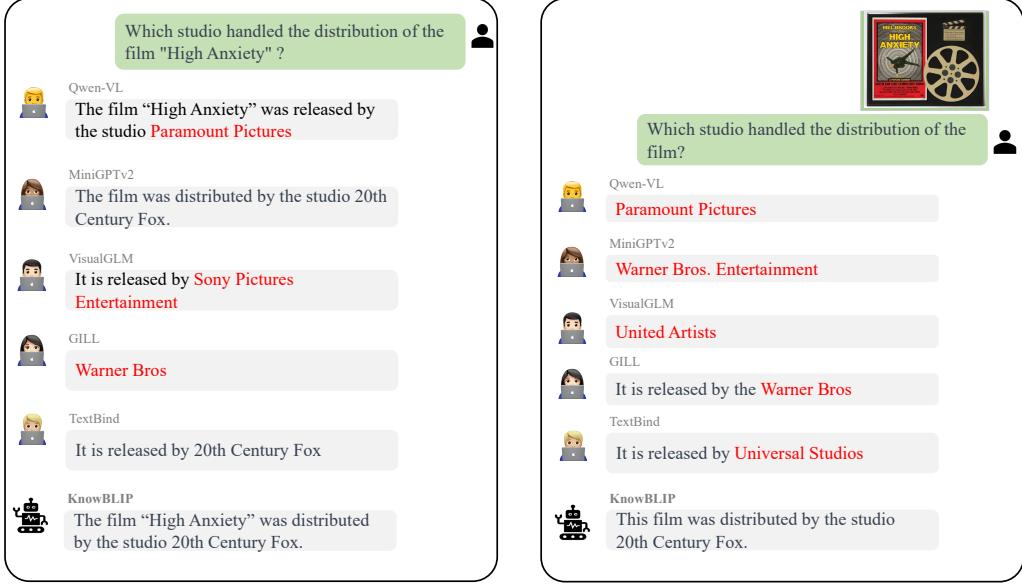
- 220 • When training IWK-300k on each model, both models demonstrated a noticeable improvement com-  
 221 pared to the zero-shot setting. This improvement suggests that our dataset effectively incorporates  
 222 world knowledge into this models.
- 223 • In the text generation task, all models showed a significant improvement, indicating that textual  
 224 generation tasks are relatively easy for existing LVLMs. However, it was also observed that the  
 225 I2T task lagged behind the T2T task, highlighting the challenge of utilizing LVLMs to understand  
 226 images of world knowledge entities.
- 227 • In the image generation task, the model TextBind performed poorly in both the zero-shot and  
 228 training settings. In the TextBind experiment setting, it encodes the generated text between two  
 229 special tokens using the text encoder of the SD model to generate an image. However, most of the  
 230 generated texts are special phrases in IWK-300k, such as the names of celebrities or cities. The  
 231 images generated by SD still appear to be disconnected from the world knowledge entity, which  
 232 accounts for the poor performance of TextBind.
- 233 • Our **KnowBLIP** exhibited substantial improvement compared to TextBind in all three tasks,  
 234 especially T2I and I2I. Through the retrieval method, we are able to provide real-world images  
 235 instead of purely imaginary ones as much as possible.

## 236 4.5 Ablation Study

237 In order to thoroughly investigate our framework and dataset, we conducted ablation studies by  
 238 training four subsets individually in all models.

239 The results of the ablation studies are presented in Table 5, and the following observations were  
 240 made:

- 241 • Training individually on each subset task led to a certain degree of improvement in both models  
 242 compared to the zero-shot setting, indicating that our dataset is capable of injecting world knowledge  
 243 into the models for each individual task.
- 244 • A comparison with the results of mixed training tasks revealed a slight increase in the performance  
 245 of the I2T task during individual training. Conversely, the T2T task showed a minor decrease  
 246 in individual training. These observations suggest that mixed training enhances the ability to  
 247 comprehend input images and better aligns world knowledge entities with the input images.
- 248 • However, this setting also has a negative impact on how well models master world knowledge  
 249 through text-only. Similarly, similar performance differences were observed between I2I and T2I  
 250 tasks.



(a) Text2Text task.

(b) Image2Text task.

Figure 7: The visualization of several models in text generation tasks.

251 **4.6 Visualization**

252 To thoroughly examine our framework’s processes during execution, we have employed visualizations  
 253 to elucidate, as shown in Figure 7 and Figure 8. 1) For the text generation task, all models have  
 254 the capability to understand the content of the text or image and perform the basic task of question  
 255 and answering. This is evident in their ability to provide appropriate responses when asked about a  
 256 film distribution company, such as “Warner Bros” or “Sony Pictures Entertainment”. However, their  
 257 performance is unsatisfactory when faced with world knowledge in the <KNO> slot that requires a  
 258 response. This indicates that while all models can easily master the correct response format, they  
 259 struggle with the correct world knowledge entity slot. Similarly, GILL and TextBind make mistakes  
 260 in image generation tasks, leading to a chain of errors in the retrieved or generated images. 2) In the  
 261 image generation task, we observed that the world knowledge of the <KNO> slot is correct in both  
 262 TextBind and KnowBLIP in the T2I task. However, TextBind generates an image that matches the  
 263 literal meaning (such as a strange **fox**, possibly caused by the modifiers of **20th Century**). More  
 264 examples are available in the appendix.

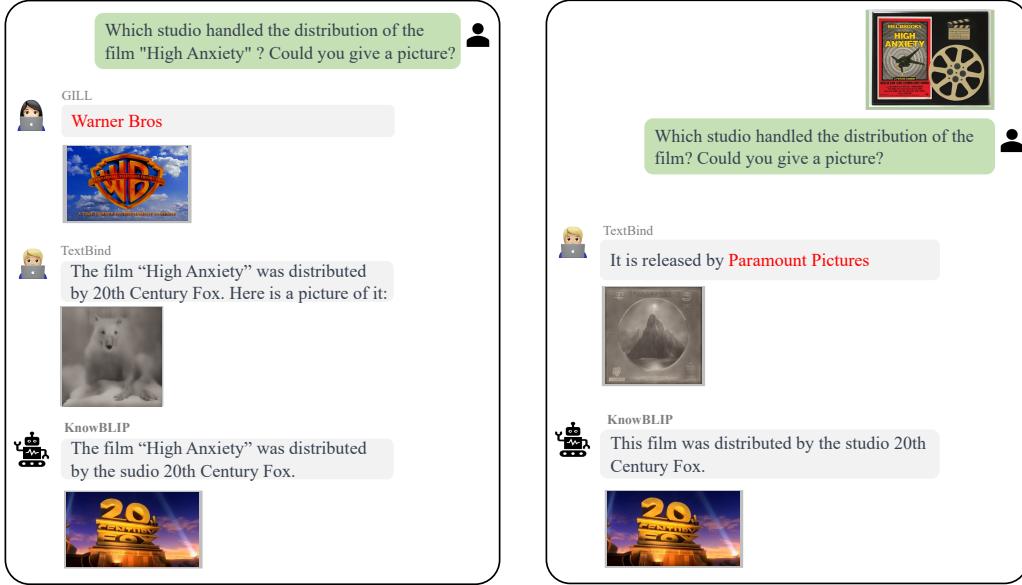
265 **5 Conclusions**

266 In this study, we aim to explore a collection of datasets on world knowledge in order to enhance the  
 267 contextual understanding of LVLMs. Our proposed approach, KnowBLIP, facilitates a comprehensive,  
 268 world-based “any2any” paradigm that bridges the gap between LVLMs and real-world applications.  
 269 Through zero-shot testing and comprehensive comparative analysis, we have demonstrated the  
 270 effectiveness of our method.

271 Moving forward, we plan to expand the scope of this research by increasing the size of the dataset  
 272 and enhancing the parameter size. We believe that augmenting the volume of data and parameters  
 273 will enable the model to better align with real-world contexts.

274 **References**

- 275 [1] Dosovitskiy, A., L. Beyer, A. Kolesnikov, et al. An image is worth 16x16 words: Transformers  
 276 for image recognition at scale, 2021.
- 277 [2] Liu, H., C. Li, Q. Wu, et al. Visual instruction tuning, 2023.



(a) Text2Image task.

(b) Image2Image task.

Figure 8: The visualization of several models in image generation tasks.

- 278 [3] Zhu, D., J. Chen, X. Shen, et al. Minigpt-4: Enhancing vision-language understanding with  
279 advanced large language models. *arXiv preprint arXiv:2304.10592*, 2023.
- 280 [4] Chen, K., Z. Zhang, W. Zeng, et al. Shikra: Unleashing multimodal llm's referential dialogue  
281 magic. *arXiv preprint arXiv:2306.15195*, 2023.
- 282 [5] Lai, X., Z. Tian, Y. Chen, et al. Lisa: Reasoning segmentation via large language model. *arXiv  
283 preprint arXiv:2308.00692*, 2023.
- 284 [6] Zheng, K., X. He, X. E. Wang. Minigpt-5: Interleaved vision-and-language generation via  
285 generative vokens, 2023.
- 286 [7] Dong, R., C. Han, Y. Peng, et al. DreamLLM: Synergistic multimodal comprehension and  
287 creation. In *The Twelfth International Conference on Learning Representations*. 2024.
- 288 [8] Li, H., S. Li, D. Cai, et al. Textbind: Multi-turn interleaved multimodal instruction-following in  
289 the wild, 2024.
- 290 [9] Wu, S., H. Fei, L. Qu, et al. Next-gpt: Any-to-any multimodal llm, 2023.
- 291 [10] Koh, J. Y., D. Fried, R. Salakhutdinov. Generating images with multimodal language models.  
292 *NeurIPS*, 2023.
- 293 [11] Peng, B., C. Li, P. He, et al. Instruction tuning with gpt-4. *arXiv preprint arXiv:2304.03277*,  
294 2023.
- 295 [12] Chen, X., H. Fang, T.-Y. Lin, et al. Microsoft coco captions: Data collection and evaluation  
296 server. *arXiv preprint arXiv:1504.00325*, 2015.
- 297 [13] Changpinyo, S., P. Sharma, N. Ding, et al. Conceptual 12m: Pushing web-scale image-text  
298 pre-training to recognize long-tail visual concepts, 2021.
- 299 [14] Schuhmann, C., R. Vencu, R. Beaumont, et al. Laion-400m: Open dataset of clip-filtered 400  
300 million image-text pairs, 2021.
- 301 [15] Goyal, Y., T. Khot, D. Summers-Stay, et al. Making the V in VQA matter: Elevating the role of  
302 image understanding in Visual Question Answering. In *Conference on Computer Vision and  
303 Pattern Recognition (CVPR)*. 2017.
- 304 [16] Marino, K., M. Rastegari, A. Farhadi, et al. Ok-vqa: A visual question answering benchmark  
305 requiring external knowledge. In *Conference on Computer Vision and Pattern Recognition  
(CVPR)*. 2019.

- 307 [17] Hudson, D. A., C. D. Manning. Gqa: A new dataset for real-world visual reasoning and  
308 compositional question answering. In *Proceedings of the IEEE/CVF conference on computer*  
309 *vision and pattern recognition*, pages 6700–6709. 2019.
- 310 [18] Krishna, R., Y. Zhu, O. Groth, et al. Visual genome: Connecting language and vision using  
311 crowdsourced dense image annotations. 2016.
- 312 [19] Lu, P., S. Mishra, T. Xia, et al. Learn to explain: Multimodal reasoning via thought chains for  
313 science question answering. In *The 36th Conference on Neural Information Processing Systems*  
314 (*NeurIPS*). 2022.
- 315 [20] Singh, A., V. Natarajan, M. Shah, et al. Towards vqa models that can read. In *Proceedings*  
316 *of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8317–8326.  
317 2019.
- 318 [21] Feng, J., Q. Sun, C. Xu, et al. MMDialog: A large-scale multi-turn dialogue dataset towards  
319 multi-modal open-domain conversation. In *Proceedings of the 61st Annual Meeting of the Asso-*  
320 *ciation for Computational Linguistics (Volume 1: Long Papers)*, pages 7348–7363. Association  
321 for Computational Linguistics, Toronto, Canada, 2023.
- 322 [22] Zeng, Y., D. Cao, X. Wei, et al. Multi-modal relational graph for cross-modal video moment  
323 retrieval. In *Proceedings of the CVPR*, pages 2215–2224. IEEE, 2021.
- 324 [23] Touvron, H., T. Lavig, G. Izacard, et al. Llama: Open and efficient foundation language models,  
325 2023.
- 326 [24] Chiang, W.-L., Z. Li, Z. Lin, et al. Vicuna: An open-source chatbot impressing gpt-4 with 90%\*  
327 chatgpt quality, 2023.
- 328 [25] Bai, J., S. Bai, S. Yang, et al. Qwen-vl: A versatile vision-language model for understanding,  
329 localization, text reading, and beyond. *arXiv preprint arXiv:2308.12966*, 2023.
- 330 [26] Du, Z., Y. Qian, X. Liu, et al. Glm: General language model pretraining with autoregressive  
331 blank infilling. In *Proceedings of the 60th Annual Meeting of the Association for Computational*  
332 *Linguistics (Volume 1: Long Papers)*, pages 320–335. 2022.
- 333 [27] Ouyang, L., J. Wu, X. Jiang, et al. Training language models to follow instructions with human  
334 feedback, 2022.
- 335 [28] Rombach, R., A. Blattmann, D. Lorenz, et al. High-resolution image synthesis with latent  
336 diffusion models, 2022.
- 337 [29] Bordes, A., N. Usunier, A. Garcia-Durán, et al. Translating embeddings for modeling multi-  
338 relational data. In *Advances in Neural Information Processing Systems (NIPS 26)*. 2013.
- 339 [30] Radford, A., J. W. Kim, C. Hallacy, et al. Learning transferable visual models from natural  
340 language supervision, 2021.
- 341 [31] OpenAI, J. Achiam, S. Adler, et al. Gpt-4 technical report, 2024.
- 342 [32] Sanh, V., A. Webson, C. Raffel, et al. Multitask prompted training enables zero-shot task  
343 generalization, 2022.
- 344 [33] Li, J., D. Li, S. Savarese, et al. Blip-2: Bootstrapping language-image pre-training with frozen  
345 image encoders and large language models, 2023.
- 346 [34] Wei, J., M. Bosma, V. Y. Zhao, et al. Finetuned language models are zero-shot learners, 2022.
- 347 [35] He, K., H. Fan, Y. Wu, et al. Momentum contrast for unsupervised visual representation learning,  
348 2020.
- 349 [36] Rasley, J., S. Rajbhandari, O. Ruwase, et al. Deepspeed: System optimizations enable training  
350 deep learning models with over 100 billion parameters. In *KDD '20: The 26th ACM SIGKDD*  
351 *Conference on Knowledge Discovery and Data Mining*. 2020.
- 352 [37] Chen, J., D. Zhu, X. Shen, et al. Minigpt-v2: large language model as a unified interface for  
353 vision-language multi-task learning. *arXiv preprint arXiv:2310.09478*, 2023.
- 354 [38] Liu, H., C. Li, Y. Li, et al. Improved baselines with visual instruction tuning, 2024.
- 355 [39] Wang, W., Q. Lv, W. Yu, et al. Cogvlm: Visual expert for pretrained language models, 2024.
- 356 [40] Team, G., M. Reid, N. Savinov, et al. Gemini 1.5: Unlocking multimodal understanding across  
357 millions of tokens of context, 2024.

- 358 [41] Wang, Y., Y. Kordi, S. Mishra, et al. Self-instruct: Aligning language models with self-generated  
359 instructions, 2023.
- 360 [42] Sun, Q., Y. Wang, C. Xu, et al. Multimodal dialogue response generation, 2022.
- 361 [43] Oñoro-Rubio, D., M. Niepert, A. García-Durán, et al. Answering visual-relational queries in  
362 web-extracted knowledge graphs. In *AKBC*. 2019.

363 **A Details of IWK-300K Collection**

364 In this section, we will provide a comprehensive overview of dataset collection.

365 **A.1 Prompt for Data Construction**

366 To construct a high-quality dataset and reduce labor costs, we divide the dataset construction pipeline  
367 into three subtasks and create detailed prompt templates for semi-automatic data production. The  
368 three parts include 1) transforming triplets into natural sentences, 2) writing diverse questions, and 3)  
369 constructing multi-turn conversations.

370 **A.1.1 Step 1 - Transforming Triplets into Natural Sentences**

371 We begin by extracting the original data from a multi-modal knowledge graph in the format of  
372 “`<head_id, relation_chain, tail_id>`” (e.g. “`</m/027rn, /location/country/form_of_government,`  
373 `/m/06cx9>`”). In this format, “`/m/027rn`” and `/m/06cx9`” represent the IDs of the head entity and  
374 tail entity respectively. The “`/location/country/form_of_government`” signifies the relation chain  
375 between the head entity and the tail entity. By post-processing the correspondence relationship of  
376 the ID-entity, we can obtain plaintext in the format of “`<Dominican Republic, /location/country/-`  
377 `form_of_government, Republic>`”. Due to creating directly from triples to QA pairs, we firstly  
378 transforming triplets into natural sentences. Specifically, we randomly select 10 distinct triples and  
379 translate into natural sentences by human experts. Then these triples are treated as few shot templates  
380 and the following prompt is given to GPT-4 to produce natural sentences[41].

**Input:**

**Task definition:** there is triple format information from a knowledge graph. You are required  
to transform it in natural sentences.

**Templates:**

Triple 1:< Boston College > < /location/location/containedby > < United States of America >  
Answer 1:

1. Boston College is located within the United States of America.
2. Boston College is contained by the United States of America.

.....

Triple n: <Keio University> </location/location/containedby><Japan>

Answer n:

**Output:**

1. Keio University is located within the United States of America.
2. Keio University contained by the United States of America.

.....

382 For each triples that answers with fewer than 5 sentences, we will call GPT4’s API again with the  
383 following prompt to further enrich it. In order to improve the quality and richness of the answers,  
384 we will include a more detailed prompt in this repeated API call. In this way, GPT4 can further  
385 explore and integrate relevant information based on new prompts, generating richer and more detailed  
386 answers. This answer not only includes more sentences and details, but also better demonstrates the  
387 knowledge domain and contextual environment involved in triplets.

**Input:**

**Task definition:** Could you give me some synonymous description for this sentence:

**Templates:**

Sentence 1:Boston College is located within the United States of America. Answer 1:

1. Boston College is contained by the United States of America.
2. Boston College is situated in the United States of America.

.....

Sentence n: Ithaca College is located in New York.

Answer n:

**Output:**

1. Ithaca College is situated in New York.
2. The location of Ithaca College is New York.

.....

389 **A.1.2 Step 2 - Writing Diverse Questions**

390 In this step, we have obtained world knowledge presented in natural sentence form from the previous  
391 step. To ensure the diversity of the final multi round dialogue[42], we generated question and  
392 answer pairs (QA pairs) for each natural sentence using prompts in the main text. The final generated  
393 question and answer pairs will have rich content and diverse forms, covering various possible dialogue  
394 scenarios and topics. These question and answer pairs can be used as materials for constructing  
395 multiple rounds of dialogue and as datasets for training dialogue models, improving the generation  
396 ability and interaction experience of dialogue models. These QA-pairs are in the following format.

```
{ "Q": "What is the form of government in the Dominican Republic?",  
  "A": "The Dominican Republic has a form of government that is a republic." }  
{ "Q": "Where is the location of Ithaca College?",  
  "A": "Ithaca College is located in New York." }  
...
```

398 **A.1.3 Step 3 - Constructing Multi-Turn Conversations**

399 To construct engaging and informative multi-turn conversations that effectively promote learning for  
400 each entity, we embark on a meticulous process. Initially, we gather all the question-answer (QA)  
401 pairs related to each entity, ensuring a comprehensive understanding of the subject matter. The heart  
402 of our approach lies in the careful selection and composition of questions. We randomly select a  
403 batch of questions from the QA pairs set of a particular entity. This ensures that the questions cover a  
404 wide range of topics related to the entity, making the conversation diverse and engaging. To further  
405 enrich the conversational experience, we introduce an innovative element: the random replacement  
406 of entity concepts with corresponding images. With a 50% probability, we substitute the textual  
407 description of an entity with a visually arresting image. This not only adds a visual dimension to the  
408 conversation but also aids in the comprehension of abstract or complex concepts. Moreover, this  
409 approach ensures that each entity is represented in a unique and engaging manner, fostering adequate  
410 learning and understanding. By combining textual and visual elements, we create a immersive and  
411 interactive learning experience that is tailored to the needs and preferences of the learners.

412 **A.2 Image Filter for Constructed Data**

413 For each entity in the meticulously constructed dataset IWK300K, we have the capability to retrieve  
414 the original images associated with it. On average, each entity boasts approximately 61.23 images,  
415 providing a rich visual representation of the entity's diverse attributes and characteristics. However,  
416 upon closer inspection, we observe that a significant portion of these images contain excessive  
417 irrelevant information and are highly disconnected from the semantic essence of the entity. To ensure  
418 that Large Vision-Language Models (LVLMs) can effectively learn from the images of these entities,  
419 it becomes imperative to meticulously filter out these “dirty” images. This filtering process not  
420 only enhances the quality of the dataset but also ensures that the LVLMs can focus on extracting  
421 meaningful and relevant features from the images. As Illustrated in Figure 3, the process of selecting  
422 an appropriate image from the Multimedia Knowledge Graph (MMKG) involves three distinct steps:  
423 1) Filtering Gate, 2) Semantic Gate and 3) Refining Gate. The following is the detailed introduction  
424 for these three steps.

425 **A.2.1 Filtering Gate**

426 In this gate, we filter images from an image perspective by computing similarity scores between the  
427 images of an entity. We utilize two algorithms, dHash and CLIP, to extract features and calculate  
428 similarity scores. Specifically, for an image set of one entity,  $IMG = (img_1, \dots, img_n)$ , where  
429  $img_{(.)}$  represents a candidate image of a entity. Then, for each candidate image, we compute its  
430 similarity to the other images and sum these similarities. The dHash algorithm can be computed  
431 using the following formula.

$$S_{dHash}^{img_i} = \sum_{t=1, t \neq i}^n matmul(dHash(img_i), dHash(img_t)^T). \quad (6)$$

432 Similarly, the CLIP algorithm could be computed with the following formula.

$$S_{CLIP}^{img_i} = \sum_{t=1, t \neq i}^n matmul(CLIP(img_i), CLIP(img_t)^T). \quad (7)$$

433 Through the calculations of the above two algorithms, we could find the image set exists the situation  
 434 where two images are identical. To ensure that each selected image only appears once and avoid the  
 435 high ranking of weighted similarity caused by duplicate images. We retain only one of every two  
 436 images with a similarity of 1 in two algorithms, and filter out the rest. Then we can obtain the final  
 437 score by normalizing the mean of the two similarities with the following formula, where the  $\lambda$  is set  
 438 to 0.2.

$$S_{final}^{img_i} = \lambda S_{dHash}^{img_i} + (1 - \lambda) S_{CLIP}^{img_i}. \quad (8)$$

439 Finally, our filter gate utilize the  $S_{final}^{img_i}$  as the score for sorting and take the top 10 image for the next  
 440 step filtering.

#### 441 A.2.2 Semantic Gate

442 In semantic gate, we filter images from its corresponding entity texts perspective. For this purpose,  
 443 we use CLIP to obtain text and image features, and we compute the similarity between text and image  
 444 features to represent the semantic similarity. CLIP is a multi-modal pre-trained model that uses a  
 445 dual-stream encoding framework and has trained on large-scale data comprising four hundred million  
 446 image and text pairs collected from the internet. It is highly effective on a wide range of downstream  
 447 tasks and possesses a powerful ability to project text and image into the same semantic space.

448 Specifically, given a entity text  $t$  and its first step filtered images set  $IMG = (img_1, \dots, img_{10})$ , we  
 449 would encode entity text  $t$  by CLIP text encoder while encode image  $img_{(.)}$  by CLIP visual encoder.  
 450 The similarity between entity text  $t$  and its candidate image  $img_i$  could be computed in the following  
 451 formula.

$$S_{semantic}^{t-img_i} = matmul(CLIP_{visual}(img_i), CLIP_{text}(t)^T). \quad (9)$$

452 Then we could obtain the order of candidate image set based on the calculated semantic similarity  
 453 score and take top 5 image for the last step selecting.

#### 454 A.2.3 Refining Gate

455 On the account of the limited filter ability of existing model or algorithm, we utilize human  
 456 experts to select a image from the previous step. Given a entity text  $t$ , its filtered image set  
 457  $IMG = (img_1, \dots, img_5)$ , the corresponding set of constructed natural sentence (Appx 1.1) and  
 458 corresponding Wikipedia, human experts would firstly select image based on their understanding  
 459 of entity text and its constructed single turn conversations. (e.g. entity “20th century fox” and its  
 460 constructed natural sentence “The flim xxxx is released by the studio 20th century fox”). There are 5  
 461 human experts in total. Each human expert will select an appropriate image and finally count the  
 462 number of times each image has been selected. When the number of times is greater than or equal to  
 463 2, this image will become the final selection result. For most entities in the dataset, human experts  
 464 will have a consistent final result. The rest of entities will be combined to Wikipedia, origin images  
 465 set and Google search engine to find a appropriate image.

## 466 B Dataset Diversity

467 To assess the diversity of sentences within the IWK-300k datasets, we conducted a topic frequency  
 468 analysis of the world knowledge encompassed in the IWK-300k corpus. As depicted in Figure 9,  
 469 we initially quantified the overall frequency of each topic across the entities present. Notably, we  
 470 identified a total of 63 distinct topic types. Subsequently, through a rigorous process of human  
 471 classification and summarization, these topics were grouped into 12 distinct areas, collectively  
 472 covering a broad and diverse range of real-world domains.

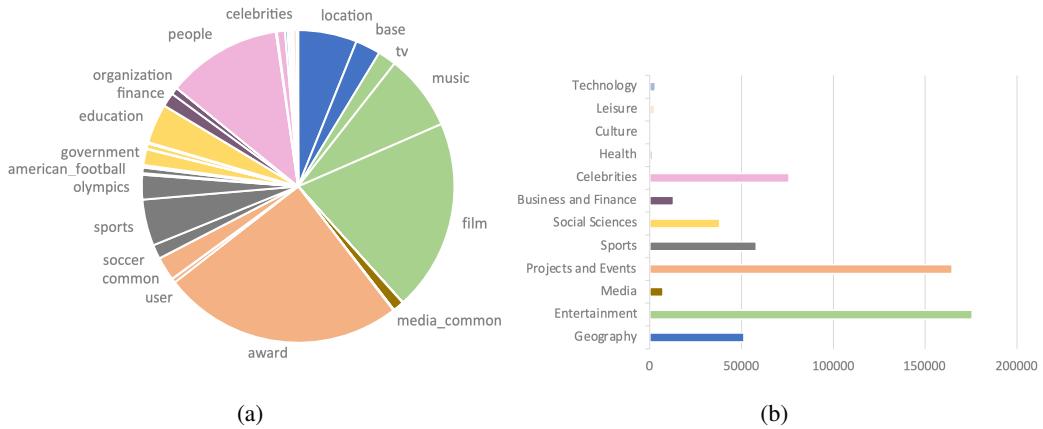


Figure 9: The distribution of topics in our dataset.

#### Example1 JSON of our IWK-300k dataset

```
{
  "id": "m0291hr0",
  "conversations": [
    {
      "from": "human",
      "value": "Who_is_the_director_of_this_film?\n<image>",
      "image_list": ["m.0291hr.jpg"]
    },
    {
      "from": "model",
      "value": "The_film_\"High Anxiety\"_was_produced_by_Mel_Brooks .",
      "image_list": []
    },
    {
      "from": "human",
      "value": "What_studio_handled_the_distribution_of_the_film_\"High Anxiety\"? Could_you_give_a_picture_of_it?",
      "image_list": []
    },
    {
      "from": "model",
      "value": "The_film_\"High Anxiety\"_was_distributed_by_the_studio_20th_Century_Fox . Here_is_a_picture_of_it: \n<image>",
      "image_list": ["m.016tt2.jpg"]
    },
    {
      "from": "human",
      "value": "Do_you_have_any_similar_moments?"
    },
    {
      "from": "model",
      "value": "20th_Century_Fox_had_distributed_the_film_\"The Happening .\""
      "image_list": []
    }
  ]
}
```

## Example2 JSON of our IWK-300k dataset

---

```
{  
    "id": "m01_x6v0",  
    "conversations": [  
        {  
            "from": "human",  
            "value": "Which\u00a0university\u00a0did\u00a0<image>\u00a0graduate\u00a0from?"  
            "image_list": ["m.01_x6v.jpg"]  
        },  
        {  
            "from": "model",  
            "value": "Trey\u00a0Parker\u00a0was\u00a0educated\u00a0at\u00a0Berklee\u00a0College\u00a0of\u00a0Music."  
            "image_list": []  
        },  
        {  
            "from": "human",  
            "value": "Where\u00a0is\u00a0the\u00a0Berklee\u00a0College\u00a0of\u00a0Music\u00a0situated?\u00a0Could\u00a0you\u00a0give\u00a0a\u00a0picture\u00a0of\u00a0it?\u00a0"  
            "image_list": []  
        },  
        {  
            "from": "model",  
            "value": "The\u00a0Berklee\u00a0College\u00a0of\u00a0Music\u00a0is\u00a0located\u00a0in\u00a0Boston.\u00a0Here\u00a0is\u00a0a\u00a0picture\u00a0of\u00a0it:\u00a0\n<image>"  
            "image_list": ["m.02g839.jpg"]  
        }  
    ]  
}
```

---

## 473 C Case study

474 The four formats of the IWK-300K dataset endow KnowBLIP with a diverse array of functionalities,  
475 enabling it to handle a wide range of tasks and scenarios related to real world knowledge. To illustrate  
476 these capabilities, we present several concrete user cases in Figure 11-15, each showcasing a unique  
477 aspect of KnowBLIP’s versatility.

### 478 C.1 Text Perception and Generation

479 One core innovation of KnowBLIP lies in its remarkable ability to perceive and generate real-world  
480 knowledge based on the contextual nuances of user conversations. This aspect is particularly valuable  
481 in open-world scenarios, where individuals often have intricate and multifaceted ideas about the  
482 world, transcending mere general knowledge. Consider a conversation between a user and KnowBLIP.  
483 The user might express an interest in learning about a particular topic, say, the film “High Anxiety”.  
484 Our KnowBLIP with its unique capability, not only understands the user’s intention but also responds  
485 with accurate real-world knowledge related to the entity “High Anxiety.” It might offer details about  
486 the film’s director, plot and so on. Moreover, KnowBLIP excels in handling interleaved text and  
487 images with different entities in multi-turn conversations. Imagine a conversation where the user  
488 first asks about “High Anxiety” and then transitions to another entity, such as “20th Century Fox”,  
489 the studio that produced it. KnowBLIP seamlessly shifts its focus, maintaining coherence and depth  
490 throughout the conversation.

### 491 C.2 Visual Perception and Generation

492 The other core creation of KnowBLIP lies in its ability to not only comprehend but also retrieve  
493 images related to world knowledge entities based on the contextual conversation flow provided by  
494 users. For instance, consider entities like “Richard Rodgers”, “New York, I Love You”, “University  
495 of Oxford”, and “Germany”. These entities are rich in contextual information and often form the  
496 backbone of conversations about various topics. In an open-world scenario, this feature of KnowBLIP

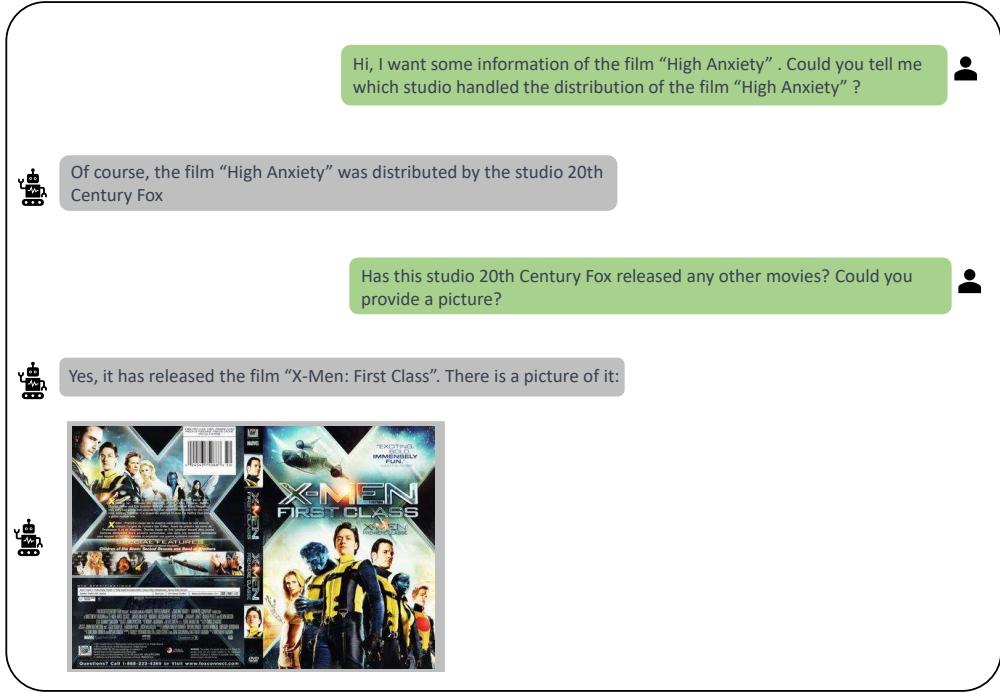


Figure 10: The visualization of KnowBLIP’s text understanding, i.e. Text2Text and Text2Image tasks.

497 becomes particularly handy. Many times, users may simply input an image containing some form  
 498 of world knowledge and inquire about its associated information. Such inquiries could range from  
 499 identifying landmarks in a photograph to understanding the historical significance of a painting.  
 500 Furthermore, there are instances where users actively seek images related to specific world knowledge  
 501 entities. For example, a user might be interested in exploring images related to the movie “X-Men:  
 502 First Class” to gain a deeper understanding of its characters, settings, or themes. In such cases,  
 503 KnowBLIP’s ability to retrieve relevant images becomes invaluable. However, it’s worth noting that  
 504 not all image generation schemes are equipped to handle world knowledge effectively. As illustrated  
 505 in Figure 8 (a), the use of a stable diffusion model for image generation can sometimes fall short in  
 506 this regard. While the model may produce visually appealing images, it may lack the contextual depth  
 507 and world knowledge that KnowBLIP is capable of incorporating. To address this gap, KnowBLIP  
 508 incorporates a robust knowledge base and advanced retrieval mechanisms. This allows it to not  
 509 only generate images but also ensure that these images are informed by relevant world knowledge.  
 510 Whether it’s retrieving images of famous landmarks or exploring the visual representations of fictional  
 511 worlds, KnowBLIP provides users with a richer and more informative visual experience.

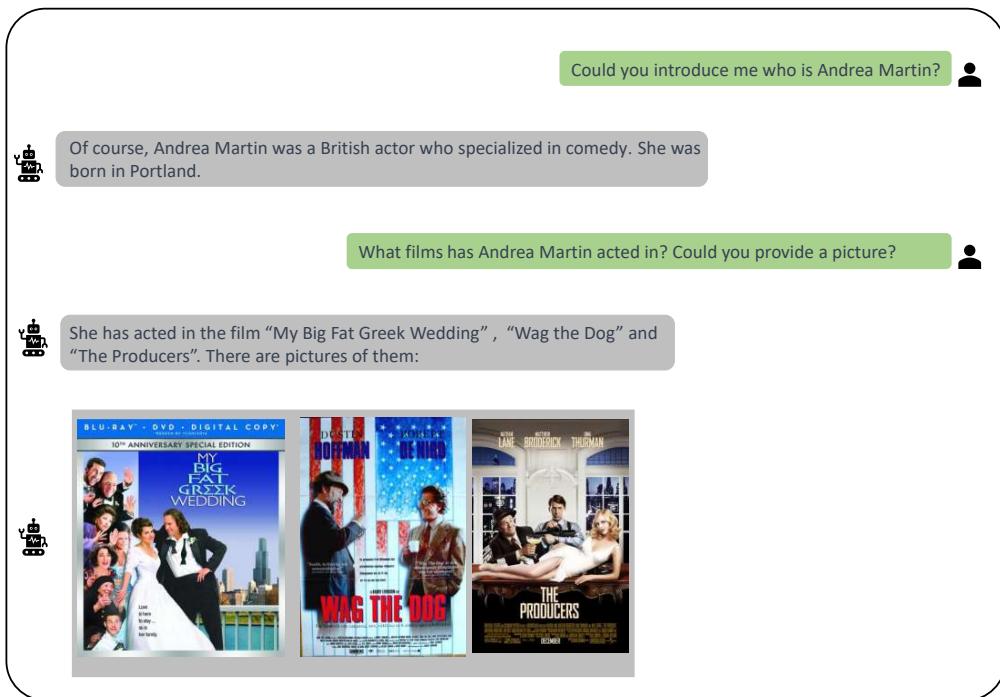


Figure 11: The visualization of KnowBLIP's text understanding, i.e. Text2Text and Text2Image tasks.

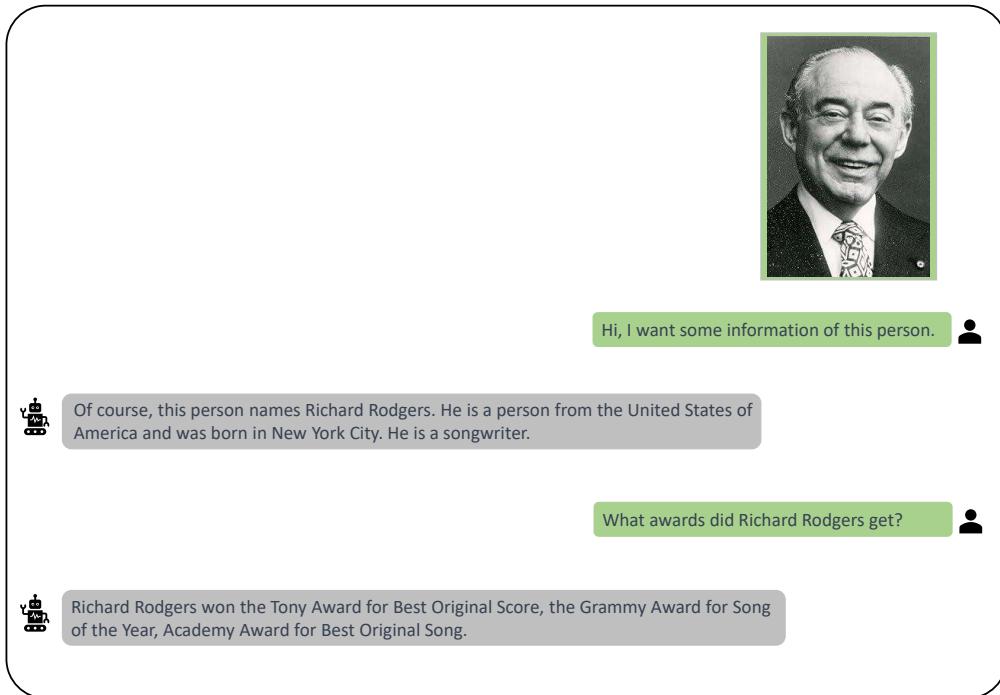


Figure 12: The visualization of KnowBLIP's image understanding, i.e. Image2Text and Image2Image tasks.

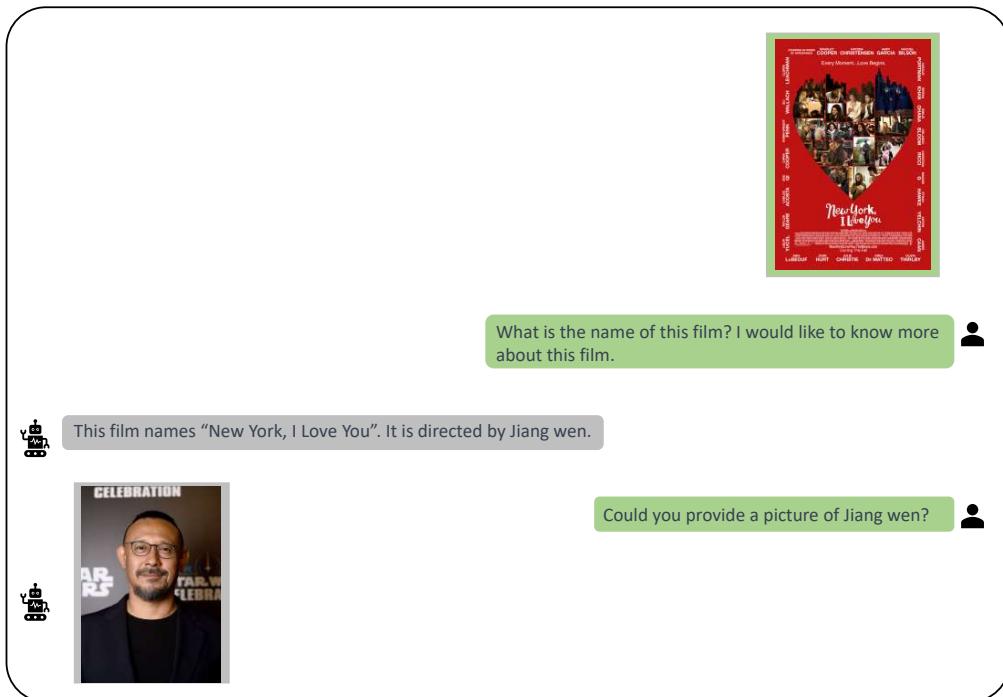


Figure 13: The visualization of KnowBLIP’s image understanding, i.e. Image2Text and Image2Image tasks.

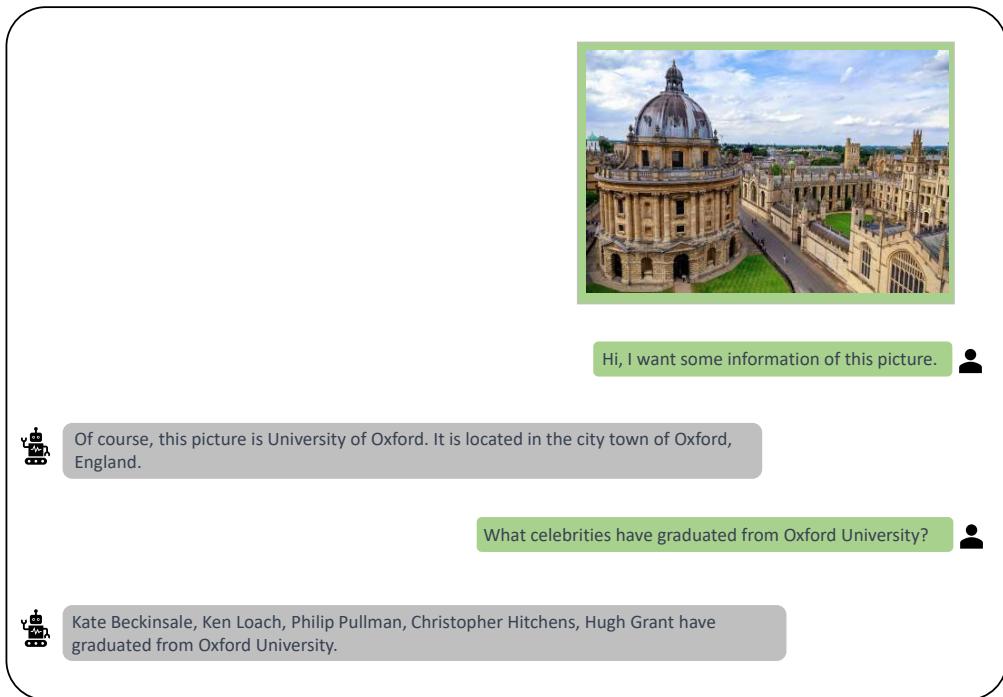


Figure 14: The visualization of KnowBLIP’s image understanding, i.e. Image2Text and Image2Image tasks.

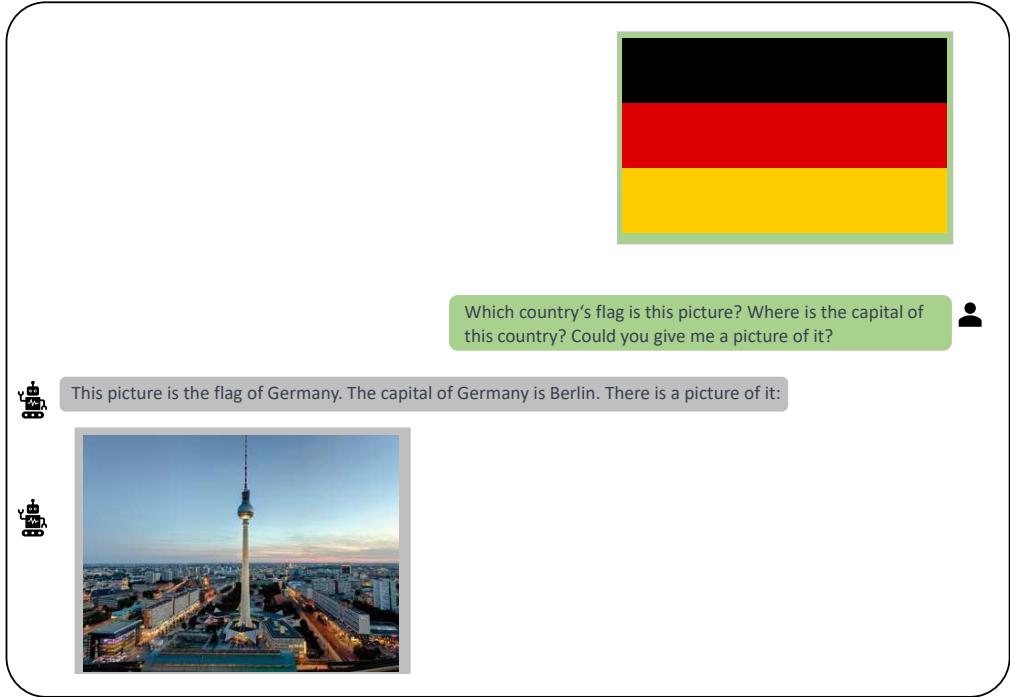


Figure 15: The visualization of KnowBLIP’s image understanding, i.e. Image2Text and Image2Image tasks.

## 512 D Datasheet of IWK-300k

513 In this section, we provide a datasheet about our dataset.

### 514 D.1 Motivation For Datasheet Creation

515 **Why was the datasheet created? (e.g., was there a specific task in mind? was there a specific  
516 gap that needed to be filled?)**

517 This dataset was developed to facilitate research on Multi-modal Large Language models for the world  
518 knowledge based on their image-text understanding ability. While there are existing datasets for Multi-  
519 modal Large Language models to perceive image and text, none of them include a comprehensive set  
520 of image-text pair for learning world knowledge. By creating a multi-turn and interveled annotated  
521 dataset, we aim to fill this gap and provide researchers with a resource that can be used to develop  
522 more complete ablities multi-modal models for world knowledge.

523 **Has the dataset been used already? If so, where are the results so others can compare (e.g.,  
524 links to published papers)?**

525 Our dataset is proposed for the first time. This dataset was built upon a subset of existing multi-modal  
526 knowledge graph (FB15K[29]), which is consisted by structured triples with the relation in the real  
527 world. Our new dataset, IWK-300K, serves as an extension to the this multi-modal knowledge graph  
528 dataset, providing more comprehensive world knowledge annotations in conversation format for  
529 multi-modal models, which can be used to improve the complete abilities of multi-modal models.

530 **What (other) tasks could the dataset be used for?**

531 This dataset could be used for diverse tasks under world knowledge, including but not limited to tasks  
532 such as image caption, VQA, image generation, cross-modal retrieval.

533 **Who funded the creation dataset?**

534 None. We built this dataset at our own expense to promote the development of academia and industry.

535 **Any other comment?**

536 None.

## 537 **D.2 Datasheet Composition**

538 **What are the instances?(that is, examples; e.g., documents, images, people, countries) Are there  
539 multiple types of instances? (e.g., movies, users, ratings; people, interactions between them;  
540 nodes, edges)**

541 The instances in this dataset are images and triples from the FB15K dataset[43]. The original images  
542 and triples contain information about the the image of nodes and relation between nodes. Our  
543 new dataset is densely translateed the original images and triples into multi-turn, multi-modal and  
544 interleaved conversations established through GPT-4 and human experts, which contains with 4  
545 formats (text-text, image-text, image-image, text-image)

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

547 IWK-300k is a dataset of 300, 199 single rounds of conversation, averaging 5.95 rounds per complete  
548 multi-turn conversation. Each single turn conversation includes two entity and their relation. The  
549 average number of images per entity is 61.23.

550 **What data does each instance consist of ? “Raw” data (e.g., unprocessed text or images)?  
551 Features/attributes? Is there a label/target associated with instances? If the instances related to  
552 people, are subpopulations identified (e.g., by age, gender, etc.) and what is their distribution?**

553 The IWK-500K dataset consists of multiple types of instance. The types consist of: (a) text to text  
554 conversation, (b) text to image conversation, (c) image to text conversation, (d) image to image  
555 conversation.

556 **Is there a label or target associated with each instance? If so, please provide a description.**

557 Yes, a single turn conversation often constructed by question of head entity in origin triple, while the  
558 label is the response of conversation constructed by tail entity in origin triple.

559 **Is any information missing from individual instances? If so, please provide a description,  
560 explaining why this information is missing (e.g., because it was unavailable). This does not  
561 include intentionally removed information, but might include, e.g., redacted text.**

562 The dataset does not include all complete information on a entity, as this information is also not  
563 available from the original dataset.

564 **Are relationships between individual instances made explicit (e.g., users’ movie ratings, social  
565 network links)? If so, please describe how these relationships are made explicit.**

566 Not applicable.

567 **Does the dataset contain all possible instances or is it a sample (not necessarily random) of  
568 instances from a larger set? If the dataset is a sample, then what is the larger set? Is the sample  
569 representative of the larger set (e.g., geographic coverage)? If so, please describe how this  
570 representativeness was validated/verified. If it is not representative of the larger set, please  
571 describe why not (e.g., to cover a more diverse range of instances, because instances were  
572 withheld or unavailable).**

573 No. IWK-500K only contains entity and its relation between entities on public wikipedia.

574 **Are there recommended data splits (e.g., training, development/validation, testing)? If so, please  
575 provide a description of these splits, explaining the rationale behind them.**

576 Our dataset is split into 3 non-overlapping subsets, where 0.6, 0.2 and 0.2 are used for training,  
577 validation and testing. Specifically, for one triplet that has four formats in different formats of  
578 question-answer pair, we would split its different formats into training, validation and testing sets,  
579 which ensures the world knowledge of this triplet having trained.

580 **Are there any errors, sources of noise, or redundancies in the dataset? If so, please provide a  
581 description.**

582 Noise and errors are inevitable in datasets. The most likely source of errors is incorrect labels or  
583 relation in the origin dataset or in the semi-automatic data production from GPT-4 caused by its  
584 hallucination. However, the IWK-500K dataset collection process has multiple quality assurance  
585 steps by human experts that aims to substantially reduce the prevalence of noise and errors.

586 **Is the dataset self-contained, or does it link to or otherwise rely on external resources (e.g.,  
587 websites, tweets, other datasets)? If it links to or relies on external resources, a) are there  
588 guarantees that they will exist, and remain constant, over time; b) are there official archival  
589 versions of the complete dataset (i.e., including the external resources as they existed at the time  
590 the dataset was created); c) are there any restrictions (e.g., licenses, fees) associated with any  
591 of the external resources that might apply to a future user? Please provide descriptions of all  
592 external resources and any restrictions associated with them, as well as links or other access  
593 points, as appropriate.**

594 Our dataset is developed from the existing FB15K dataset. The FB15k dataset contains knowl-  
595 edge base relation triples and textual mentions of Freebase entity pairs. It has a total of 592,213  
596 triplets with 14,951 entities and 1,345 relationships. The FB15K dataset can be accessed here:  
597 [https://github.com/mniepert/mmkb/blob/master/FB15K/FB15K\\_EntityTriples.txt](https://github.com/mniepert/mmkb/blob/master/FB15K/FB15K_EntityTriples.txt).

598 **Any other comments?**

599 None.

### 600 D.3 Collection Process

601 **What mechanisms or procedures were used to collect the data (e.g., hardware apparatus or  
602 sensor, manual human curation, software program, software API)? How were these mechanisms  
603 or procedures validated?**

604 We use GPT-4 API to semi-automatic data production. Prompts had detailed in main content.

605 **How was the data associated with each instance acquired? Was the data directly observable  
606 (e.g., raw text, movie ratings), reported by subjects (e.g., survey responses), or indirectly  
607 inferred/derived from other data (e.g., part-of-speech tags, model-based guesses for age or  
608 language)? If data was reported by subjects or indirectly inferred/derived from other data, was  
609 the data validated/verified? If so, please describe how.**

610 A full description appears in the associated paper and its appendix. The three steps include 1)  
611 transforming triplets into natural sentences, 2) writing diverse questions, and 3) constructing multi-  
612 turn conversations. 4) Human experts review. However, briefly:

- 613 • Transforming Triplets into Natural Sentences. Some few shot templates, the prompt and required  
614 translated entity is given to GPT-4 to produce data.
- 615 • Writing Diverse Questions. generated question and answer pairs (QA pairs) for each natural  
616 sentence using prompts given to GPT-4.
- 617 • Constructing Multi-Turn Conversations. Randomly select a batch of questions from the QA pairs  
618 set of a particular entity
- 619 • Human experts review.

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

622 All images from the FB15K had filtered by a sampling strategy, which including three distinct steps:  
623 1) Filtering Gate, 2) Semantic Gate and 3) Refining Gate. In the main content is involved the detailed  
624 introduction for these three steps.

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

627 The Artificial Intelligent students from South China Technology of University were involved in  
628 reviewed of the data. These students were compensated based on the number of single conversation  
629 they reviewed, with payment free.

630 **Over what timeframe was the data collected? Does this timeframe match the creation timeframe  
631 of the data associated with the instances (e.g., recent crawl of old news articles)? If not, please  
632 describe the timeframe in which the data associated with the instances was created.**

633 According to the original paper [29], the collection of FB15K images took place in 2015. The  
634 attribute annotations were conducted in December 2023 - Third 2024.

635 **D.4 Data Preprocessing**

636 **Was any preprocessing/cleaning/labeling of the data done (e.g., discretization or bucketing,  
637 tokenization, part-of-speech tagging, SIFT feature extraction, removal of instances, processing  
638 of missing values)? If so, please provide a description. If not, you may skip the remainder of  
639 the questions in this section.**

640 No preprocessing. The labeling of the dataset was done manually by human annotators.

641 **Was the “raw” data saved in addition to the preprocessed/cleaned/labeled data (e.g., to support  
642 unanticipated future uses)? If so, please provide a link or other access point to the “raw” data.**

643 Not applicable.

644 **Is the software used to preprocess/clean/label the instances available? If so, please provide a  
645 link or other access point.**

646 Not applicable.

647 **Does this dataset collection/processing procedure achieve the motivation for creating the dataset  
648 stated in the first section of this datasheet? If not, what are the limitations?**

649 Not applicable.

650 **Any other comments**

651 None.

652 **D.5 Dataset Distribution**

653 **How will the dataset be distributed? (e.g., tarball on website, API, GitHub; does the data have  
654 a DOI and is it archived redundantly?)**

655 The FB15K dataset can be accessed through the following link: <sup>3</sup>. Additionally, the newly annotated  
656 dataset, IWK-500K, is uploaded here: <sup>4</sup>.

657 **When will the dataset be released/first distributed? What license (if any) is it distributed under?**

658 The dataset is already available.

659 **Are there any copyrights on the data?**

660 The FB15K is released under BSD-3-Clause License.

661 **Are there any fees or access/export restrictions?**

662 There are no fees. Our dataset is also intended to be used for research and educational purposes.

663 **Any other comments?** None.

664 **D.6 Dataset Maintenance**

665 **Who is supporting/hosting/maintaining the dataset?**

666 The dataset and the website where the annotations are released will be maintained by the authors of  
667 the manuscript.

668 **Will the dataset be updated? If so, how often and by whom?**

669 We do not have concrete plans as of yet; we will announce any updates on the dataset website.

---

<sup>3</sup><https://github.com/mnietpert/mmkb/tree/master>

<sup>4</sup>[https://drive.google.com/drive/folders/1ZgfApnIhzXg8n3R\\_kwix5C0Klv6UP8](https://drive.google.com/drive/folders/1ZgfApnIhzXg8n3R_kwix5C0Klv6UP8)

670 **How will updates be communicated? (e.g., mailing list, GitHub)**

671 Updates will be communicated through the dataset website: <sup>5</sup>.

672 **If the dataset becomes obsolete how will this be communicated?**

673 Through our github website: <sup>6</sup>.

674 **Is there a repository to link to any/all papers/systems that use this dataset?**

675 There is a repository, maintained by the authors of the manuscript at <sup>7</sup>.

676 **If others want to extend/augment/build on this dataset, is there a mechanism for them to do so?**

677 **If so, is there a process for tracking/assessing the quality of those contributions. What is the  
678 process for communicating/distributing these contributions to users?**

679 The dataset is under the MIT licence so anyone has freedom to do so. Currently, we do not have  
680 mechanisms in place; however, others may contact us to discuss potential use cases if they prefer.

---

<sup>5</sup>[https://drive.google.com/drive/folders/1ZgfApnIhzXg8n3R\\_kwix5C0Klv6UP8](https://drive.google.com/drive/folders/1ZgfApnIhzXg8n3R_kwix5C0Klv6UP8)

<sup>6</sup><https://github.com/sdjhshbswp/IWK500K>

<sup>7</sup>

681 **NeurIPS Paper Checklist**

682 **1. Claims**

683 Question: Do the main claims made in the abstract and introduction accurately reflect the  
684 paper's contributions and scope?

685 Answer: [Yes]

686 Justification: We have summarized the main contributions at the end of the introduction.

687 **2. Limitations**

688 Question: Does the paper discuss the limitations of the work performed by the authors?

689 Answer: [Yes]

690 Justification: We have summarized the limitations in the experiment and appendix.

691 **3. Theory Assumptions and Proofs**

692 Question: For each theoretical result, does the paper provide the full set of assumptions and  
693 a complete (and correct) proof?

694 Answer: [NA]

695 Justification: Our paper does not include theoretical results.

696 **4. Experimental Result Reproducibility**

697 Question: Does the paper fully disclose all the information needed to reproduce the main ex-  
698 perimental results of the paper to the extent that it affects the main claims and/or conclusions  
699 of the paper (regardless of whether the code and data are provided or not)?

700 Answer: [Yes]

701 Justification: We have provided the available URL of all codes and datasets.

702 **5. Open access to data and code**

703 Question: Does the paper provide open access to the data and code, with sufficient instruc-  
704 tions to faithfully reproduce the main experimental results, as described in supplemental  
705 material?

706 Answer: [Yes]

707 Justification: We have provided the available URL of all codes and datasets.

708 **6. Experimental Setting/Details**

709 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-  
710 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the  
711 results?

712 Answer: [Yes]

713 Justification: We have displayed it in the experiment.

714 **7. Experiment Statistical Significance**

715 Question: Does the paper report error bars suitably and correctly defined or other appropriate  
716 information about the statistical significance of the experiments?

717 Answer: [No]

718 Justification: We will add this experiment in the future.

719 **8. Experiments Compute Resources**

720 Question: For each experiment, does the paper provide sufficient information on the com-  
721 puter resources (type of compute workers, memory, time of execution) needed to reproduce  
722 the experiments?

723 Answer: [Yes]

724 Justification: We have explained it in the setting.

725 **9. Code Of Ethics**

726 Question: Does the research conducted in the paper conform, in every respect, with the  
727 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

- 728                  Answer: [Yes]  
729                  Justification: We have confirmed the Ethics Guidelines.
- 730                  **10. Broader Impacts**  
731                  Question: Does the paper discuss both potential positive societal impacts and negative  
732                  societal impacts of the work performed?  
733                  Answer: [Yes]  
734                  Justification: We believe our work will inspire the multi-modal field.
- 735                  **11. Safeguards**  
736                  Question: Does the paper describe safeguards that have been put in place for responsible  
737                  release of data or models that have a high risk for misuse (e.g., pretrained language models,  
738                  image generators, or scraped datasets)?  
739                  Answer: [NA]  
740                  Justification: Our paper poses no such risks.
- 741                  **12. Licenses for existing assets**  
742                  Question: Are the creators or original owners of assets (e.g., code, data, models), used in  
743                  the paper, properly credited and are the license and terms of use explicitly mentioned and  
744                  properly respected?  
745                  Answer: [Yes]  
746                  Justification: We have provided the available URL of datasets.
- 747                  **13. New Assets**  
748                  Question: Are new assets introduced in the paper well documented and is the documentation  
749                  provided alongside the assets?  
750                  Answer: [NA]  
751                  Justification: Our paper does not release new assets.
- 752                  **14. Crowdsourcing and Research with Human Subjects**  
753                  Question: For crowdsourcing experiments and research with human subjects, does the paper  
754                  include the full text of instructions given to participants and screenshots, if applicable, as  
755                  well as details about compensation (if any)?  
756                  Answer: [NA]  
757                  Justification: Our paper does not involve crowdsourcing nor research with human subjects.
- 758                  **15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human  
759                  Subjects**  
760                  Question: Does the paper describe potential risks incurred by study participants, whether  
761                  such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)  
762                  approvals (or an equivalent approval/review based on the requirements of your country or  
763                  institution) were obtained?  
764                  Answer: [NA]  
765                  Justification: Our paper does not involve crowdsourcing nor research with human subjects.