EXPLORING MULTI-GRAINED CONCEPT ANNOTATIONS FOR MULTIMODAL LARGE LANGUAGE MODELS

Anonymous authors

Paper under double-blind review

ABSTRACT

Multimodal Large Language Models (MLLMs) excel in vision–language tasks by pre-training solely on coarse-grained concept annotations (*e.g.*, image captions). We hypothesize that integrating fine-grained concept annotations (*e.g.*, object labels and object regions) will further improve performance, as both data granularities complement each other in terms of breadth and depth in concept representation.

We introduce a new dataset featuring Multimodal Multi-Grained Concept annotations (MMGIC) for MLLMs. In constructing MMGIC, we explore the impact of different data recipes on multimodal comprehension and generation. Our analyses reveal that multi-grained concept annotations integrate and complement each other, under our structured template and autoregressive discrete framework.

We definitively show that multi-grained concepts do facilitate MLLMs to better locate and learn concepts, aligning vision and language at multiple granularities. We further validate our hypothesis by investigating the comparison and collaboration between MMGIC and image–caption data on 12 multimodal comprehension and generation benchmarks, *e.g.*, their appropriate combination achieve 3.95% and 2.34% accuracy improvements on POPE and SEED-Bench. Code, data and models will be made openly available.

025

004

010 011

012

013

014

015

016

017

018

019

021

1 INTRODUCTION

With the rapid development of Large Language Models (LLMs) (Brown et al., 2020; Chowdhery et al., 2023; Touvron et al., 2023; and Visual Foundation Models (VFMs) (Radford et al., 2021; Rombach et al., 2022; Dehghani et al., 2023), researchers have started to explore the potential of unifying them into Multimodal Large Language Models (MLLMs) to perform various Vision–Language (VL) tasks, such as image captioning, visual question answering, and text-to-image generation. MLLMs, such as Emu (Sun et al., 2023c), SEED-LLaMA (Ge et al., 2023b) and LaVIT (Jin et al., 2023), follow a similar framework that integrates the capabilities of LLMs and VFMs under an autoregressive objective of predicting the next visual or textual token, showing impressive performance on VL tasks.

Despite their success, existing MLLMs typically do not make full use of *concepts* in VL learning, 040 relying on coarse-grained concept annotations (e.g., image captions) but ignoring fine-grained concept annotations (e.g., object labels and object regions). This leads to superficial and incomplete 041 understanding of concepts, limiting VL alignment. Specifically, by "concept" we mean an abstraction 042 and generalization of a group of things having common characteristics (Goguen, 1969; Carey, 2000; 043 Blouw et al., 2016). Concepts can be categorized into concrete and abstract concepts by whether they 044 can be sensed by five human senses (Shevade et al., 2005; Connell et al., 2018). Concrete concepts, such as objects, attributes, and relationships, are not only easy to collect and annotate, but also semanti-046 cally consistent when expressed across modalities (Chen et al., 2019; Xu et al., 2020; Xie et al., 2020). 047 Hence, many traditional Vision–Language Models (VLMs) combine coarse- and fine-grained concept 048 annotations to better learn cross-modal consistent concrete concepts, thus improving VL alignment (Li et al., 2020; Zeng et al., 2021; Shen et al., 2022; Menon & Vondrick, 2023). However, they rely on additional components and loss functions to leverage different-grained concept annotations (e.g., 051 bounding box prediction), and optimize the multimodal comprehension ability of VLMs at different granularities separately through multitask learning. Moreover, they seldom explore the potential of 052 multi-grained concept annotations in multimodal generation tasks, such as text-to-image generation, let alone exploring both multimodal comprehension and generation tasks under the same framework. 054 As prior work demonstrate that concepts are crucial for VL alignment and combining coarse- and fine-055 grained concept annotations can better learn concepts, we argue that existing MLLMs should make bet-056 ter use of concepts by incorporating multi-grained concept annotations into their training. To this end, 057 we first merge four public, object detection datasets to construct a new multimodal dataset, MMGIC. 058 Our goal is not to replace or surpass existing image-caption datasets, but to address the lack of datasets with multi-grained concept annotations for MLLMs and explore their potential. Different from previous VLM work, we 1) provide multimodal annotations for images, including both textual forms 060 (caption, labels and label descriptions) and visual form (object regions), to enrich multi-grained con-061 cept annotations; 2) design a structured template to integrate multimodal multi-grained concept anno-062 tations into image-text interleaved documents, to leverage the complex context processing capability 063 of MLLMs; 3) instead of additional components or loss functions used in VLMs, train MLLMs with 064 an autoregressive discrete framework and predict the next visual or textual token in a multimodal dis-065 crete token sequence, to improve multimodal comprehension and generation ability across multiple 066 granularities **simultaneously**. This can reuse existing LLM training regimes and ensure the generality 067 and applicability of MMGIC to different MLLM frameworks. The key contributions are as follows: 068

- We introduce MMGIC, a new dataset with multimodal multi-grained concept annotations (Sec. 2). Under a general MLLM framework (Sec. 3), we show that MMGIC can help MLLMs better locate and learn concepts, thus aligning vision and language across multiple granularities (Sec. 4).
- We explore different data recipes for multi-grained concept annotations (Sec. 4.1 & 4.4). Our analyses show that multi-grained annotations can integrate and complement each other to help MLLMs ground concepts in the textual annotations to corresponding regions in images, thus improving the ability to understand and generate concepts.
- Through the evaluation of 12 benchmarks for multimodal comprehension and generation in both pre-training (Sec. 4.2) and supervised fine-tuning stages (Sec. 4.3), we explore the comparison and collaboration between MMGIC and coarse-grained image–caption data. We find that they each have their own strengths in depth and breadth of concept representation, and that appropriate curriculum learning strategies can effectively combine their strengths to further improve performance.
- 081 082

083

084

085

087

2 MMGIC DATASET

To be clear, our goal is not to supplant existing image–caption datasets, but to build a multimodal dataset with multi-grained concept annotations to address the lack of such datasets, and then explore its potential in MLLMs. We now introduce its collection, pre-processing, complement and construction.

087 2.1 COLLECTION AND PRE-PROCESSING

In this work, we focus on concrete concepts, especially objects, attributes of objects, and relationships
between objects. They are fundamental elements in VL learning and widely annotated in object detection datasets. Therefore, we collect four public large-scale human-annotated object detection datasets,
including Open Images (Kuznetsova et al., 2020), Objects365 (Shao et al., 2019), V3Det (Wang
et al., 2023a), and Visual Genome (Krishna et al., 2017). Images in these datasets are uploaded to
Flickr by real-world users and collected by dataset providers. They typically show complex scenes
with multiple objects, and are annotated with fine-grained category labels and object bounding boxes.
Comparing with widely-used coarse-grained image-caption datasets, fine-grained object annotations
provided in these datasets can help MLLMs locate and learn concepts in images.

098 Object Annotation Pre-processing. Fine-grained object annotations includes category labels and 099 bounding box coordinates for each object region. To accommodate varying aspect ratios of bounding 100 boxes and the requirement for a square input image, we crop a new larger square region S_i that 101 contains the original object region $R_i, R_i \subseteq S_i$, with their centers aligned as closely as possible. We 102 then update the annotations of S_i by integrating the category label of R_i with the category labels 103 of surrounding object regions. Notably, instead of transforming bounding box coordinates into tokens in the text (Chen et al., 2021; Liu et al., 2023c; Peng et al., 2023) or visual markers in the 104 image (Shtedritski et al., 2023; Yang et al., 2023; Yao et al., 2024), for each object, MMGIC directly 105 provides visual tokens of the cropped region S_i and textual tokens of fine-grained category labels and 106 location descriptions (Figure 1 3). Fine-grained cropped regions can help MLLMs locate and align 107 concepts in images and in annotations at a detailed level.

108	# Detailed Analysis of Objects in the Image		# Detailed Analysis of Objects in the Image	
109	Image: [image]	1	Image: [image]	100 miles
110	Caption: [caption]	1	Caption: a woman holding a baby in a chair a	at the beach.
111	Objects and their descriptions:		Objects and their descriptions:	
112	- [object_label]: [object_label_des]		- sunglasses: a type of protective eyewear	designed to prevent bright sunlight
113	Attributes of objects and their descriptions:		Attributes of objects and their descriptions:	
114	- [attribute_label] {[object_label]}: [attribute_label_des]	2	- sitting {woman, baby}: a posture where o	one's weight is supported by the buttocks
115	Relationships between objects and their descriptions:		Relationships between objects and their des	scriptions:
116	- [relationship_label] {[subject_label]-[object_label]}:		- wears {woman-sunglasses}: a situation i	n which a person (subject) is putting on or
117	[relationship_label_des]		carrying a certain item (object), such as cloth	ning or accessories
118	## Overview of Selected Object Regions		## Overview of Selected Object Regions	
110	### Overview of a Selected Object Region		### Overview of a Selected Object Region	### Overview of a Selected Object Region
119	Region: [image]	<u>.</u>	Region: [image] -	Region: [image] -
120	Location: [location]	3	Location: Bottom Right	Location: Top Center
121	Objects:		Objects:	Objects:
122	- [object_label]	H.	- baby	- sunglasses

125

126 127 128

129

Figure 1: Structured template (*Left*) and data example (*Right*) of MMGIC. Different colored text indicates template text, image placeholders, annotation placeholders and multi-grained concept annotations, respectively. Each image–text interleaved data sample will be tokenized into discrete tokens.

2.2 Multi-Grained Concept Annotation Complement

We follow LAION-COCO (Schuhmann et al., 2023) to automatically synthesize captions for all images (Figure 1 1) by BLIP-2 (Li et al., 2023b) and CLIP (Radford et al., 2021), since partial images are not annotated with captions. We do not use or synthesize captions for object regions to avoid introducing hallucinations.

134

135 Label Description Generation. Label descriptions are corresponding concept descriptions of con-136 crete concepts in the image, which convey understanding about a concept by visually observed details 137 and relevant knowledge. Previous works (Shen et al., 2022; Yao et al., 2022; Menon & Vondrick, 2023) successfully improve concept understanding by introducing label descriptions from WordNet (Miller, 138 1992) and LLMs (Brown et al., 2020). Inspired by them, we design prompt templates and several 139 human-annotated examples for each type of category labels (object, object attribute, and relationship 140 between objects), and then generate label-description pairs (Figure 1 $\begin{bmatrix} 2 \end{bmatrix}$) by GPT-4 (Achiam et al., 141 2023). We manually check them to ensure quality. Moreover, for better differentiation of polysemous 142 category labels, we manually check them based on the definitions in WordNet, and update them based 143 on the specific data samples, *e.g.*, "batter" \rightarrow "batter (ballplayer)" or "batter (cooking)".

144 145

146 2.3 CONSTRUCTION

147 Above steps transform four object detection datasets into MMGIC, a multimodal dataset with multi-148 grained concept annotations. Furthermore, we carefully design a structured template to integrate 149 multi-grained concept annotations into an image-text interleaved document. As shown in Figure 1 150 (*Left*), the structured template consists of: 1 coarse-grained image-annotation part: each image is 151 annotated with a short and general description of the whole image; 2 fine-grained image-annotation 152 part: concrete concepts (objects, attributes and relationships) present in the image are annotated with corresponding category labels and label descriptions; 3 fine-grained object-annotation part: 153 each object in the image is annotated with a cropped object region, object labels in the region, and a 154 location description. A data example of MMGIC is shown in Figure 1 (Right). 155

Different from previous VLM work that provide multiple isolated annotations for each image,
 MMGIC provides richer multimodal multi-grained concept annotations in both textual forms (caption,
 labels and label descriptions) and visual form (object regions) for each image, and integrates them into
 an image-text interleaved document by our structured template. This can leverage MLLMs' complex
 context processing capability to facilitate VL alignment across multiple granularities simultaneously
 under our MLLM framework. In a nutshell, MMGIC fills the gap in the MLLM field for datasets
 with multi-grained concept annotations. It contains 3.5M unique images, 23.9M unique object

regions, and 61.8M category label-description pairs. Based on MMGIC, we explore and analyse
 different data recipes for multi-grained concept annotations (Sec. 4.1 & 4.4), and further compare
 MMGIC with image-caption data (Sec. 4.2 & 4.3). More data details are shown in Appendix F.

3 FRAMEWORK

We introduce a general MLLM framework and its two training stages. Our goal is not to develop new frameworks, training objectives or benchmark SOTAs, but to explore the potential of multi-grained concept annotations for MLLMs under the general framework.

171 172 173

178

195

170

166

167 168

3.1 AN AUTOREGRESSIVE DISCRETE MLLM FRAMEWORK

Based on previous works (Ge et al., 2023b; Jin et al., 2023), we standardize a framework consisting
of several visual modules and a LLM with an extended VL vocabulary (Fig. 2). It is trained with an
autoregressive objective to generate predictions of the next token in a discrete sequence of image-text
interleaved tokens, and can support our exploration in multimodal comprehension and generation.

179 Visual Modules. Inherited from LaVIT (Jin et al., 2023), the visual modules consist of a visual encoder, 181 a visual tokenizer, a visual decoder and a diffusion model. The visual encoder is a pre-trained vision 182 transformer (Dosovitskiy et al., 2020; Sun et al., 183 2023b), which encodes an image into a sequence of visual embeddings. The visual tokenizer quantizes 185 these embeddings into a sequence of discrete visual tokens by a visual codebook (van den Oord et al., 187 2017). The visual decoder reconstructs predicted 188 visual tokens into a sequence of visual embeddings, 189 which are then taken as the condition of the diffusion 190 model (Sohl-Dickstein et al., 2015; Ho et al., 2020; 191 Rombach et al., 2022) to progressively generate target image pixels from a Gaussian noise. Overall, 192 visual modules can tokenize an image into a sequence 193



Figure 2: Illustration of our general MLLM framework. Only the LLM are loaded and partially fine-tuned during training.

194 of discrete visual tokens as input and decode predicted visual tokens back into an image.

LLM with an Extended VL Vocabulary. The LLM is inherited from LLaMA-2-7B (Touvron et al., 2023b), and its vocabulary is extended to support both textual and visual tokens. Since visual tokens in the VL vocabulary correspond one-to-one with visual latent codes in the visual codebook, instead of initializing visual token embeddings with the distribution of original textual token embeddings (Ge et al., 2023b; Jin et al., 2023), we directly initialize them with visual latent codes and a projection matrix (More details in Appendix D.1). The LLM can process a sequence of interleaved visual and textual tokens, and predict the next visual or textual token in an autoregressive manner.

To pursue simplicity, efficiency, and scalability (Ge et al., 2023b), instead of additional components or 203 loss functions used in VLMs, we train MLLMs with a single autoregressive objective. This allows for a 204 fair comparison between MMGIC and coarse-grained image-caption data under the same framework, 205 and ensures the generality and applicability of MMGIC to different MLLM frameworks. More impor-206 tantly, our framework facilitates VL alignment across multiple granularities simultaneously through 207 MMGIC and MLLMs' complex context processing capability. In all training stages, we follow com-208 mon practice (Ge et al., 2023a; Zhu et al., 2023a) to freeze most of the parameters and only tune partial 209 parameters of the LLM: a VL vocabulary, additional LoRA modules (Hu et al., 2021), norm layers and 210 a language modeling head layer, to greatly improve efficiency. We pre-tokenize images into discrete 211 visual token sequences and do not load all visual modules during training since they are frozen.

212 213

214

3.2 TRAINING STAGE 1: PRE-TRAINING (PT)

215 Similar to LLMs that tokenize text-only documents into discrete textual token sequences, for each data sample in MMGIC shown in Fig. 1 (*Right*) or an image–caption pair, our framework tokenize it

into a discrete token sequence consisting of interleaved visual and textual token sequences. To help
the LLM distinguish between two types of sequences, following LaVIT, we add two special tokens
[IMG] and [/IMG] to the vocabulary, and insert them before and after each visual token sequence,
respectively. Our framework is trained with an autoregressive objective to maximize the likelihood of
predicting the next visual or textual token. Details of training settings are provided in Appendix E.1.

- 221 222
- 3.3 TRAINING STAGE 2: SUPERVISED FINE-TUNING (SFT)

To align pre-trained MLLMs with natural language instructions, following previous works (Sun et al., 224 2023c;a; Liu et al., 2023c;b; 2024a; Zhu et al., 2023a; Hu et al., 2024a), we collect 1.21M samples 225 from several public datasets for supervised fine-tuning, including multimodal instruction datasets (Liu 226 et al., 2023b; Yu et al., 2023c; Zhang et al., 2023; Chen et al., 2023; LAION, 2024b; Brooks et al., 227 2023; Zhang et al., 2024) and text-only instruction datasets (Taori et al., 2023; Steven Tey, 2023), 228 and 1M samples from an aesthetic image-caption dataset (LAION, 2024a). We also play back 1M 229 samples from MMGIC to avoid forgetting the knowledge learned in the pre-training stage. Following 230 LLaVA (Liu et al., 2023c;b; 2024a), all datasets are transformed into a unified format, which consists 231 of a general system message and multiple instruction-answer pairs. Only answer tokens are accounted 232 in loss calculation. Details of instruction data are provided in Appendix H.

233 234

235

245

253

254

255

256

257

258

259

4 EXPERIMENT

236 Based on MMGIC and a general MLLM framework, we can explore the potential of multi-grained 237 concept annotations for MLLMs on various VL benchmarks in both pre-training and SFT stages. 238 Specifically, we follow previous works (Ge et al., 2023b; Liu et al., 2023c; Zhu et al., 2023a) to focus 239 on the zero-shot multimodal comprehension and generation capabilities, including image captioning 240 (COCO (Chen et al., 2015), NoCaps (Agrawal et al., 2019)), text-to-image generation (COCO (Chen et al., 2015), VIST (Huang et al., 2016)), visual question answering (VQAv2 (Goyal et al., 2017), 241 GQA (Hudson & Manning, 2019), VizWiz (Gurari et al., 2018)), comprehensive multi-choice bench-242 marks (POPE (Li et al., 2023c), MME (Yin et al., 2023), MMBench (Liu et al., 2023d), ScienceQA (Lu 243 et al., 2022b), SEED-Bench (Li et al., 2023a)). More evaluation details are provided in Appendix E.2. 244

246 4.1 DATA RECIPES FOR MULTI-GRAINED CONCEPT ANNOTATIONS

As mentioned in Sec. 2.3 and Fig. 1, multi-grained concept annotations include four components: coarse-grained image captions (C), fine-grained category labels (L), label descriptions (D) and object regions (R). Hence, we design four data recipes to investigate the importance of each component on zero-shot image captioning and image generation tasks and find the best recipe. As shown in Table 1:

- {0, 1, 2}-th rows: simply appending category labels to image captions does not help and even hurts the performance on both tasks. MLLMs may struggle to understand the association between images, captions and category labels, leading to confusion and treating category labels as additional noise. Whereas, label descriptions can strengthen the association between images and annotations, to help MLLMs align the concepts represented by labels with the concepts in the image, thus mitigating the performance drop. As shown in Figure 3 (*Top*), MMGIC(C) incorrectly identifies "accordion" as "electronic keyboard". While for MMGIC(CLD), label descriptions can help MLLMs better understand labels by **visual details** "pleated bellows and keyboard, box-like" and **relevant knowledge** "portable".
- {2,3}-th rows: object regions can further complement other annotations to help MLLMs better locate and align concepts in images and in annotations at a fine-grained level, thus significantly improving the performance on both tasks. As shown in Figure 3 (*Bottom*), MMGIC(CLD) incorrectly identifies "laying" as "sitting" and "on top of a toilet" as "on the floor next to a toilet". While for MMGIC(CLDR), object regions can help MLLMs ground concepts in the textual annotations to corresponding regions in the image, especially in terms of instance interaction and spatial relationship here, which is also consistent with meso analyses in Section 4.4.

We further analyse generated images by MLLMs pre-trained with different data recipes in Figure 4 (*Left*). The input prompt and the label description of "bagel" are shown in the top and generated images are shown in the bottom. We found that while MMGIC(C) could accurately generate "on top of a white plate on top of a table", it fails to correctly generate "a bagel", let alone "a blueberry bagel". Table 1: Zero-shot evaluation of different data recipes for MMGIC. C: image caption; L: category labels; D: label descriptions; R: object regions; CLIP-T/I: image-text or image-image similarity via CLIP; \downarrow : lower is better, otherwise higher is better. The best results are **bold**.

Data Recipes	Image C COCO	Captioning NoCaps	MS	Imag -COCO-3	ge Genera 0K	ation VIST			
CLDR	CIDEr	CIDEr	FID (\downarrow)	CLIP-T	CLIP-I	$FID~(\downarrow)$	CLIP-I		
0 🗸	113.64	99.11	7.20	30.81	71.62	67.61	62.22		
1 🗸 🗸	113.67	96.80	10.52	30.43	71.09	71.67	61.96		
2 🗸 🇸 🇸	116.02	98.61	8.92	30.89	71.60	51.89	64.30		
3 / / / /	118.30	102.01	7.36	31.57	72.24	35.33	66.10		

ASS TOO	MMGIC(C): MMGIC(CLD):	an older man with glasses playing an electronic keyboard. an older man playing an accordion in a restaurant.
	 electronic keyboa accordion: a po 	ard: a versatile musical instrument that imitates the sound which consists of a set of keys rtable musical instrument known for its pleated bellows and keyboard, its box-like shape
	MMGIC(CLD): MMGIC(CLDR):	a cat sitting on the floor next to a toilet in a bathroom. a cat laying on top of a toilet in a bathroom.
P	 sitting: a po laying: a po 	osture where one's weight is supported by the buttocks rather than the feet osture where an individual or an animal is reclining or resting in a horizontal posture

Figure 3: Comparison of generated captions by MLLMs pre-trained with different data recipes. MMGIC(C), MMGIC(CLD) and MMGIC(CLDR) denote the $\{0, 2, 3\}$ -th data recipes in Table 1, respectively. The bottom right of each example shows associated label-description pairs from MMGIC.

Label descriptions in MMGIC(CLD) can provide more visual details about both "bagel" and "blue-berry", which helps MLLM better understand and generate the concept of "bagel" with the feature of "round bread product", and also add four "blueberries" to th "bagel". Furthermore, with the help of object regions in MMGIC(CLDR), MLLMs can correctly understand and generate a "bagel" with the feature of "its hole in the center". We also show some emergent abilities of MLLMs pre-trained with MMGIC in Figure 4 (Right). Notably, MMGIC does not contain any samples about image editing and in-context image synthesis. Leveraging image-text interleaved documents with multi-grained concept annotations, the top two examples show that MLLMs can precisely understand editing instructions and perform appropriate editing, while the bottom example shows that MLLMs can synthesize image pre-cisely based on the image-text interleaved sequences. More analyses are provided in Appendix C.2.

In summary, multi-grained concept annotations are not isolated from each other. With our structured template and MLLM framework, they can integrate into multimodal documents and complement each other to help MLLMs better learn concepts, thus improving the ability to understand and generate concepts. We take the 3-rd data recipe as the default data recipe and denote it simply as MMGIC.

4.2 COMPARISON AND COLLABORATION BETWEEN MMGIC AND IMAGE-CAPTION DATA

MMGIC is constructed from public object detection datasets mainly covers common concepts. Widely-used coarse-grained image-caption data, e.g., Conceptual Captions (Sharma et al., 2018) and LAION-5B (Schuhmann et al., 2022a), typically cover diverse, scalable but noisy concepts. To explore the potential of MMGIC, we then investigate the comparison and collaboration between MMGIC with image-caption data. To strike a balance between increasing concept breadth, reducing data noise and improving efficiency, we collect several large-scale public image-caption datasets (Ordonez et al., 2011; Sharma et al., 2018; Changpinyo et al., 2021; Schuhmann et al., 2022a;b; Sun et al., 2024), and follow LLaVA (Liu et al., 2023c) to first filter them by the frequency of noun-phrases extracted by spaCy from their given synthesized captions, and then automatically synthesize captions same as MMGIC. We name this dataset as IC, where 52M unique images are collected and selected, almost 15 times more than MMGIC. More details are provided in Appendix G. As shown in Table 2:



Figure 4: Comparison of generated images by MLLMs pre-trained with different data recipes (*Left*) and image editing and multimodal in-context image synthesis examples (*Right*).

Table 2: Zero-shot evaluation of comparison and collaboration between MMGIC and IC. IC-PART: select the same number of samples as MMGIC from IC; MMGIC+IC: joint training; MMGIC \rightarrow IC: train on MMGIC first and then on IC. The best results are **bold** and the second-best are underlined.

	1			I	-	~					
			Image Captioning		Image Generation						
	Training Data	COCO	NoCaps	MS-COCO-30K			VIST				
			CIDE					CLIDI			
		CIDEr	CIDEr	FID (↓)	CLIP-I	CLIP-I	FID (↓)	CLIP-I			
0	IC-Part	95.74	85.00	11.62	29.94	68.21	64.03	63.41			
1	IC	104.15	92.24	7.65	31.40	70.27	41.65	65.06			
2	MMGIC	118.30	102.01	7.36	31.57	72.24	35.33	66.10			
3	MMGIC+IC	106.45	92.98	7.11	31.65	70.93	36.54	65.89			
4	$MMGIC \rightarrow IC$	105.62	93.77	7.29	31.57	70.26	37.45	65.88			
5	$IC \rightarrow MMGIC$	120.84	105.59	7.13	31.96	71.54	37.13	65.62			
6	$MMGIC+IC \rightarrow MMGIC$	121.22	105.33	7.22	<u>31.91</u>	71.75	36.23	65.79			

• {0,1,2}-th rows: comparing with IC, MMGIC can help MLLMs achieve **significantly** better performance on both tasks even with a much **smaller** number of samples, which demonstrates the effectiveness of multi-grained concept annotations in concept understanding and generation.

• {3, 4, 5, 6}-th rows: naturally, we try to improve the performance and concept breadth of MLLMs by exploring the collaboration between MMGIC and IC. However, simply joint training MMGIC and IC achieves better performance than IC, but significantly worse than MMGIC, especially on image captioning and VIST. We then follow the curriculum learning strategy (Mc-Cann et al., 2019) to train them in different orders. Interestingly, training on IC first and then on MMGIC achieves significantly better performance than all the above strategies on image captioning and partial metrics of image generation. This is consistent with recent findings (Hu et al., 2024b; Li et al., 2024a) that training with **high-quality** data **later** in the pre-training phase leads to better performance. Moreover, considering that the noise in IC still cause a slight performance drop on VIST ({2, 5}-th rows), we first jointly train on MMGIC and IC to alleviate the effect of noise in IC, and then on MMGIC, eventually achieving the best average performance (6-th row).

Overall, by exploring comparison and collaboration between MMGIC with large-scale coarse-grained
 image-caption data in the pre-training stage, we demonstrate that MMGIC achieves better performance than IC on both tasks, and their appropriate collaboration can further improve the average performance. We present {1, 2, 6}-th rows as three baselines, and denote them as MLLM-{IC, MMGIC, MMGIC & IC}. More discussions about collaboration strategies are provided in Appendix B.2.

4.3 EVALUATION ON DOWNSTREAM VISION–LANGUAGE BENCHMARKS AFTER SFT

To further explore the potential of MMGIC dataset on various downstream VL benchmarks, we then
 perform SFT (Section 3.3) on our three baselines. Technically, different training datasets, training set tings, framework, etc., lead to non-comparable and unfair comparisons of our baselines with existing
 MLLMs. Hence, we show SOTA MLLMs in gray as upper bound references, and their computation
 and data resources are extremely expensive and large (well over 10 times that of our work).

Table 3: Zero-shot evaluation on multimodal comprehension benchmarks after SFT. MLLMs in Group a are for comprehension only, while MLLMs in Group b are for both comprehension and generation. MMB: MMBench; SQA^I: ScienceQA-IMG; SEED^I: SEED-Bench-IMG; *: w/o SFT.

Model	Image C COCO	Captioning NoCaps	VQAv2	VQA GQA	VizWiz	POPE	Multi-Ch MME	oice Be MMB	nchmar SQA ^I	k SEED ^I
SOTA MLLMs as upper b	ound refe	rences, with	h more tr	aining a	lata or tr	ainable	paramete	ers, not	compar	able
LLaVA-1.5-7B			78.50	62.00	50.00	85.90	1826.80	65.20	66.80	65.80
a Emu-I-14B	120.40	108.80	62.00	46.00	38.30					58.00
Emu2-Chat-37B			84.90	65.10	54.90			62.40		68.90
VL-GPT-I-7B	133.70		67.20	51.50	38.90					
b LaVIT-v2-7B*	133.30	112.00	68.30	47.90	41.00					
SEED-LLaMA-I-8B	124.50	97.78	66.20	52.24	55.10	79.92	1497.53	52.58	60.24	51.50
Our comparable baseline	Our comparable baselines									
0 MLLM-IC	108.13	92.71	70.28	56.02	52.62	81.14	1646.71	59.54	65.94	58.41
1 MLLM-MMGIC	119.35	104.19	70.13	<u>56.84</u>	51.14	83.25	1636.47	58.51	65.79	60.03
2 MLLM-MMGIC & IC	122.31	106.97	70.57	56.97	52.66	85.09	1668.19	59.88	66.24	60.75

Table 4: Zero-shot evaluation on multimodal generation benchmarks after SFT. MLLMs in Group a are for generation only, while MLLMs in Group b are for both comprehension and generation.

Madal	MS	-COCO-3	VIST		
Model	$FID(\downarrow)$	CLIP-T	CLIP-I	$FID(\downarrow)$	CLIP-I
SOTA MLLMs as upper bound references, r	iot compa	vrable			
KOSMOS-G (Pan et al., 2023)	10.99				
a GILL (Koh et al., 2024)	12.20				64.10
Emu2-Gen-37B (Sun et al., 2023a)		29.70	68.60		
VL-GPT-I-7B (Zhu et al., 2023a)	11.53				
b LaVIT-v2-7B* (Jin et al., 2023)	7.10	31.93	71.06	34.76	68.41
SEED-LLaMA-I-8B (Ge et al., 2023b)	16.66	29.52	69.22	43.69	65.21
Our comparable baselines					
0 MLLM-IC	8.11	30.90	70.72	38.19	65.37
1 MLLM-MMGIC	6.79	31.63	72.44	34.32	67.66
2 MLLM-MMGIC & IC	7.29	<u>31.54</u>	72.03	<u>34.39</u>	<u>67.19</u>

Zero-shot Multimodal Comprehension. As shown in Table 3, MLLM-MMGIC & IC achieves the best performance on all 10 benchmarks compared to the other two baselines, and even outperforms some SOTA MLLMs with more training data or full-param training or larger LLMs on some bench-marks. Moreover, even with less than 4M pre-training data compared to MLLM-IC with 52M pre-training data, MLLM-MMGIC significantly outperforms MLLM-IC on the benchmarks that inspect in-depth understanding of common concrete concepts, e.g., COCO, NoCaps, POPE, SEED-Bench. In contrast, MLLM-IC outperforms MLLM-MMGIC on benchmarks that require a broader understand-ing of concrete concepts, e.g., VizWiz, MME, MMBench. More importantly, MLLM-MMGIC & IC can effectively **combine** the strengths of both in terms of **depth** and **breadth** of concept representation and further improve performance, e.g., 3.95% and 2.34% accuracy improvements on POPE and SEED-Bench. We further analyse in terms of dataset statistics and concept overlap in App. C.1.

Zero-shot Multimodal Generation. As shown in Table 4, for two text-to-image generation bench-marks that focus on common concrete concepts, MLLM-MMGIC achieves the best performance, and matches or even outperforms some SOTA MLLMs on some metrics. This demonstrates that fine-grained category labels, label descriptions and object regions can help MLLMs better learn and generate concepts. Besides, the noise introduced by IC (discussed in Section 4.2) in the pre-training stage may not be well alleviated by SFT, thus MLLM-MMGIC & IC achieves the second-best performance. More results and analyses on image editing are provided in Appendix C.2.



Figure 5: Analysis on 8 dimensions of SEED-Bench-IMG. Left: the performance of MLLM-MMGIC trained with different-grained concept annotations from MMGIC. Right: corresponding case studies. CG, FG, and MG denote MLLMs trained with coarse-, fine-, and multi-grained concept annotations from MMGIC, respectively.

4.4 THE IMPACT OF DIFFERENT-GRAINED CONCEPT ANNOTATIONS

To delve into the potential of multi-grained concept annotations for MLLMs, we conduct meso analyses on SEED-Bench (Li et al., 2023a), a large-scale multi-choice benchmark for multimodal comprehension. We follow common practice (Ge et al., 2023b) to select 9 image evaluation dimensions, *i.e.*, about 14K questions, and denote it as SEED-Bench-IMG. We ignore the "*Text Recognition*" dimension due to noisy data and small scale (only 85), leading to result fluctuations. "*Scene Understanding, Visual Reasoning*" dimensions focus on the holistic understanding and cross-modal reasoning of the image, while the other 6 dimensions focus on in-depth understanding of concrete concepts in the image. "Coarse-grained" (CG) means that only image captions from MMGIC are used. "Fine-grained" (FG) uses category labels, label descriptions and object regions from MMGIC. "Multi-grained" (MG) means that both above are used.

Quantitative Analysis. We first explore the impact of different-grained concept annotations from MMGIC on each evaluation dimension of SEED-Bench-IMG in Figure 5 (*Left*). For overall accuracy, FG significantly improves performance over CG by 1.39 points, and multi-grained can further improve by 1.4 points. For each evaluation dimension:

- FG provides **deeper** understanding of concepts than CG. Compared with coarse-grained image captions, fine-grained concept annotations can help MLLMs better **understand** and **locate** concepts in images, and **recognize** relationships between concepts, especially on the "*Instance Identity, Spatial Relation, Instance Counting, Instance Interaction, Visual Reasoning*" dimensions.
- MG can facilitate **collaboration** between concept annotations of different granularities, thus fully **integrating** each other's strengths and achieving further improvements in all dimensions, especially on the "*Scene Understanding, Instance Attribute, Instance Counting, Instance Interaction, Visual Reasoning*" dimensions. This demonstrates that our structured template for MMGIC can help MLLMs better utilize multi-grained concept annotations to learn concepts and thus promote vision–language alignment **across** multiple granularities **simultaneously**.

478 Qualitative Analysis. To better analyse the advantages of different-grained concept annotations in
479 MMGIC for MLLMs, we provide corresponding qualitative analysis in Figure 5 (*Right*). More case
480 studies are provided in Appendix C.4.

Example 1 "Instance Identity": FG provides deeper understanding of concepts than CG. While CG provides a holistic description of the image, it cannot help MLLMs distinguish between "Heels" and "Boots". However, FG provides visual details "lifting the heel above the toes" and "covers the foot and extends up the leg" by label–description pairs, which help MLLMs distinguish between similar concepts. Object region-annotation pairs further help MLLMs better locate and learn these concepts. Hence, FG helps capture visual details and identify concepts correctly.

• Example 2 "Scene Understanding": CG provides more **holistic** understanding of the image than FG. While the lady in the image is holding an "umbrella", her surroundings show that the sky is covered with clouds, and it is not raining. MLLMs trained with FG seem to be too focused on the "umbrella" and ignore the overall scene of the image, while CG can help MLLMs predict the correct answer by better understanding the **global context**.

491 • Example 3 "Visual Reasoning": MG can **combine** the advantages of CG and FG to improve 492 visual reasoning. This image shows a football match with a player in "yellow uniforms" lying on the ground and a player in "white uniforms" celebrating as the audience cheers him on. 493 Rich objects and details confuse MLLMs trained with FG to confirm the main focus of the 494 image, while MLLMs trained with CG perceive the visually more prominent "yellow" player 495 as the main focus. However, with our structured template and MLLM framework, MG enables 496 different-grained concept annotations to **integrate** and **complement** each other, thus better 497 understanding and reasoning about the image from both global context and local details. 498 MLLMs trained with MG correctly identify the "white" player as the main focus of the image. 499

We also conduct experiments on self-generated annotations in evaluation, MMGIC directly as SFT data, the impact of the nature of image-text interleaved, text-only performance in Appendix C.3 & C.

5 RELATED WORK

505 MLLMs have emerged recently including ones (Sun et al., 2023c; Ge et al., 2023b; Jin et al., 2023; Zhu et al., 2023a; Dong et al., 2024) that propose multimodal generalists capable of both multimodal 506 comprehension and generation. However, existing MLLMs often overlook the importance of concepts 507 and only utilize coarse-grained concept annotations, e.g., image captions, which may limit vision-508 language alignment. Factually, many traditional VLMs recognized the importance of concepts in 509 vision-language learning. To better utilize concepts in vision-language learning and improve perfor-510 mance, they incorporated fine-grained concept annotations into coarse-grained image captions. For ex-511 ample, object labels (Li et al., 2020; Zhang et al., 2021), label descriptions (Shen et al., 2022; Yao et al., 512 2022; Menon & Vondrick, 2023; Li et al., 2024b), region descriptions (Zeng et al., 2021) and object 513 regions (Zeng et al., 2021; Li et al., 2022b). Especially, Oscar (Li et al., 2020) appends fine-grained ob-514 ject labels detected in the image after the coarse-grained image caption to simplify semantic alignment 515 between vision and language. X-VLM (Zeng et al., 2021) proposes to align visual concepts (images 516 and object regions) with coarse-grained image captions and fine-grained object labels and object re-517 gion descriptions in multi-granularity. In this paper, different from previous VLM work, we collect and construct rich multimodal multi-grained concept annotations in MMGIC dataset. Without additional 518 components or loss functions, we design a structured template to leverage the advantages of MLLMs 519 under an autoregressive discrete framework. Through evaluation of 12 multimodal comprehension and 520 generation benchmarks, as well as the comparison, combination, and analysis of MMGIC and image-521 caption data, for the first time, we explore and demonstrate the potential of MMGIC in MLLMs. 522

523 524

535

486

487

488

489

490

500

501 502

504

6 CONCLUSION AND FUTURE WORK

We introduce a new multimodal dataset, MMGIC, providing multi-grained concept annotations 526 in both textual and visual form. MMGIC allows us to measure the potential of appropriate use of 527 concepts. With our structured template and autoregressive discrete framework, multi-grained concept 528 annotations can integrate and complement each other to help MLLMs better locate and learn concepts, 529 thereby aligning vision and language across multiple granularities simultaneously. Furthermore, 530 MMGIC and coarse-grained image-caption data each have their own strengths in depth and breadth 531 of concept representation, and appropriately combining them can effectively integrate their strengths 532 to further improve performance. We hope to inspire future research to further explore the potential 533 of multi-grained concept annotations by incorporating more different types of annotations, scaling 534 up data by automatic synthesis, scaling up (concrete and even abstract) concepts, and other VL tasks.

536 537 REFERENCES

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman,
 Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report.
 ArXiv preprint, abs/2303.08774, 2023. URL https://arxiv.org/abs/2303.08774.

- Harsh Agrawal, Peter Anderson, Karan Desai, Yufei Wang, Xinlei Chen, Rishabh Jain, Mark Johnson,
 Dhruv Batra, Devi Parikh, and Stefan Lee. nocaps: novel object captioning at scale. In 2019 *IEEE/CVF International Conference on Computer Vision, ICCV 2019, Seoul, Korea (South), October 27 November 2, 2019*, pp. 8947–8956. IEEE, 2019. doi: 10.1109/ICCV.2019.00904.
 URL https://doi.org/10.1109/ICCV.2019.00904.
- Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel
 Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. Flamingo: a visual language
 model for few-shot learning. Advances in neural information processing systems, 35:23716–23736, 2022.
- Anas Awadalla, Le Xue, Oscar Lo, Manli Shu, Hannah Lee, Etash Kumar Guha, Matt Jordan, Sheng
 Shen, Mohamed Awadalla, Silvio Savarese, et al. Mint-1t: Scaling open-source multimodal data
 by 10x: A multimodal dataset with one trillion tokens. *arXiv preprint arXiv:2406.11271*, 2024.
- Peter Blouw, Eugene Solodkin, Paul Thagard, and Chris Eliasmith. Concepts as semantic pointers: A
 framework and computational model. *Cognitive science*, 40(5):1128–1162, 2016.
- Tim Brooks, Aleksander Holynski, and Alexei A Efros. Instructpix2pix: Learning to follow image editing instructions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 18392–18402, 2023.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhari-559 wal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, 561 Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. In Hugo Larochelle, Marc'Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, 565 and Hsuan-Tien Lin (eds.), Advances in Neural Information Processing Systems 33: Annual 566 Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 567 2020, virtual, 2020. URL https://proceedings.neurips.cc/paper/2020/hash/ 568 1457c0d6bfcb4967418bfb8ac142f64a-Abstract.html.
- Likun Cai, Zhi Zhang, Yi Zhu, Li Zhang, Mu Li, and Xiangyang Xue. Bigdetection: A large-scale benchmark for improved object detector pre-training. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 4777–4787, 2022.
- 573 Susan Carey. The origin of concepts. *Journal of Cognition and Development*, 1(1):37–41, 2000.

578

579

580

- Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and
 Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 9650–9660, 2021.
 - Soravit Changpinyo, Piyush Sharma, Nan Ding, and Radu Soricut. Conceptual 12m: Pushing web-scale image-text pre-training to recognize long-tail visual concepts. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 3558–3568, 2021.
- Hui Chen, Guiguang Ding, Zijia Lin, Sicheng Zhao, and Jungong Han. Cross-modal image-text retrieval with semantic consistency. In *Proceedings of the 27th ACM International Conference on Multimedia, MM 2019, Nice, France, October 21-25, 2019*, pp. 1749–1757, 2019. doi: 10.1145/3343031.3351055. URL https://doi.org/10.1145/3343031.3351055.
- Lin Chen, Jisong Li, Xiaoyi Dong, Pan Zhang, Conghui He, Jiaqi Wang, Feng Zhao, and Dahua
 Lin. Sharegpt4v: Improving large multimodal models with better captions. *ArXiv preprint*, abs/2311.12793, 2023. URL https://arxiv.org/abs/2311.12793.
- Ting Chen, Saurabh Saxena, Lala Li, David J Fleet, and Geoffrey Hinton. Pix2seq: A language modeling framework for object detection. In *International Conference on Learning Representations*, 2021.
- Xinlei Chen, Hao Fang, Tsung-Yi Lin, Ramakrishna Vedantam, Saurabh Gupta, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco captions: Data collection and evaluation server. ArXiv preprint, abs/1504.00325, 2015. URL https://arxiv.org/abs/1504.00325.

594 595 596	Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. Palm: Scaling language modeling with pathways. <i>Journal of Machine Learning Research</i> , 24(240):1–113,
597	2023.
598	Louise Connell Dermot Lynott, and Brigny Banks. Interoception: the forgotten modality in perceptual
599	grounding of abstract and concrete concepts. <i>Philosophical Transactions of the Royal Society B</i> :
600	Biological Sciences, 373(1752):20170143, 2018.
602	
602	Wenliang Dai, Nayeon Lee, Boxin Wang, Zhuoling Yang, Zihan Liu, Jon Barker, Tuomas Rintamaki,
604 605	Mohammad Shoeybi, Bryan Catanzaro, and Wei Ping. Nvlm: Open frontier-class multimodal llms. <i>arXiv preprint arXiv:2409.11402</i> , 2024.
606	Mostafa Dehghani, Josip Djolonga, Basil Mustafa, Piotr Padlewski, Jonathan Heek, Justin Gilmer,
607	Andreas Peter Steiner, Mathilde Caron, Robert Geirhos, Ibrahim Alabdulmohsin, et al. Scaling
608	vision transformers to 22 billion parameters. In International Conference on Machine Learning,
609	pp. 7480–7512. PMLR, 2023.
610	Runnei Dong Chunnui Han Yuang Peng Zekun Oi, Zheng Ge, Jinrong Vang, Liang Zhao, Jianjian
611	Sun, Hongyu Zhou, Haoran Wei, Xiangwen Kong Xiangyu Zhang Kaisheng Ma and Li Yi
612	DreamLLM: Synergistic multimodal comprehension and creation. In <i>The Twelfth International</i>
613	Conference on Learning Representations, 2024. URL https://openreview.net/forum?
614	id=y01KGvd9Bw.
615	
616	Alexey Dosovilskiy, Lucas Beyer, Alexander Kolesnikov, Dirk weissendorn, Alaonua Zhai, Thomas Unterthiner Mostafa Debahani Matthias Minderer Georg Heigold Sylvain Gelly et al. An image
617	is worth 16x16 words: Transformers for image recognition at scale. In <i>International Conference</i>
618	on Learning Representations, 2020.
619	
620	Yuying Ge, Yixiao Ge, Ziyun Zeng, Xintao Wang, and Ying Shan. Planting a seed of vision in large
621 622	language model. ArXiv preprint, abs/2307.08041, 2023a. URL https://arxiv.org/abs/ 2307.08041.
623 624 625 626	Yuying Ge, Sijie Zhao, Ziyun Zeng, Yixiao Ge, Chen Li, Xintao Wang, and Ying Shan. Making llama see and draw with seed tokenizer. <i>ArXiv preprint</i> , abs/2310.01218, 2023b. URL https://arxiv.org/abs/2310.01218.
627 628	Joseph A Goguen. The logic of inexact concepts. Synthese, 19(3/4):325-373, 1969.
629 630 631 632 633	Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. Making the V in VQA matter: Elevating the role of image understanding in visual question answering. In 2017 <i>IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017</i> , pp. 6325–6334. IEEE Computer Society, 2017. doi: 10.1109/CVPR.2017.670. URL https://doi.org/10.1109/CVPR.2017.670.
634	Agrim Gupta, Piotr Dollár, and Ross B. Girshick. LVIS: A dataset for large vocabulary instance
635	segmentation. In <i>IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long</i>
636	Beach, CA, USA, June 16-20, 2019, pp. 5356–5364. Computer Vision Foundation / IEEE, 2019. doi:
637	10.1109/CVPR.2019.00550. URL http://openaccess.thecvf.com/content_CVPR_
638	2019/html/Gupta_LVIS_A_Dataset_for_Large_Vocabulary_Instance_
639	Segmentation_CVPK_2019_paper.html.
640	Danna Gurari, Qing Li, Abigale J. Stangl, Anhong Guo, Chi Lin, Kristen Grauman, Jiebo Luo, and
041	Jeffrey P. Bigham. Vizwiz grand challenge: Answering visual questions from blind people. In
642	2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018, Salt Lake City,
643	UT, USA, June 18-22, 2018, pp. 3608–3617. IEEE Computer Society, 2018. doi: 10.1109/CVPR.
645	2018.00380. UKL http://openaccess.thecvf.com/content_cvpr_2018/html/
646	Gurarr_vrzwrz_Grand_Challenge_CVPK_2018_paper.html.
647	Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob

647Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob
Steinhardt. Measuring massive multitask language understanding, 2021.

- Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In Hugo Larochelle, Marc'Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin (eds.), Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual, 2020. URL https://proceedings.neurips.cc/paper/2020/hash/ 4c5bcfec8584af0d967f1ab10179ca4b-Abstract.html.
- Edward J Hu, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. Lora: Low-rank adaptation of large language models. In *International Conference on Learning Representations*, 2021.
- Jinyi Hu, Yuan Yao, Chongyi Wang, SHAN WANG, Yinxu Pan, Qianyu Chen, Tianyu Yu, Hanghao Wu, Yue Zhao, Haoye Zhang, Xu Han, Yankai Lin, Jiao Xue, dahai li, Zhiyuan Liu, and Maosong Sun. Large multilingual models pivot zero-shot multimodal learning across languages. In *The Twelfth International Conference on Learning Representations*, 2024a. URL https://openreview.net/forum?id=Kuh5qgCGCp.
- Shengding Hu, Yuge Tu, Xu Han, Chaoqun He, Ganqu Cui, Xiang Long, Zhi Zheng, Yewei Fang,
 Yuxiang Huang, Weilin Zhao, et al. Minicpm: Unveiling the potential of small language models
 with scalable training strategies. ArXiv preprint, abs/2404.06395, 2024b. URL https://arxiv.
 org/abs/2404.06395.
- 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. Visual storytelling. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 1233–1239, San Diego, California, 2016. Association for Computational Linguistics. doi: 10.18653/v1/N16-1147. URL https://aclanthology.org/N16-1147.
- Drew A. Hudson and Christopher D. Manning. 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*, pp. 6700–6709. Computer Vision Foundation / IEEE, 2019. doi: 10.1109/CVPR.2019.00686. URL http://openaccess.thecvf.com/content_CVPR_2019/html/Hudson_GQA_A_New_Dataset_for_Real-World_Visual_Reasoning_and_Compositional_CVPR_2019_paper.html.
- Yang Jin, Kun Xu, Liwei Chen, Chao Liao, Jianchao Tan, Bin Chen, Chenyi Lei, An Liu, Chengru Song, Xiaoqiang Lei, et al. Unified language-vision pretraining with dynamic discrete visual tokenization. *ArXiv preprint*, abs/2309.04669, 2023. URL https://arxiv.org/abs/2309.04669.
- Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete
 Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 4015–4026, 2023.
 - Jing Yu Koh, Daniel Fried, and Russ R Salakhutdinov. Generating images with multimodal language models. *Advances in Neural Information Processing Systems*, 36, 2024.
- Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie
 Chen, Yannis Kalantidis, Li-Jia Li, David A Shamma, et al. Visual genome: Connecting language
 and vision using crowdsourced dense image annotations. *International journal of computer vision*,
 123(1):32–73, 2017.
- Alina Kuznetsova, Hassan Rom, Neil Alldrin, Jasper Uijlings, Ivan Krasin, Jordi Pont-Tuset, Shahab
 Kamali, Stefan Popov, Matteo Malloci, Alexander Kolesnikov, Tom Duerig, and Vittorio Ferrari.
 The open images dataset v4: Unified image classification, object detection, and visual relationship
 detection at scale. *IJCV*, 2020.
- 700 701

689

690

691

LAION. laion-coco-aesthetic. https://huggingface.co/datasets/guangyil/ laion-coco-aesthetic, 2024a.

702 703 704	GPT-4V LAION. Laion-gpt4v. https://huggingface.co/datasets/laion/gpt4v-dataset,2024b.
705 706 707 708	Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut. ALBERT: A lite BERT for self-supervised learning of language representations. In 8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020. OpenReview.net, 2020. URL https://openreview.net/forum?id=H1eA7AEtvS.
709 710 711 712	Hugo Laurençon, Lucile Saulnier, Léo Tronchon, Stas Bekman, Amanpreet Singh, Anton Lozhkov, Thomas Wang, Siddharth Karamcheti, Alexander Rush, Douwe Kiela, et al. Obelics: An open web-scale filtered dataset of interleaved image-text documents. <i>Advances in Neural Information</i> <i>Processing Systems</i> , 36, 2024.
713 714 715 716	Bo Li, Hao Zhang, Kaichen Zhang, Dong Guo, Yuanhan Zhang, Renrui Zhang, Feng Li, Ziwei Liu, and Chunyuan Li. Llava-next: What else influences visual instruc- tion tuning beyond data?, May 2024a. URL https://llava-vl.github.io/blog/ 2024-05-25-llava-next-ablations/.
717 718 719 720	Bohao Li, Rui Wang, Guangzhi Wang, Yuying Ge, Yixiao Ge, and Ying Shan. Seed-bench: Bench- marking multimodal llms with generative comprehension. <i>ArXiv preprint</i> , abs/2307.16125, 2023a. URL https://arxiv.org/abs/2307.16125.
721 722 723	Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image pre- training for unified vision-language understanding and generation. <i>ArXiv preprint</i> , abs/2201.12086, 2022a. URL https://arxiv.org/abs/2201.12086.
724 725 726 727	Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In <i>International conference on machine learning</i> , pp. 19730–19742. PMLR, 2023b.
728 729	Liunian Li, Zi-Yi Dou, Nanyun Peng, and Kai-Wei Chang. Desco: Learning object recognition with rich language descriptions. <i>Advances in Neural Information Processing Systems</i> , 36, 2024b.
730 731 732 733	Liunian Harold Li, Pengchuan Zhang, Haotian Zhang, Jianwei Yang, Chunyuan Li, Yiwu Zhong, Lijuan Wang, Lu Yuan, Lei Zhang, Jenq-Neng Hwang, et al. Grounded language-image pre-training. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pp. 10965–10975, 2022b.
734 735 736	Xiaotong Li, Fan Zhang, Haiwen Diao, Yueze Wang, Xinlong Wang, and Ling-Yu Duan. Densefusion-1m: Merging vision experts for comprehensive multimodal perception. <i>arXiv preprint arXiv:2407.08303</i> , 2024c.
737 738 739 740	Xiujun Li, Xi Yin, Chunyuan Li, Pengchuan Zhang, Xiaowei Hu, Lei Zhang, Lijuan Wang, Houdong Hu, Li Dong, Furu Wei, et al. Oscar: Object-semantics aligned pre-training for vision–language tasks. In <i>European Conference on Computer Vision</i> , pp. 121–137. Springer, 2020.
741 742 743 744 745	Yifan Li, Yifan Du, Kun Zhou, Jinpeng Wang, Xin Zhao, and Ji-Rong Wen. Evaluating object hallucination in large vision-language models. In Houda Bouamor, Juan Pino, and Kalika Bali (eds.), <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pp. 292–305, Singapore, 2023c. Association for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-main.20. URL https://aclanthology.org/2023.emnlp-main.20.
746 747 748 749	Zhang Li, Biao Yang, Qiang Liu, Zhiyin Ma, Shuo Zhang, Jingxu Yang, Yabo Sun, Yuliang Liu, and Xiang Bai. Monkey: Image resolution and text label are important things for large multi-modal models. <i>ArXiv preprint</i> , abs/2311.06607, 2023d. URL https://arxiv.org/abs/2311.06607.
750 751 752 753	Ji Lin, Hongxu Yin, Wei Ping, Yao Lu, Pavlo Molchanov, Andrew Tao, Huizi Mao, Jan Kautz, Mohammad Shoeybi, and Song Han. Vila: On pre-training for visual language models. <i>ArXiv</i> preprint, abs/2312.07533, 2023. URL https://arxiv.org/abs/2312.07533.
754 755	Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In <i>European conference on computer vision</i> , pp. 740–755. Springer, 2014.

756	Fuxiao Liu, Tianrui Guan, Zongxia Li, Lichang Chen, Yaser Yacoob, Dinesh Manocha, and Tianvi
757	Zhou. Hallusionbench: You see what you think? or you think what you see? an image-context
758	reasoning benchmark challenging for gpt-4v (ision), llava-1.5, and other multi-modality models.
759	arXiv preprint arXiv:2310.14566, 2023a.

- Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning. ArXiv preprint, abs/2310.03744, 2023b. URL https://arxiv.org/abs/2310.03744.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. In
 Thirty-seventh Conference on Neural Information Processing Systems, 2023c. URL https:
 //openreview.net/forum?id=w0H2xGH1kw.
- Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee.
 Llava-next: Improved reasoning, ocr, and world knowledge, 2024a. URL https://llava-vl.github.io/blog/2024-01-30-llava-next/.
- Qian Liu, Xiaosen Zheng, Niklas Muennighoff, Guangtao Zeng, Longxu Dou, Tianyu Pang, Jing
 Jiang, and Min Lin. Regmix: Data mixture as regression for language model pre-training. *arXiv preprint arXiv:2407.01492*, 2024b.
- Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang Zhang, Wangbo Zhao, Yike Yuan, Jiaqi
 Wang, Conghui He, Ziwei Liu, et al. Mmbench: Is your multimodal model an all-around player?
 ArXiv preprint, abs/2307.06281, 2023d. URL https://arxiv.org/abs/2307.06281.
- Jiasen Lu, Christopher Clark, Rowan Zellers, Roozbeh Mottaghi, and Aniruddha Kembhavi. Unifiedio: A unified model for vision, language, and multi-modal tasks. In *The Eleventh International Conference on Learning Representations*, 2022a.
- Jiasen Lu, Christopher Clark, Sangho Lee, Zichen Zhang, Savya Khosla, Ryan Marten, Derek Hoiem, and Aniruddha Kembhavi. Unified-io 2: Scaling autoregressive multimodal models with vision language audio and action. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 26439–26455, 2024a.
- Pan Lu, Swaroop Mishra, Tanglin Xia, Liang Qiu, Kai-Wei Chang, Song-Chun Zhu, Oyvind Tafjord, Peter Clark, and Ashwin Kalyan. Learn to explain: Multimodal reasoning via thought chains for science question answering. *Advances in Neural Information Processing Systems*, 35:2507–2521, 2022b.
- Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng, Kai-Wei Chang, Michel Galley, and Jianfeng Gao. Mathvista: Evaluating mathematical reasoning of foundation models in visual contexts. In *The Twelfth International Conference on Learning Representations*, 2024b. URL https://openreview.net/forum?id=KUNzEQMWU7.
- Bryan McCann, Nitish Shirish Keskar, Caiming Xiong, and Richard Socher. The natural language
 decathlon: Multitask learning as question answering, 2019. URL https://openreview.
 net/forum?id=BllfHhR9tm.
- Brandon McKinzie, Zhe Gan, Jean-Philippe Fauconnier, Sam Dodge, Bowen Zhang, Philipp Dufter, Dhruti Shah, Xianzhi Du, Futang Peng, Floris Weers, et al. Mm1: Methods, analysis & insights from multimodal llm pre-training. *ArXiv preprint*, abs/2403.09611, 2024. URL https://arxiv.org/abs/2403.09611.
- Sachit Menon and Carl Vondrick. Visual classification via description from large language models.
 In *The Eleventh International Conference on Learning Representations*, 2023. URL https: //openreview.net/forum?id=jlAjNL8z5cs.
- Christian M Meyer and Iryna Gurevych. Wiktionary: A new rival for expert-built lexicons? Exploring
 the possibilities of collaborative lexicography. na, 2012.
- George A. Miller. WordNet: A lexical database for English. In Speech and Natural Language:
 Proceedings of a Workshop Held at Harriman, New York, February 23-26, 1992, 1992. URL https://aclanthology.org/H92-1116.

- Chancharik Mitra, Brandon Huang, Trevor Darrell, and Roei Herzig. Compositional chain-ofthought prompting for large multimodal models. *ArXiv preprint*, abs/2311.17076, 2023. URL https://arxiv.org/abs/2311.17076.
- Niklas Muennighoff, Alexander Rush, Boaz Barak, Teven Le Scao, Nouamane Tazi, Aleksandra
 Piktus, Sampo Pyysalo, Thomas Wolf, and Colin A Raffel. Scaling data-constrained language
 models. Advances in Neural Information Processing Systems, 36, 2024.
- Vicente Ordonez, Girish Kulkarni, and Tamara L. Berg. Im2text: Describing images using 1 million captioned photographs. In John Shawe-Taylor, Richard S. Zemel, Peter L. Bartlett, Fernando C. N. Pereira, and Kilian Q. Weinberger (eds.), Advances in Neural Information Processing Systems 24: 25th Annual Conference on Neural Information Processing Systems 2011. Proceedings of a meeting held 12-14 December 2011, Granada, Spain, pp. 1143–1151, 2011. URL https://proceedings.neurips.cc/paper/2011/hash/ 5dd9db5e033da9c6fb5ba83c7a7ebea9-Abstract.html.
- Xichen Pan, Li Dong, Shaohan Huang, Zhiliang Peng, Wenhu Chen, and Furu Wei. Kosmos-g: Generating images in context with multimodal large language models. *ArXiv preprint*, abs/2310.02992, 2023. URL https://arxiv.org/abs/2310.02992.
- Zhiliang Peng, Li Dong, Hangbo Bao, Qixiang Ye, and Furu Wei. Beit v2: Masked image modeling with vector-quantized visual tokenizers. *ArXiv preprint*, abs/2208.06366, 2022. URL https://arxiv.org/abs/2208.06366.
- Zhiliang Peng, Wenhui Wang, Li Dong, Yaru Hao, Shaohan Huang, Shuming Ma, and Furu
 Wei. Kosmos-2: Grounding multimodal large language models to the world. *ArXiv preprint*, abs/2306.14824, 2023. URL https://arxiv.org/abs/2306.14824.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*, pp. 8748–8763. PMLR, 2021.
- Hanoona Rasheed, Muhammad Maaz, Sahal Shaji, Abdelrahman Shaker, Salman Khan, Hisham Cholakkal, Rao M Anwer, Erix Xing, Ming-Hsuan Yang, and Fahad S Khan. Glamm: Pixel grounding large multimodal model. *ArXiv preprint*, abs/2311.03356, 2023. URL https:// arxiv.org/abs/2311.03356.
- Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF confer- ence on computer vision and pattern recognition*, pp. 10684–10695, 2022.
- Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman. Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. In *Proceed-ings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 22500–22510, 2023.
- Christoph Schuhmann, Robert Kaczmarczyk, Aran Komatsuzaki, Aarush Katta, Richard Vencu, Romain Beaumont, Jenia Jitsev, Theo Coombes, and Clayton Mullis. Laion-400m: Open dataset of clip-filtered 400 million image-text pairs. In *NeurIPS Workshop Datacentric AI*, number FZJ-2022-00923. Jülich Supercomputing Center, 2021.
- Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi
 Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. Laion-5b: An
 open large-scale dataset for training next generation image-text models. *Advances in Neural Information Processing Systems*, 35:25278–25294, 2022a.
- Christoph Schuhmann, Andreas Köpf, Richard Vencu, Theo Coombes, and Romain Beaumont.
 Laion-coco. https://huggingface.co/datasets/laion/laion-coco, 2022b.
- Christoph Schuhmann, Andreas Köpf, Richard Vencu, Theo Coombes, and Romain Beaumont. Laion coco: 600m synthetic captions from laion2b-en, 2023. URL https://laion.ai/blog/laion-coco/.

864 865 866	Shuai Shao, Zeming Li, Tianyuan Zhang, Chao Peng, Gang Yu, Xiangyu Zhang, Jing Li, and Jian Sun. Objects365: A large-scale, high-quality dataset for object detection. In 2019 IEEE/CVF International Control on Commuter Vision ICCV 2010. Security Kenner (South) October 27
867	November 2, 2019, pp. 8429–8438. IEEE, 2019. doi: 10.1109/ICCV.2019.00852. URL https:
868	//doi.org/10.1109/ICCV.2019.00852.
869	
870	Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. Conceptual captions: A cleaned,
871	Annual Masting of the Association for Computational Linguistics (Volume 1: Long Papers)
872	nn 2556–2565 Melbourne Australia 2018 Association for Computational Linguistics doi:
873	10.18653/v1/P18-1238. URL https://aclanthology.org/P18-1238.
874	
875	Sheng Shen, Chunyuan Li, Xiaowei Hu, Yujia Xie, Jianwei Yang, Pengchuan Zhang, Zhe Gan, Lijuan
876 877	Wang, Lu Yuan, Ce Liu, et al. K-lite: Learning transferable visual models with external knowledge. <i>Advances in Neural Information Processing Systems</i> , 35:15558–15573, 2022.
878	Bagesbree Shevade, Hari Sundaram, and Min Yen-Kan. A collaborative annotation framework. In
879 880	2005 IEEE International Conference on Multimedia and Expo, pp. 1346–1349. IEEE, 2005.
881	Aleksandar Shtedritski, Christian Rupprecht, and Andrea Vedaldi. What does clip know about a
882	red circle? visual prompt engineering for vlms. In Proceedings of the IEEE/CVF International
883	Conference on Computer Vision, pp. 11987–11997, 2023.
884	Jascha Sohl-Dickstein Fric A Weiss Niru Maheswaranathan and Surva Ganguli Deen unsupervised
885	learning using nonequilibrium thermodynamics. In Francis R. Bach and David M. Blei (eds.).
886	Proceedings of the 32nd International Conference on Machine Learning, ICML 2015, Lille, France,
887	6-11 July 2015, volume 37 of JMLR Workshop and Conference Proceedings, pp. 2256–2265.
888	JMLR.org, 2015. URL http://proceedings.mlr.press/v37/sohl-dickstein15.
800	ntml.
891	Petru Soviany, Radu Tudor Ionescu, Paolo Rota, and Nicu Sebe. Curriculum learning: A survey.
892	International Journal of Computer Vision, 130(6):1526–1565, 2022.
893 894	ChatGPT Steven Tey. Sharegpt. https://sharegpt.com/, 2023.
895 896	Keqiang Sun, Junting Pan, Yuying Ge, Hao Li, Haodong Duan, Xiaoshi Wu, Renrui Zhang, Aojun Zhou, Zipeng Qin, Yi Wang, et al. Journeydb: A benchmark for generative image understanding. <i>Advances in Neural Information Processing Systems</i> , 36, 2024.
898	
899	Quan Sun, Yufeng Cui, Xiaosong Zhang, Fan Zhang, Qiying Yu, Zhengxiong Luo, Yueze Wang,
900	Yongming Rao, Jingjing Liu, Tiejun Huang, et al. Generative multimodal models are in-context learners. ArXiv preprint abs/2312 13286 2023a URL https://arxiv.org/abs/2312
901	13286.
902	
903	Quan Sun, Yuxin Fang, Ledell Wu, Xinlong Wang, and Yue Cao. Eva-clip: Improved training
904	techniques for clip at scale. ArXiv preprint, abs/2303.15389, 2023b. URL https://arxiv.
905	019/abs/2303.13389.
906	Quan Sun, Qiying Yu, Yufeng Cui, Fan Zhang, Xiaosong Zhang, Yueze Wang, Hongcheng Gao,
907	Jingjing Liu, Tiejun Huang, and Xinlong Wang. Generative pretraining in multimodality. ArXiv
908	preprint, abs/2307.05222, 2023c. URL https://arxiv.org/abs/2307.05222.
909	Rohan Taori Ishaan Gulrajani Tianyi Zhang Yann Dubois Xuechen Li Carlos Guestrin Percy
910	Liang, and Tatsunori B. Hashimoto. Stanford alpaca: An instruction-following llama model
311 012	https://github.com/tatsu-lab/stanford alpaca, 2023.
JIZ	
913	Chameleon Team. Chameleon: Mixed-modal early-fusion foundation models. <i>arXiv preprint</i>
915	arxiv:2405.09818, 2024.
916	Changyao Tian, Xizhou Zhu, Yuwen Xiong, Weiyun Wang, Zhe Chen. Wenhai Wang, Yuntao Chen.
917	Lewei Lu, Tong Lu, Jie Zhou, Hongsheng Li, Yu Qiao, and Jifeng Dai. Mm-interleaved: Interleaved image-text generative modeling via multi-modal feature synchronizer, 2024.

918 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée 919 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and 920 efficient foundation language models. ArXiv preprint, abs/2302.13971, 2023a. URL https: 921 //arxiv.org/abs/2302.13971. 922 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay 923 Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation 924 and fine-tuned chat models. ArXiv preprint, abs/2307.09288, 2023b. URL https://arxiv. 925 org/abs/2307.09288. 926 Aäron van den Oord, Oriol Vinyals, and Koray Kavukcuoglu. Neural discrete representation learning. 927 In Isabelle Guyon, Ulrike von Luxburg, Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. 928 Vishwanathan, and Roman Garnett (eds.), Advances in Neural Information Processing Systems 929 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, 930 Long Beach, CA, USA, pp. 6306-6315, 2017. URL https://proceedings.neurips.cc/ 931 paper/2017/hash/7a98af17e63a0ac09ce2e96d03992fbc-Abstract.html. 932 Jiaqi Wang, Pan Zhang, Tao Chu, Yuhang Cao, Yujie Zhou, Tong Wu, Bin Wang, Conghui He, and 933 Dahua Lin. V3det: Vast vocabulary visual detection dataset. In Proceedings of the IEEE/CVF 934 International Conference on Computer Vision, pp. 19844–19854, 2023a. 935 936 Weiyun Wang, Min Shi, Qingyun Li, Wenhai Wang, Zhenhang Huang, Linjie Xing, Zhe Chen, Hao 937 Li, Xizhou Zhu, Zhiguo Cao, et al. The all-seeing project: Towards panoptic visual recognition 938 and understanding of the open world. In The Twelfth International Conference on Learning 939 Representations, 2023b. 940 Haoning Wu, Zicheng Zhang, Erli Zhang, Chaofeng Chen, Liang Liao, Annan Wang, Chunyi Li, 941 Wenxiu Sun, Qiong Yan, Guangtao Zhai, et al. Q-bench: A benchmark for general-purpose 942 foundation models on low-level vision. arXiv preprint arXiv:2309.14181, 2023a. 943 Yifan Wu, Pengchuan Zhang, Wenhan Xiong, Barlas Oguz, James C Gee, and Yixin Nie. The role 944 of chain-of-thought in complex vision-language reasoning task. ArXiv preprint, abs/2311.09193, 945 2023b. URL https://arxiv.org/abs/2311.09193. 946 947 De Xie, Cheng Deng, Chao Li, Xianglong Liu, and Dacheng Tao. Multi-task consistency-preserving 948 adversarial hashing for cross-modal retrieval. IEEE Transactions on Image Processing, 29:3626-949 3637, 2020. 950 Xing Xu, Tan Wang, Yang Yang, Lin Zuo, Fumin Shen, and Heng Tao Shen. Cross-modal attention 951 with semantic consistence for image-text matching. IEEE transactions on neural networks and 952 learning systems, 31(12):5412-5425, 2020. 953 954 Jianwei Yang, Hao Zhang, Feng Li, Xueyan Zou, Chunyuan Li, and Jianfeng Gao. Set-of-mark prompting unleashes extraordinary visual grounding in gpt-4v. ArXiv preprint, abs/2310.11441, 955 2023. URL https://arxiv.org/abs/2310.11441. 956 957 Lewei Yao, Jianhua Han, Youpeng Wen, Xiaodan Liang, Dan Xu, Wei Zhang, Zhenguo Li, Chunjing 958 Xu, and Hang Xu. Detclip: Dictionary-enriched visual-concept paralleled pre-training for open-959 world detection. Advances in Neural Information Processing Systems, 35:9125–9138, 2022. 960 Yuan Yao, Ao Zhang, Zhengyan Zhang, Zhiyuan Liu, Tat-Seng Chua, and Maosong Sun. Cpt: 961 Colorful prompt tuning for pre-trained vision-language models. AI Open, 5:30–38, 2024. 962 963 Shukang Yin, Chaoyou Fu, Sirui Zhao, Ke Li, Xing Sun, Tong Xu, and Enhong Chen. A survey 964 on multimodal large language models. ArXiv preprint, abs/2306.13549, 2023. URL https: //arxiv.org/abs/2306.13549. 965 966 Lili Yu, Bowen Shi, Ramakanth Pasunuru, Benjamin Muller, Olga Golovneva, Tianlu Wang, Arun 967 Babu, Binh Tang, Brian Karrer, Shelly Sheynin, et al. Scaling autoregressive multi-modal models: 968 Pretraining and instruction tuning. arXiv preprint arXiv:2309.02591, 2(3), 2023a. 969 Qiying Yu, Quan Sun, Xiaosong Zhang, Yufeng Cui, Fan Zhang, Xinlong Wang, and Jingjing Liu. 970 Capsfusion: Rethinking image-text data at scale. ArXiv preprint, abs/2310.20550, 2023b. URL 971 https://arxiv.org/abs/2310.20550.

972 973 974	Tianyu Yu, Jinyi Hu, Yuan Yao, Haoye Zhang, Yue Zhao, Chongyi Wang, Shan Wang, Yinxv Pan, Jiao Xue, Dahai Li, et al. Reformulating vision–language foundation models and datasets towards universal multimodal assistants. <i>ArXiv preprint</i> , abs/2310.00653, 2023c. URL https:
975	//arxiv.org/abs/2310.00653.
976	Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens,
977	Dongfu Jiang, Weiming Ren, Yuxuan Sun, et al. Mmmu: A massive multi-discipline multimodal
978	understanding and reasoning benchmark for expert agi. ArXiv preprint, abs/2311.16502, 2023.
979	URL https://arxiv.org/abs/2311.16502.
900	Van Zang, Vinsong Zhang, and Hang Li. Multi grained vision language pre-training: Aligning texts
082	with visual concepts ArXiv preprint abs/2111 08276 2021 URL https://arxiv.org/
983	abs/2111.08276.
984	
985	Yan Zeng, Xinsong Zhang, Hang Li, Jiawei Wang, Jipeng Zhang, and Wangchunshu Zhou. X 2-vlm:
986	All-in-one pre-trained model for vision-language tasks. <i>IEEE Transactions on Pattern Analysis</i>
987	ana Machine Intelligence, 2023.
988	Kai Zhang, Lingbo Mo, Wenhu Chen, Huan Sun, and Yu Su. Magicbrush: A manually annotated
989	dataset for instruction-guided image editing. Advances in Neural Information Processing Systems,
990	36, 2024.
991	Penochuan Zhang Xiujun Li Xiaowei Hu Jianwei Yang Lei Zhang Lijuan Wang Veiin Choi and
992	Jianfeng Gao, Vinyl: Revisiting visual representations in vision–language models. In <i>Proceedings</i>
993	of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 5579–5588, 2021.
994	
995	Yanzhe Zhang, Ruiyi Zhang, Jiuxiang Gu, Yufan Zhou, Nedim Lipka, Diyi Yang, and Tong Sun.
996	Liavar: Enhanced visual instruction tuning for text-rich image understanding. ArXiv preprint,
997	aus/2500.1/10/, 2025. UKL https://atxiv.org/abs/2500.1/10/.
998	Jinguo Zhu, Xiaohan Ding, Yixiao Ge, Yuying Ge, Sijie Zhao, Hengshuang Zhao, Xiaohua Wang, and
999	Ying Shan. VI-gpt: A generative pre-trained transformer for vision and language understanding
1000	and generation. ArXiv preprint, abs/2312.09251, 2023a. URL https://arxiv.org/abs/
1001	2312.09251.
1002	Wanrong Zhu, Jack Hessel, Anas Awadalla, Samir Yitzhak Gadre, Jesse Dodge, Alex Fang, Youngjae
1003	Yu, Ludwig Schmidt, William Yang Wang, and Yejin Choi. Multimodal c4: An open, billion-scale
1004	corpus of images interleaved with text, 2023b.
1005	
1007	
1007	
1009	
1010	
1011	
1012	
1013	
1014	
1015	
1016	
1017	
1018	
1019	
1020	
1021	
1022	
1023	
1024	
1025	

1026		
1027		
1028	Thi	s appendix is organized as follows.
1029		
1030		• In Section A, we discuss the broader impact of our work.
1031		• In Section B, we discuss the limitations of our work and potential future research directions.
1032 1033 1034 1035		• In Section C, we provide more experimental results and analyses, including concept overlap, image generation and editing, self-generated annotations in evaluation, more cases on SEED-Bench, MMGIC directly as SFT data, the impact of the nature of image-text interleaved on VIST, evaluation on Q-Bench and HallusionBench, and text-only performance.
1036 1037 1038		 In Section D, we provide discussions on our framework and existing MLLMs, applicability of MMGIC, annotation synthesis in MMGIC, MMGIC and existing interleaved datasets. In Section E (referred by Section 3.3 & 4) we provide more training and evaluation details of
1039 1040		our experiments.
1041		• In Section F (referred by Section 2), we provide complete details of MMGIC dataset.
1042		• In Section G (referred by Section 4.2), we provide more details of IC dataset.
1043 1044 1045		• In Section H (referred by Section 3.3 & 4.3), we provide more details of SFT data we used in our experiments.
1046 1047 1048 1049 1050	A	BROADER IMPACT

In this paper, we introduce MMGIC, a new dataset with multimodal multi-grained concept annotations 1051 for MLLMs. We first collect, pre-process and complement four public large-scale human-annotated 1052 object detection datasets. With a well-designed structured template, we transform these datasets 1053 into our MMGIC dataset to address the lack of such datasets in the MLLM field and support our 1054 exploration. Under an autoregressive discrete framework, we explore the potential of MMGIC for 1055 MLLMs by evaluations and analyses on various downstream vision-language benchmarks. We 1056 demonstrate that MMGIC can provide multi-grained concept annotations to help MLLMs better 1057 locate and learn concepts, thus facilitating vision-language alignment across multiple granularities simultaneously. We will open-source the models and the code of the complete pipeline from data 1058 processing, model training to evaluation for facilitating more reproducible research. We hope that our 1059 work can inspire more research on multi-grained concept annotations for MLLMs. 1060

We do not anticipate a specific negative impact, but, as with any multimodal large language model (MLLM), we recommend to exercise caution. MMGIC dataset is constructed from four public widely-used object detection datasets and will not bring any potential ethical or societal issues. We will follow the usage rules and copyright of original datasets. Newly generated data parts will follow the usage rules and copyrights of the corresponding open-source or closed-source models. We will require researchers who use our code and models to follow the principles of positive AI research to avoid model abuse and negative societal impact.

1068 1069

1070

1072

B LIMITATION AND FUTURE WORK

1071 B.1 AUTOMATIC ANNOTATION SYNTHESIS

Nowadays, both LLMs and MLLMs can continuously improve the performance and generalization ability by scaling up training data. Although our MMGIC dataset, based on several human-annotated object detection datasets and our well-designed data pipeline, can provide high-quality multi-grained concept annotations, it is still limited by the scale of the original datasets. A natural idea is to automatically synthesize a variety of different types of annotations for any image. This is a promising direction for future research, which can scale up to more concepts and continuously improve MLLMs.

1079 Recent works (Peng et al., 2023; Pan et al., 2023; Wang et al., 2023b; Rasheed et al., 2023; Li et al., 2023d; Yang et al., 2023; Li et al., 2024c) have proposed to automatically synthesize various

1080 annotations for large-scale web images with external open/close source models and complicated pipelines. KOSMOS-2 (Peng et al., 2023) extracted noun phrases from image captions and detected 1082 bounding boxes from images to synthesize Grounded image-text pairs (GRIT) to improve the grounding capability of MLLMs. KOSMOS-G (Pan et al., 2023) synthesized image captions and 1084 segment object regions for images to construct compositional image generation instruction tuning data to improve the zero-shot subject-driven image generation capability of MLLMs. All-Seeing (Wang et al., 2023b) and GLaMM (Rasheed et al., 2023) designed complicated pipelines to combine multiple 1086 external models to synthesize various annotations for large-scale web images from SA-1B (Kirillov 1087 et al., 2023), including detailed image captions, object labels, object regions, question-answer pairs 1088 about objects in the image, scene graphs, and so on, to improve the image/region-level captioning 1089 and grounding capability of MLLMs. SoM (Yang et al., 2023) synthesized semantic segmentation 1090 masks for images to improve the image segmentation and grounding capability of MLLMs. 1091

Although these works have shown that automatically synthesized annotations can improve the 1092 performance of MLLMs, cascading multiple external models and complicated pipelines may introduce 1093 noise and bias to the synthesized annotations. Existing works introduced complex data pipelines 1094 and costly manual reviews (especially for pre-training data exceeding 1M) to reduce noise, but 1095 their released data still contains quite a bit of hallucinations, noise, and redundancy. Moreover, these works mainly focus on exploring the potential of automatically synthesized annotations in the multimodal comprehension or generation capability of MLLMs, and have not explored both under a unified framework at the same time. Among them, All-Seeing and GLaMM are similar to 1099 MMGIC proposed in this paper. Both of them synthesize various annotations for images, including 1100 image captions, category labels and object regions provided in MMGIC. They don't provide label 1101 descriptions but provide visual question-answering pairs or scene graphs, or other annotations. However, the final data forms of them are similar to the related works in traditional VLMs, where 1102 each image is paired with isolated different types of annotations. They need to introduce additional 1103 components and loss functions to utilize different annotations, such as object region annotations, to 1104 optimize the model's comprehension capability at different granularities separately. Besides, their 1105 experiments mainly explore the object recognition and segmentation, image/region-level captioning 1106 and grounding capability of MLLMs, which are actually downstream tasks that obviously benefit 1107 from their synthesized/annotated different types of annotations. 1108

In contrast, in this paper, from the perspective of concept annotations, we explore the potential of multi-grained concept annotations for MLLMs for the first time, in both multimodal comprehension and multimodal generation. We mainly focus on general multimodal benchmarks instead of specific benchmarks that fine-grained annotations are good at (*e.g.*, object recognition and grounding benchmarks). We want to evaluate and analyze the potential of multi-grained concept annotations for MLLMs in a more comprehensive, general and fair way, by comparing and collaborating with widely-used coarse-grained image-caption data.

We believe that the automatic synthesis of multi-grained concept annotations for any image is a 1116 promising direction for future research. In addition, incorporating more different types of annotations, 1117 including not only various annotations found in existing works, but also annotations for text-rich 1118 images and table/chart images, would be valuable for future research. We have explored the potential 1119 of multi-grained concept annotations for MLLMs in both multimodal comprehension and generation, 1120 including image captioning, text-to-image generation, visual question answering, multi-choice 1121 benchmarks, and image editing. It is interesting to explore more benchmarks on object recognition 1122 and segmentation, grounding capability and other vision-language tasks in the future. Finally, in 1123 this work, we focus on concrete concepts, especially objects, attributes of objects, and relationships 1124 between objects. It is very interesting to explore more concrete concepts and also abstract concepts, such as emotions, events, and so on, in the future. 1125

- 1126
- 1127
- 1128 1129

B.2 COLLABORATION STRATEGY OF MMGIC AND IMAGE-CAPTION DATA

In Section 4.2, we investigate comparison and collaboration between MMGIC and coarse-grained
 image-caption dataset IC on 4 benchmarks for multimodal comprehension and generation. Simply
 joint training with MMGIC and IC cannot bring performance improvements compared to training
 with MMGIC only, and even leads to obvious performance degradation on COCO, NoCaps and
 VIST. Considering that IC is ~ 15 times larger than MMGIC, we also try to balance the two datasets



Figure 6: Visualization of the original and reconstructed images by the visual tokenizer inherited from
 LaVIT (Jin et al., 2023) used in this paper. Text-rich and chart images cannot be well reconstructed.

by using $2 \sim 4$ times of MMGIC¹. However, repeating MMGIC, *e.g.*, 4MMGIC+IC, introduces only minor performance fluctuations and results in lower average performance.

1154 Therefore, we follow the curriculum learning strategy (McCann et al., 2019) to train them in different 1155 orders and achieves significantly better average performance by training on IC first and then on 1156 MMGIC. This is consistent with recent findings (Hu et al., 2024b; Li et al., 2024a) that training with 1157 high-quality data late in the pre-training phase leads to better performance. Moreover, considering that the noise in IC still cause a slight performance drop on VIST ($\{2, 5\}$ -th rows in Table 2), and 1158 jointly training on MMGIC and IC outperforms IC on both tasks ($\{1,3\}$ -th rows in Table 2). Hence, 1159 we first jointly train on MMGIC and IC to alleviate the effect of noise in IC, and then on MMGIC, 1160 eventually achieving the best average performance (6-th row in Table 2). Similarly, further increasing 1161 MMGIC repetitions, e.g., $3MMGIC+IC \rightarrow MMGIC$, only yields minor performance fluctuations 1162 and leads to noticeable average performance drops. We speculate that repetition of MMGIC leads 1163 to overfitting of MMGIC and insufficient learning of IC, resulting in performance declines and 1164 fluctuation. We believe that the collaboration strategy of MMGIC and image-caption data is an 1165 interesting research question about data mixing (Liu et al., 2024b), data repetition (Muennighoff et al., 1166 2024) and curriculum learning (Soviany et al., 2022). Automatic annotation synthesis for large-scale 1167 image-caption data we discussed in Appendix B.1 is also a promising solution to scale up to more concepts and continuously improve MLLMs without suffering from the problem of data duplication. 1168

1169 1170

1171

1151

B.3 A MLLM FRAMEWORK FOR MULTIMODAL COMPREHENSION AND GENERATION

Parameter-Efficient Training. In this work, we standardize a generative vision–language framework based on existing MLLMs, which consists of a LLM and several visual modules. We follow
previous works (Ge et al., 2023a; Zhu et al., 2023a) to freeze all visual modules and most of the LLM
parameters during training to greatly improve efficiency. Although the performance of the framework
is competitive on various vision–language benchmarks, it is limited by the frozen visual modules and
parameter-efficient fine-tuning of the LLM. In the future, we will explore full-parameter fine-tuning
of the LLM and visual modules to further improve the performance of the framework.

1179

Visual Tokenizer. Unlike traditional visual tokenizers like VQ-VAE (van den Oord et al., 2017) that reconstruct from original image pixels, SEED (Ge et al., 2023a) followed BEIT v2 (Peng et al., 2022) to train the visual tokenizer by reconstructing from visual embeddings with high-level semantics, and found that the latter strategy can better capture high-level semantics instead of low-level image details in the former strategy, which is more effective for multimodal comprehension. As the concurrent work of SEED, LaVIT also found that high-level semantics matters for downstream tasks, and further supported dynamic sequence length varying from the image. Compare SEED and LaVIT, we take the visual tokenizer from LaVIT, which can better preserve image details and semantic information.

1187

¹We follow the advice from Muennighoff et al. (2024) to repeat data no more than 4 times.



Figure 7: Comparison of performance (CIDEr Score) and concept overlap (%) for different training data settings on image captioning datasets COCO and NoCaps.

However, no matter which visual tokenizer is used, the vector quantization process is a trade-off
between the cost/difficulty of learning and the capacity of representing visual information. For
example, as shown in Figure 6, the visual tokenizer of LaVIT we used in this work cannot well
reconstruct the original image by the visual decoder and diffusion model, especially text-rich images
and chart images. Chameleon (Team, 2024) also found such limitations of traditional visual tokenizers
(reconstructed from original image pixels) on heavy OCR-related tasks.

1208 We also conduct experiments on downstream VL benchmarks containing text-rich images and 1209 table/chart images to further validate the limitation of the visual tokenizer. Our MLLM-MMGIC & IC 1210 achieves 32.0 accuracy on MMMU (Yue et al., 2023) validation set, where the human performance, 1211 the best model performance and random performance are 88.6, 59.4, 22.1, respectively; and achieves 26.0 accuracy on MathVista (Lu et al., 2024b) testmini set, where the human performance, the best 1212 model performance and random performance are 60.3, 63.9, 17.9, respectively. In the future, we will 1213 explore more advanced visual tokenizers to better capture high-level semantics and low-level image 1214 details for multimodal comprehension on general, text-rich and table/chart images. 1215

1216 LLM Part. We conduct experiments with LLaMA-2-7B as the LLM part in our framework. 1217 LLaMA-2-7B is widely used in the LLM and MLLM fields due to its good performance and 1218 generalizability, ensuring the reliability, reproducibility and scalability of our experimental results. 1219 Especially, most of the existing SOTA MLLMs (e.g., EMU (Sun et al., 2023c), LaVIT (Jin et al., 1220 2023), SEED-LLaMA (Ge et al., 2023b), DreamLLM (Dong et al., 2024), VL-GPT (Zhu et al., 1221 2023a)) all use the LLaMA series or its variants to initialize the LLM part of their models. Since our 1222 research aims to explore the potential of MMGIC as a new data paradigm for MLLMs, we follow 1223 this common practice to use LLaMA-2-7B in our framework. It will be interesting to explore more 1224 advanced LLMs in the future to further improve the performance of our framework, considering recent findings that the model size scaling of LLM is more effective than image encoder in yielding 1225 improved performance (Li et al., 2024a). 1226

1227

1199

1200 1201

¹²²⁸ C MORE EXPERIMENTAL RESULTS AND ANALYSIS

1230 C.1 CONCEPT OVERLAP 1231

In Section 4.3, we find that MMGIC and IC have their own strengths in depth and breadth of concept representation. To further verify our findings, inspired by K-LITE (Shen et al., 2022), we investigate the concept overlap between different training data settings and downstream datasets in terms of dataset statistics. The concept overlap is computed as the percentage of concepts in a downstream dataset that are covered by the training data. We select two image captioning datasets, COCO and NoCaps, and directly take noun chunks extracted from their ground-truth captions as concepts. We select the $\{0, 1, 2, 6\}$ -th data recipes in Table 2 and also include MMGIC(C) from $\{0$ -th data recipe in Table 1 for better comparison. As shown in Figure 7, we can see that:

- 1239
- Comparing IC-PART and IC: with the number of image–caption pairs increasing from less than 4M to almost 52M, the concept overlap increases significantly (10.2% and 11.2%) and the performance also improves obviously (8.5 and 7.2 points) on both datasets.

Table 5: Zero-shot evaluation on the image editing benchmark DreamBench (Ruiz et al., 2023) after
SFT. DINO: cosine similarity between generated and real images via DINO (Caron et al., 2021). The
best results are **bold** and the second-best are underlined.

	Model	D	reamBenc	h
	Model	CLIP-T	CLIP-I	DINO
SC	OTA MLLMs as upper bound references, no	ot compara	ble	
	KOSMOS-G (Pan et al., 2023)	28.70	84.70	69.40
a	Emu2-Gen-37B (Sun et al., 2023a)	28.70	85.00	76.60
b	SEED-LLaMA-I-8B (Ge et al., 2023b)	27.06	79.27	54.42
Οı	ır comparable baselines			
0	MLLM-IC	27.50	82.92	67.41
1	MLLM-MMGIC	27.69	82.62	67.36
2	MLLM-MMGIC & IC	27.62	83.06	68.30

• Comparing MMGIC(C) and MMGIC: the performance improvement (4.7 and 2.9 points) may not come from improved concept overlap (only 2.5% and 2.4%). This provides side evidence that fine-grained concept annotations in MMGIC can integrate and complement coarse-grained concept annotations in MMGIC(C), to help MLLMs learn concepts deeply and thus improve performance.

- Comparing MMGIC and MMGIC+IC \rightarrow MMGIC: MMGIC can collaborate with IC to achieve higher concept overlap (15.4% and 17.7%). Appropriate curriculum learning strategies can effectively integrate the strengths of both in terms of depth and breadth of concept representation, thus improving performance on both datasets (2.9 and 3.3 points).
- Comparing IC-PART, IC and MMGIC(C): although they are all image–caption pairs and captions are all synthesized by the same pipeline (Appendix F.4), MMGIC(C) can achieve better performance on both datasets with lower concept overlap and the same number of image–caption pairs as IC-PART. We attribute this to the fact that the images in MMGIC(C) have higher quality, *i.e.*, contain more concepts and are more complex than IC-PART and IC collected from the web. Specially, the number of unique noun chunks in the three training data are: 1.7M, 10.9M, and 0.4M, respectively, while the average number of unique noun chunks per image are: 3.0, 3.0, and 8.8, respectively. This is consistent with recent findings (Dai et al., 2024; Li et al., 2024a) that the impact of data quality is greater than data scale.

C.2 ANALYSIS ON IMAGE EDITING

To further investigate how MMGIC can better help MLLMs utilize concepts to align vision and language and improve the capability of generating concepts we conduct both qualitative and quantitative analyses on image editing task.

Although MMGIC does not contain any samples about image editing, with our structured template, each sample of MMGIC can be seen as an image–text interleaved document with an image, multiple cropped regions from the image, textual annotations, and also template instruction text. Such design can help MLLMs better understand each region in the image, as well as the relationships within and between regions, thus better learning and generating concepts, and aligning images and instructions.
We hypothesize that this may help MLLMs learn more diverse image generation abilities.

Quantitative Analysis. To verify this, we include about 0.32M image editing instruction samples from Instructpix2pix (Brooks et al., 2023) and MagicBrush (Zhang et al., 2024) during the SFT stage, and evaluate the performance of our baselines on the image editing benchmark DreamBench (Ruiz et al., 2023), as shown in Table 5. Compared MMGIC with IC, IC performs better on image (subject) fidelity (CLIP-I and DINO), and MMGIC performs better on text fidelity (CLIP-T, i.e., adherence to editing instructions). We attribute this to the fact that the structured template and rich multimodal annotations in MMGIC can help MLLMs better understand and follow editing instructions, while the massive image–caption pairs and higher concept breadth in IC can help MLLMs better generate



Figure 8: Detailed evaluation of MLLM-MMGIC & IC with different evaluation strategies on 8 dimensions of SEED-Bench-IMG.

- and edit images. Furthermore, we also find that the collaboration of MMGIC and IC combines
 the advantages of both, and achieves better average performance on both image and text fidelity.
 Compared with existing SOTA MLLMs, MMGIC & IC not only significantly outperforms SEEDLLaMA (Ge et al., 2023b) (using more pre-training data, SFT data, and full-param fine-tuning), but
 also achieves comparable performance to generation-only MLLMs (KOSMOS-G (Pan et al., 2023))
 and Emu2-Gen-37B (Sun et al., 2023a)) trained with far more image generation and editing data.
- 1323 1324

Qualitative Analysis. We also show some qualitative examples of image editing from MLLMs 1325 pre-trained with MMGIC in Figure 4 (*Right*). The top two examples show that MMGIC can precisely 1326 understand editing instructions and perform appropriate editing, while the bottom example shows that 1327 MMGIC can synthesize image precisely based on the image-text interleaved sequences. Considering 1328 that MMGIC dataset does not contain any samples about above two emergent abilities of image 1329 editing and multimodal in-context image synthesis, we hypothesize that integrating multi-grained 1330 concept annotations into image-text interleaved documents through our structured template can help 1331 MLLMs better understand and locate objects in the image, to more accurately edit the objects in the 1332 image, and synthesize images based on the image-text interleaved sequences. 1333

- 1334
- 1335 1336 1337
- C.3 SELF-GENERATED MULTI-GRAINED CONCEPT ANNOTATIONS IN EVALUATION

Recent works (Sun et al., 2023c; Wu et al., 2023b; Mitra et al., 2023) have proposed to ask MLLMs 1338 to answer questions with the help of image annotations (e.g., captions, or scene graphs) generated by 1339 themselves or external models during evaluation. Inspired by these works, we naturally ask MLLMs 1340 trained with MMGIC to self-generate coarse- and fine-grained concept annotations for images, like 1341 annotations provided in MMGIC. Considering the inference efficiency and avoiding noise, we only 1342 generate category labels as fine-grained concept annotations. Specifically, we explore different evaluation strategies on SEED-Bench-IMG in Fig. 8. After training with MMGIC, fine-grained 1344 concept annotations can help MLLMs generate better coarse-grained image captions to improve 1345 performance, especially on the "Instance Location, Instance Interaction, Visual Reasoning" dimensions. Furthermore, self-generate multi-grained concept annotations can integrate the advantages of both granularities, achieving 1.68 points improvement in overall accuracy compared with the 1347 zero-shot evaluation, especially on the "Instance Identity, Instance Counting, Instance Interaction, 1348 Visual Reasoning" dimensions. This demonstrates that multi-grained concept annotations can help 1349 MLLMs better understand and reason about concepts in the image during evaluation.



Figure 9: Case study of MLLM-MMGIC trained with different-grained concept annotations from MMGIC on SEED-Bench-IMG. CG, FG, and MG denote MLLMs trained with coarse-, fine-, and multi-grained concept annotations from MMGIC, respectively. ✓ denote the ground truth; × denote incorrect prediction(s). The bottom right of the first three examples show the associated label– description pairs from MMGIC.

Table 6: Zero-shot evaluation on multimodal comprehension and generation benchmarks after SFT.
 MMB: MMBench; SQA^I: ScienceQA-IMG; SEED^I: SEED-Bench-IMG. The best results are **bold**.

	Model	Image C COCO	Captioning NoCaps	VQAv2	VQA 2 GQA	VizWiz	POPE	Multi-Ch MME	oice Be MMB	enchma SQA ^I	rk SEED ^I	MS FID (↓)	-COCO-) CLIP-T	30K CLIP-I	FID (↓)	IST CLIP-
0 1	MLLM-IC + MMGIC	108.13 119.03	92.71 105.49	70.28 70.33	56.02 56.6 7	2 52.62 7 52.63	81.14 82.31	1646.71 1656.29	59.54 59.64	65.94 66.04	58.41 59.60	8.11 7.65	30.90 31.26	70.72 71.66	38.19 34.51	65.37 66.61
С	4 O UA	LITAT	IVE AN	IALYS	15 0	Б ТНЕ	Ιмρ	АСТ ОР	7 Dif	FERE	NT-G	RAINI	ED CO	NCEP	Г	
	Ann	ΟΤΑΤΙ	ONS ON	N ML	LMs	5			2						-	
W	e have pr	ovided	the qua	antitat	ive r	esults	of M	LLM-N	MMC	GIC w	vith di	fferen	t-grain	ed cor	ncept a	anno
.ic fu	rther pro	vide co	orrespoi	nding	qual	itative	anal	ysis to	bette	r ana	lyse tl	ne adv	in Sec	es of c	.4. пе oarse-	, fin
an	d multi-g	graineo	d conce	pt ann	otati	ions in	MM	GIC f	or M	LLM	s. As	show	n in Fi	gure 9	, we c	an s
tha	at:			_												
	• Exam	ples 1	$\int and \int 2$	fron	n the	"Insta	ince I	dentity	y" dir	nensi	on: w	hile C	CG can	provi	de a h	olist
	betwe	ption o en "Bc	of the wr	101e if d "He	nage els" i	by the	and "	ge capti Helme	ion, N ts" ar	1LLN nd "H	is trai ats" ii	$\frac{1}{2}$	ith CG Howev	cannc er FG	t disti can n	ngui rovi
	visual	details	s "covei	rs the	foot	and ex	tends	s up the	e leg"	and	releva	nt kno	wledg	e "pro	tectiv	e ge
	by lab	el-des	scription	n pair	s, wł	nich ca	in he	lp ML	LMs	better	r unde	erstan	d these	conc	epts.	Obje
	regior	is can	further	help.	MLL	Ms be	tter l	locate	and l	earn t	identi	conce	pts. T	hereto	ore, M	LLN
	traine		FG OF I	NG Ca	an ca	pture	inese	local (letan	s and	Identi		ncepts	correc	uy.	
	• Exam	ple 3	from th	e "Ins	stance	e Attri	bute"	dimer	ision:	after	locat	ing th	e "chai	ndelie	r" and	"wa
	in the	image	e, MLL	Ms ti 121 de	aine	d with	1 FG	or MC	j can labor	reco	gnize	the s	hared	attribi	ute "C	Irna
	conec	uy by		iai uc	tans	ueco	lateu	with	labol	ale ut	stalls a	anu pa	uterns	•		
	• Exam	ple 4	from t	he "Ir	istan	ce Lo	catio	n" dim	ensic	on: M	ILLM	s trai	ned wi	th FG	or M	IG c
	and th	e spati	ial relati	ionshi	iove	in the	ther	ge, and n	i then	corre	ectry a	inarys	e its st	irroun	ang	objec
		1 🗔	c			т.		···		.1	• .		1 /	.1	. 1 1	1.
	• Exam	requir	from in es the r	e ins nodel	to co	orrectl	v unc	i uime lerstan	d the	1: the local	detai	lcuon ls of f	betwee	en the	s and	ina i spat
	relatio	nships	The ir	nage s	show	s that	the ta	ble is p	blaced	1 arou	ind the	e chai	rs, rath	er that	n com	plete
	under	the tal	ole. Suc	h deta	ails c	an be	corre	ctly ca	pture	d by]	MLLN	Ms tra	ined w	ith FC	G or M	ĪG.
	• Exam	ple 6	from th	e "Sc	ene U	Jnders	tandi	ng" di	nensi	ion: w	vhile t	he lad	lv in th	e imas	ge is h	oldi
	an "ui	nbrella	a", her s	surrou	Indin	igs sho	w that	at the s	ky is	cove	red w	ith clo	ouds, b	ut it is	s not r	ainir
	which	implie	es that t	he we	eathe	r is "o	verca	st". M	LLM	s trai	ned w	ith FO	G seem	to be	too fo	ocus
	on the	t the c	orrect a	o the c	overa r by l	II scer better	e of t	ne ima standi	ige, w	e gloł	MLLI pal co	vis tra	ined w	nage	JOTN	IG C
					l Uy l		unuer	standi	ng un			Intext		mage		
	• Exam	ple 7	from th	ie "So	cene	Under	stand	ling" d	imen	sion:	simil	arly, l	MLLN	ls train	ned w	ith F
	is stan	ding n	romine	ntlv ir	the	scene.	The	lightin	g his	nage.	ion. a	nd his	stance	n uie all di	aw aff	tenti
	to him	as the	e main s	subjec	t.				6,	P	,					
	• Evam	nle	from th	o "Vi	1	Peaco	nina"	dimer	sion.	the i	maga	show	s a foo	thall n	natch	wha
	a play	er in th	he "whi	te uni	form	is ce	lebra	ting, a	nd a r	olavei	r in th	e "vel	low ur	iform	is lv	ing
	the gr	ound.	The au	dience	e in t	he bac	kgro	und is	cheer	ring f	or the	playe	er in th	e "wh	ite un	ifor
	and ta	king p	ictures	of hir	n. Si	nce th	e ima	ige cor	ntains	rich	objec	ts and	detail	s, ML	LMs t	rain
	with F	G have	e difficu	ally co	nnrn hore i	ning th	ie ma	in focu	$\frac{15}{7}$ of 1	ine in	age, v	while		is trai	ned w	ith (MC
	with c	our stru	actured	temp	late c	can he	lp M	LLMs	comb	bine t	he ad	vantag	ges of	CG an	id FG.	, ma
	differe	ent-gra	ined co	ncept	anno	otation	s con	pleme	nt ead	ch oth	ner, an	d bett	er und	erstan	d and	reas
	about	the im	age froi	m bot	h glo	bal sco	ene ai	nd loca	l deta	ails. F	Hence,	MLL	.Ms tra	ined v	with M	1G c
	correc	tly ide	ntify th	e "wh	ite"	player	as th	e mair	i focu	is of t	he im	age.				

Table 7: Zero-shot evaluation on multimodal generation benchmarks. For each block, the best result are **bold**.

	Training Data	MS	S-COCO-3	VIST		
		FID (↓)	CLIP-T	CLIP-I	FID (\downarrow)	CLIP-I
0	MMGIC(C)	7.20	30.81	71.62	67.61	62.22
1	MMGIC	7.36	31.57	72.24	35.33	66.10
2	MMC4-Pairs (Zhu et al., 2023b)	16.55	29.32	69.13	101.68	60.43
3	MMC4 (Zhu et al., 2023b)	17.81	28.68	68.48	40.31	66.00

1458

1471 C.5 MMGIC AS INSTRUCTION FINE-TUNING DATA

1472 If we regard the template text of our structured template as multiple instructions and multi-grained 1473 concept annotations as responses to the instructions, then MMGIC can be regarded as image-text 1474 interleaved instruction fine-tuning data. It can be directly used in the SFT stage of existing MLLMs 1475 pre-trained with image–caption data only. As shown in Table 6, the 0-th row corresponds to the 1476 performance of MLLM-IC, an MLLM that pre-trained with only image-caption data, IC, and fine-1477 tuned with our default SFT data in Section 3.3 (1.21M samples from public instruction datasets and 1M samples from an aesthetic image-caption dataset). To explore the potential of MMGIC directly 1478 as instruction fine-tuning data, we add 1M samples from MMGIC to the SFT stage of MLLM-IC. 1479 MLLM-IC+MMGIC in the 1-th row achieves better performance on all benchmarks, especially 1480 on the benchmarks that deeply inspect common concrete concepts, e.g., COCO, NoCaps, POPE, 1481 SEED-Bench, MS-COCO-30K, and VIST. This demonstrates MMGIC can be used not only as 1482 pre-training data, but also as instruction fine-tuning data to help MLLMs learn concepts better, which 1483 further ensure the generality of MMG₁C across different MLLM frameworks.

1484 1485

1486 C.6 IMPACT OF THE NATURE OF IMAGE–TEXT INTERLEAVED ON VIST

1487 Unlike MS-COCO-30K (Chen et al., 2015) that generates a new image based on a single caption, 1488 VIST (Huang et al., 2016) (Visual Storytelling) is a multimodal generation benchmark that requires 1489 MLLMs to generate a new image based on interleaved image-caption context within the same story. 1490 Compared with image-caption data, MMGIC has the nature of image-text interleaved, which may help MLLMs better understand the image-text context in VIST. Hence, in Table 7, we compare the 1491 performance of MMGIC(C) and MMGIC ($\{0,3\}$ -th data recipes in Table 1) on both MS-COCO-30K 1492 for simple text-to-image generation and VIST for in-context image synthesis. In addition, we also 1493 pre-train MLLMs with MMC4² (Zhu et al., 2023b) and MMC4-Pairs (split each MMC4 document 1494 into multiple individual image-text pairs) to ablate the impact of multi-grained concept annotations. 1495

1496 Comparing the top two rows in Table 7, after introducing the nature of image-text interleaved into 1497 our MMGIC, we find that most metrics increase, especially on VIST. However, for MMC4-Pairs and MMC4 in the bottom two rows, in-context image synthesis performance remains increases 1498 but text-to-image generation performance decreases. The results demonstrate that the nature of 1499 image-text interleaved can greatly benefit in-context image synthesis but may not be beneficial for 1500 simple text-to-image generation. We hypothesize that the multi-grained concept annotations in our 1501 MMG1C can work with the nature of image-text interleaved to help MLLMs better learn concepts 1502 and utilize concepts to align vision and language, thus improving performance on both benchmarks. 1503

1504 1505

C.7 EVALUATION ON Q-BENCH AND HALLUSIONBENCH

¹⁵⁰⁶ In this paper, we mainly focus on benchmarks that inspect the comprehension and generation capabilities of MLLMs on concrete concepts. Here, we further evaluate on two new benchmarks,

1508 1509

 ²We use a subset of MMC4, *i.e.*, MMC4-core-fewer-faces, which has a moderate sample size of 22.4M images and 5.5M image-text interleaved training samples, which is similar to the sample size of our MMGIC.
 Following the processing pipeline of MM-Interleaved (Tian et al., 2024), we process MMC4, and then convert each data sample into a discrete token sequence similar to our MMGIC.

	Model	Q-Bench	HallusionBench
a	LLaVA-1.5-7B (Liu et al., 2023b)	58.70	27.60
b	SEED-LLaMA-I-8B (Ge et al., 2023b)	47.22	32.26
0	MLLM-IC	56.66	32.58
1	MLLM-MMGIC	57.59	30.37
2	MLLM-MMGIC & IC	58.26	33.39

Table 8: Zero-shot evaluation on Q-Bench and HallusionBench after SFT.

Table 9: Evaluation on the MMLU (Hendrycks et al., 2021) benchmark with standard 5-shot setting and Accuracy as the metric.

	Training Setting	MMLU
0	LLaMA-2-7B-base	46.22
1	+ MMGIC	29.98
2	+ MMGIC + SFT	46.50

Q-Bench (Wu et al., 2023a) and HallusionBench (Liu et al., 2023a). They focus on low-level vision 1532 abilities (clearness, brightness, etc.), visual illusion and language hallucination (especially abstract 1533 concepts), respectively. As shown in Table 8, we can see that for the low level vision abilities in 1534 Q-Bench, although MMGIC mainly focus on objects, attributes and relations in images, it still helps 1535 MLLM achieve better performance than IC. We attribute this to the fact that concepts like clearness 1536 and brightness are concrete concepts that can also benefit from multi-grained concept annotations (e.g., 1537 captions, label descriptions) in MMGIC. As for HallusionBench, IC achieves significantly better 1538 performance than MMGIC since it contains more concepts, especially abstract concepts. Most im-1539 portantly, the combination of MMGIC and IC can achieve the best performance on both benchmarks, 1540 which demonstrates that their combination can effectively integrate the strengths of both in terms 1541 of depth and breadth of concept representation, thus improving performance on both benchmarks.

1542

1512

1525 1526 1527

1529

1531

1543 C.8 TEXT-ONLY PERFORMANCE

We follow common practice (Sun et al., 2023c; Ge et al., 2023b; Jin et al., 2023) to evaluate the zero-1545 shot performance of MLLMs on multimodal comprehension and generation benchmarks. Notable, 1546 some recent works also explored the text-only performance of MLLMs. For example, LaVIT (Jin 1547 et al., 2023) and MM1 (McKinzie et al., 2024) found that training MLLMs only with multimodal 1548 data could result in a significant degradation of the text-only performance. To address this issue, 1549 they mixed a substantial amount of text-only data, e.g., 66% text-only and 33% multimodal data, 1550 thus preserving the original text capability of the model. Besides, VILA (Lin et al., 2023) observed 1551 that while multimodal pre-training will hurt text-only performance, by simply incorporating some 1552 high-quality text-only instruction data during the SFT stage, the model can achieve similar text-only 1553 performance to the initial LLM.

To improve efficiency, we follow the approach of VILA to incorporate text-only instruction data, *i.e.*, ShareGPT(Steven Tey, 2023) and Alpaca (Taori et al., 2023), into our default SFT data in Section 3.3. We also conduct experiments on the de-facto text-only comprehensive multi-choice benchmark, MMLU (Hendrycks et al., 2021), to quickly evaluate the text-only performance of our models. As shown in Table C.8, we can see that there is indeed a significant performance drop after pre-training with MMGIC. However, after SFT, the MMLU performance of the model recovers well and is similar to the original LLM.

1561

562 C.9 ATTENTION MAP ANALYSIS

1563

In Section 4.4, Figure 5 and 9, we provide qualitative analysis of the impact of different-grained concept annotations in MMGIC on the performance of MLLMs. To further understand the behavior of MLLMs trained with different-grained concept annotations, we visualize the attention maps of the



Figure 10: Attention map of cases on SEED-Bench-IMG from MLLM-MMGIC trained with different-grained concept annotations in MMGIC. CG, FG, and MG denote MLLMs trained with coarse-, fine-, and multi-grained concept annotations from MMGIC, respectively. Cases are the same as in Figure 5. denote the ground truth; × denote incorrect prediction(s). The bottom of each example consists of an original image and two or three attention maps that shows the attention from the predicted answer token to the image tokens. In each attention map, the brighter the color, the higher the attention weight. The image patch with only gray shading means that it is not selected by the visual tokenizer (inherited from LaVIT (Jin et al., 2023)).

1620 cases of Figure 5 on SEED-Bench-IMG, which shows the attention from the predicted answer token
 to the image tokens. In Figure 10, we can see that:

- Example 1 "Instance Identity": Compared with CG (only focusing on the heel), FG can better understand the differences between different types of shoes in the options, to allocate attention more accurately to the image tokens related to shoes (not only the heel but also the toe and the upper part of the shoe), achieving the correct answer by capture the **local details** of the shoes.
- Example 2 "Scene Understanding": Compared with FG (only focusing on the umbrella), CG can better understand the overall scene of the image, to allocate attention more accurately to the image tokens related to the sky and the clouds, achieving the correct answer by better understanding the global context.
- Example 3 "Visual Reasoning": CG seems to be too focused on the player in the yellow uniform, while FG seems to be distracted by the audience, the player in the white uniform and the player the yellow uniform. However, MG can better combine the advantages of CG and FG, to allocate attention more accurately to the image tokens related to the player in the white uniform and the surrounding audience, achieving the correct answer (since the audience is cheering and taking pictures of the player in the white uniform) by better understanding and reasoning about the image from both global context and local details.
- 1637

1623

1624

1625

¹⁶³⁸ D DISCUSSION

1640 D.1 DIFFERENCE BETWEEN OUR FRAMEWORK WITH EXISTING MLLM FRAMEWORKS

1642 This paper mainly focuses on MLLMs that are multimodal generalists capable of both multimodal comprehension and generation. Existing MLLMs (Sun et al., 2023c; Jin et al., 2023; Ge et al., 2023b; 1643 Zhu et al., 2023a; Sun et al., 2023a; Yu et al., 2023a; Lu et al., 2022a; 2024a) are typically trained with 1644 an autoregressive objective, *i.e.*, generate the prediction of the next token in a multimodal sequence 1645 containing discrete textual tokens, and discrete visual tokens or continuous visual embeddings. 1646 MLLMs process images in different ways. For example, EMU (Sun et al., 2023c;a) and VL-1647 GPT (Zhu et al., 2023a) transform images into continuous visual embeddings by visual encoders, 1648 and then adopt a separate regression head to predict visual embeddings. In contrast, CM3Leon (Yu 1649 et al., 2023a), Unified-IO (Lu et al., 2022a; 2024a) and SEED-LLaMA (Ge et al., 2023b) tokenize 1650 images into discrete visual tokens by visual tokenizers, and then predict visual tokens by the extended 1651 vocabulary with the unified language model head. Technically, different training datasets, training 1652 settings (*e.g.*, full- or partial-param), framework, evaluation settings (*e.g.*, image resolution), etc., lead to non-comparable and unfair comparisons of our baselines with existing MLLMs. Their computation (more training data and fully fine-tune the LLM params) and data resources are extremely expensive 1654 and large (well over 10 times that of our work). Our autoregressive discrete framework for MLLMs 1655 is based on SEED-LLaMA and LaVIT, which consists of several visual modules and a LLM with 1656 an extended vision-language vocabulary, as shown in Figure 2. Next, we describe the differences 1657 between our framework and these two MLLMs. 1658

1659 • Visual Modules: our visual modules are inherited from LaVIT (Jin et al., 2023), which consists of a visual encoder, a visual tokenizer, a visual decoder and a diffusion model. The differences between the two MLLMs mainly lie in the design of the latter three modules. The visual tokenizer of SEED-LLaMA uses the Casual Q-Former from BLIP-2 (Li et al., 2023b) to compress visual 1662 representations into a fixed-length sequence of visual embeddings, while LaVIT is more flexible 1663 and designs a dynamic visual tokenizer with token selector and token merger, which can obtain 1664 a more flexible and effective sequence of visual embeddings. For the discrete visual token sequence predicted by the LLM, SEED-LLaMA uses an MLP as the visual decoder to compress and reconstruct it into the image embedding (1 token) and directly generate images through the frozen unCLIP-SD, while LaVIT uses multiple transformer blocks as the visual decoder to 1668 reconstruct the discrete visual token sequence into a continuous feature map by learned queries 1669 (256 tokens), and then uses the conditional denoising U-Net to progressively recover image pixels from a Gaussian noise with the feature map as the condition, which can retain more image information generated by the LLM, thus achieving better image generation results. We freeze 1671 all visual modules and pre-tokenize images into discrete visual token sequences, and only load 1672 vision modules to generate actual images when evaluate on downstream image generation tasks. 1673 This can greatly reduce the training cost and improve efficiency. Notably, the vision modules in

our framework, though inherits from LaVIT, only performs "Stage 1 Tokenizer Training," while
 LaVIT's vision module also performs "Stage 2: Unified Vision–Language Pre-training."

1676 • Extended VL Vocabulary: as we stated in Section 3.1, SEED-LLaMA and LaVIT directly learn 1677 new visual token embeddings and initialize them with the distribution of original textual token 1678 embeddings. Visual tokens in the VL vocabulary correspond one-to-one with visual latent codes 1679 in the visual codebook. Considering the large difference in embedding dimensions between visual token embeddings in the VL vocabulary (4096) and visual latent codes in the visual codebook (32), we follow the factorized embedding parameterization in ALBERT (Lan et al., 1681 2020) to replace visual token embeddings $|V| \times 4096$ with two smaller embedding matrices $|V| \times 32, 32 \times 4096$, where |V| is the number of visual tokens in the VL vocabulary. The former 1683 matrix is directly initialized with visual latent codes in the visual codebook. In our preliminary 1684 experiments, we found that our initialization method of the VL vocabulary can significantly 1685 improve the performance of multimodal comprehension tasks, and has similar performance on multimodal generation tasks. 1687

• Training Objective: SEED-LLaMA and our framework both treat the multimodal sequence as 1688 a discrete sequence of image-text interleaved tokens, and train the LLM with an autoregressive 1689 objective to generate predictions of the next token. To mitigate the loss of detailed information caused by vector quantization, LaVIT designs two different training objectives for image-to-text and text-to-image cases. For image-to-text, LaVIT is similar to EMU and uses the continuous visual features of images as the condition to predict all discrete textual tokens and calculate 1693 loss for textual tokens only; for text-to-image, LaVIT is similar to SEED-LLaMA and our framework, which tokenizes images into discrete visual token sequences and calculates loss for both textual and visual tokens. To pursue simplicity, efficiency, and scalability, our framework 1695 is similar to SEED-LLaMA, which tokenizes images into discrete visual token sequences, so that different types of multimodal training data can be converted into discrete sequences of image-text interleaved tokens, and trained with a unified autoregressive objective to fully learn 1698 the extended vision-language vocabulary of the LLM.

Training Data and Settings: For our framework, we use MMGIC with 4M samples and IC 1700 with 52M image–caption pairs for pre-training; we use 1.21M multimodal instruction samples, 1701 1M aesthetic image-caption pairs and 1M samples from MMGIC for supervised fine-tuning; 1702 during pre-training and SFT stages, our framework freezes all parameters of vision modules 1703 and most parameters of the LLM, and only tunes partial parameters of the LLM (<8%): the 1704 extended VL vocabulary, additional LoRA modules, norm layers and a language model head 1705 layer. Existing MLLMs typically use more pre-training and fine-tuning data, and fully fine-tune 1706 all LLM parameters and partial vision module parameters. Take SEED-LLaMA as an example, uses more pre-training data (\sim 770M vs <60M), SFT data (\sim 145M vs <4M), and full-param fine-tuning. Its computation and data resources are extremely expensive and large (well over 10 1708 times that of our work). As shown in Table 3 and 4, compared with SEED-LLaMA, our baseline 1709 MLLM-MMGIC & IC requires significantly less data ($\sim 1/15$) and partial-param fine-tuning 1710 (<8%), achieving significant advantages on most benchmarks and comparable performance 1711 on COCO and VizWiz. Comparing with LaVIT, even though our framework does not convert 1712 images into continuous visual features to trade off lower loss of image representation and better 1713 multimodal comprehension performance, we still achieve significant advantages on three VQA 1714 benchmarks, and comparable performance on other benchmarks.

It should be emphasized that it is difficult to compare our baselines with existing MLLMs fairly, and the purpose of this paper is not to develop new frameworks, training objectives or benchmark SOTAs, but to explore the potential of multi-grained concept annotations for MLLMs. To this end, we explore different data recipes for multi-grained concept annotations and investigate the comparison and collaboration between MMG1C with widely-used coarse-grained image-caption data. Evaluations and in-depth analyses on 12 benchmarks for multimodal comprehension and generation are conducted in both pre-training and supervised fine-tuning stages.

1722

1723 Comparison with Grounding MLLMs. Fine-grained object region annotations are crucial for
 1724 grounding VLMs/MLLMs. They usually convert bounding box coordinates into discrete tokens in
 1725 text or visual markers in images. Take KOSMOS-2 (Peng et al., 2023) as an example, it falls into
 1726 the former category, utilizing object bounding box coordinates through special discrete position
 1727 tokens placed after corresponding texts in the caption. Unlike KOSMOS-2, MMGIC converts object
 bounding box coordinates into corresponding object regions (cropped from images), transforming

pure textual object-level annotations into multimodal annotations. Structured templates integrate
different granularities into unified image-text interleaved documents, enhancing MLLM's ability to
locate concepts in textual annotations to corresponding regions in images. With the unique complex
context processing capabilities of MLLMs and the LM autoregressive loss, our framework can
explicitly utilize multi-grained annotations in MMGIC to optimize multimodal alignment across
multiple granularities simultaneously, improving concept understanding and generation.

1734

1735 D.2 APPLICABILITY OF MMGIC TO EXISTING MLLMS

1737 Our framework combines the advantages of LaVIT and SEED-LLaMA to allow efficient exploration 1738 on various multimodal tasks, using only LM autoregressive loss without extra modules and loss functions for MMGIC, maintaining the framework's generality and efficiency. Compared to the 1739 current image-text pairs or image-text interleaved documents used by MLLMs, MMGIC can be seen 1740 as image-text interleaved documents with textual (caption, label, label description) and visual (object 1741 region) multi-grained concept annotations. They are well-integrated through our carefully designed 1742 structured templates, forming image-text interleaved documents that can be directly used by various 1743 MLLMs. Therefore, both our MMGIC dataset and our framework exhibit good generality. MMGIC 1744 do not impose any requirements or limitations on the model architecture or loss function of MLLMs. 1745 It can be directly applied to various MLLMs, whether they process images as continuous vectors 1746 (e.g., EMU, LaVIT, VL-GPT) or discrete tokens (e.g., CM3Leon, SEED-LLaMA, Chameleon). Our 1747 experiments on 12 multimodal comprehension and generation benchmarks also ensure the reliability 1748 and credibility of our conclusions, making sure that MMGIC can be quickly applied to other MLLMs' 1749 training. Overall, our model framework and training loss follow CM3Leon and SEED-LLaMA's spirit of unified discrete sequence modeling, offering simplicity, efficiency, and scalability. We did not 1750 make any special architectural or loss design modifications for MLLMs, ensuring good generalization 1751 and reliability of our framework, training loss, dataset, and experimental explorations. 1752

1753 Moreover, our focus is to explore the potential of multi-grained concept annotations for MLLMs. 1754 Directly fine-tuning existing SOTA MLLMs cannot fully demonstrate this point. Thus, we explore 1755 comparison and collaboration between MMGIC and IC starting from the multimodal pre-training 1756 phase. Unfortunately, all SOTA MLLMs use the LAION-5B dataset (and its subsets) for pre-training. 1757 However, the official data download channels were closed by the LAION team due to safety reviews (from December 19, 2023). This prevents us from retraining existing SOTA MLLMs from the 1758 pre-training stage in a controlled environment to directly prove MMGIC's superiority over IC. 1759 Therefore, we collected \sim 250M image–caption pairs from other widely used and recognized datasets 1760 and selected \sim 52M as the IC dataset for our experiments. Through the evaluation of three baselines 1761 across 12 multimodal understanding and generation benchmarks, we demonstrate the potential of 1762 MMGIC for MLLMs. Combining MMGIC and IC can leverage their respective advantages in terms 1763 of depth and breadth of concept representation, further improving MLLM performance.

1764

1765 1766 D.3 ANNOTATION SYNTHESIS IN MMGIC

1767 MMGIC is constructed from several public human-annotated object detection datasets, which covers 1768 many common concepts but currently cannot scale up to cover more concepts in the real world. We 1769 only automatically synthesize image captions for all images with the de-facto captioning model BLIP-1770 2 (Li et al., 2023b) and ranking model CLIP (Radford et al., 2021), and synthesize label descriptions 1771 for all labels with the strong GPT-4 (Achiam et al., 2023). For synthesized captions, we randomly 1772 sample 1000 captions to manually check the quality, and find that the quality is acceptable. For synthesized label descriptions, we carefully design prompt templates and human-annotated in-context 1773 examples, and manually check all synthesized label descriptions to ensure the quality. Besides, with 1774 the help of WordNet, we check category labels, and update the polysemous labels based on the specific 1775 data samples for better differentiation, e.g., "batter" \rightarrow "batter (ballplayer)" or "batter (cooking)". 1776

In this paper, we want to use and check existing human-annotated annotations as much as possible, and minimize the automatic generation of annotations, thus reducing noise that interferes with our exploration of multi-grained concept annotations for MLLMs. Hence, we take the above strategies to ensure the quality of synthesized annotations (captions and label descriptions). Moreover, in MMGIC, fine-grained attribute and relationship labels are only available in partial images from the original datasets. Although we try to automatically generate fine-grained attributes and relationship

annotations with external open-source models based on image captions (and images), after careful
manual checks, the synthesized fine-grained annotations are found to be of low quality, with nonnegligible hallucinations and noise, which are very likely to negatively impact MLLMs. Thus, only
high-quality manually annotated fine-grained annotations were used in MMGIC to avoid unnecessary
hallucinations and noise. We focus our limited resources and effort on carefully handling existing
data with meticulous manual checks and in-depth experimental exploration to ensure the quality of
MMGIC and the reliability of our exploration.

1789 To further ensure that such partially missing fine-grained annotations do not introduce additional 1790 bias, we compare MMGIC with only object labels and label descriptions, as well as MMGIC with 1791 all category labels and label descriptions, *i.e.*, MMGIC(CLD). We find that the performance is still 1792 improved by the partially available fine-grained attribute and relationship labels. More importantly, as shown in Section 4.4 and Figure 5 (Left), compared to MMGIC with only image captions, multi-1793 grained concept annotations in MMGIC can still help MLLMs achieve significant improvements of 1794 2.15 and 10.31 points in the corresponding "Instance Attribute" and "Instance Interaction" dimensions, 1795 respectively. Similar results are also observed in Appendix C.3 and Figure 8, where MLLMs trained 1796 with MMGIC can self-generate fine-grained category labels during the inference stage to further 1797 improve the performance by 1.79 and 5.16 points in the corresponding "Instance Attribute" and "Instance Interaction" dimensions, respectively. Above experimental results demonstrate that partially 1799 missing fine-grained attribute and relationship labels do not affect the effectiveness of MMGIC and the reliability of our exploration.

- 1801
- 1802

D.4 DIFFERENCE BETWEEN MMGIC AND EXISTING IMAGE-TEXT INTERLEAVED DATASETS

MMC4 (Zhu et al., 2023b) and OBELICS (Laurençon et al., 2024) are two widely-used open-source 1805 large-scale image-text interleaved datasets. They are constructed from massive HTML documents. 1806 MINT-1T (Awadalla et al., 2024) further enhances the scale and diversity, uniquely including data from PDFs and ArXiv documents. These datasets are large-scale, diverse, and also noisy. They are 1807 important for MLLMs to reason across image and text modalities, and have been used to further 1808 improve the capabilities and performance of MLLMs, especially in image-text interleaved scenarios. 1809 Our proposed MMGIC can also be seen as image-text interleaved documents with higher-quality 1810 and better image-text alignment, which has multi-grained concept annotations to help MLLMs 1811 better locate and learn concepts, thereby promoting vision-language alignment across multiple 1812 granularities simultaneously. It is a promising direction to automatically synthesize multi-grained 1813 concept annotations for these web-scale datasets.

1814

1815 D.5 IMAGE-INDEPENDENT LABEL DESCRIPTION GENERATION

In Section 2.2 and Appendix F.5, we describe the process of generating label descriptions for category labels in MMGIC. Considering that our description generator, GPT-4, do not see the images both during training and annotation synthesis, such general image-independent generation of label descriptions may introduce hallucinations and noise. In this paper, to ensure the quality of synthesized label descriptions and avoid hallucinations, we particularly consider the following aspects:

1822

1824

1825

1826

- Effectiveness of Synthesized Label Descriptions. As mentioned in Section 2.2, many previous works use WordNet and LLM to obtain general image-independent label descriptions, and successfully improve the model's understanding of concepts corresponding to these labels. The experimental results in Table 1 of Section 4.1 also show that label descriptions can bring significant performance improvements.
- Feasible Synthesis for Common Concepts. Whether in previous works (Menon & Vondrick, 2023) or in the four open-source datasets used in this paper, since the target labels are common concepts, LLM can generate label descriptions with good generality and rich visual information even without seeing the image.
- **Disambiguate Polysemous Labels.** We also use WordNet to disambiguate the polysemous labels in MMGIC, such as "batter" \rightarrow "batter (ballplayer)" or "batter (cooking)", to further reduce the ambiguity in the label descriptions, improve the generality and quality.
- **Careful Quality Control.** We not only carefully manual-check all generated label descriptions to ensure their quality and avoid introducing additional noise or hallucinations, but also carefully

Table 10: Training hyperparameters and cost for both pre-training and supervised fine-tuning (SFT)
 stages. MMGIC and IC refer to two pre-training datasets.

Training	hyperparam	eters	
Training Data	MMGIC	IC	SFT Data
Optimizer	AdamW	AdamW	AdamW
Learning Rate	2e-4	2e-4	4e-4
Weight Decay	0.05	0.05	0.05
Training Epochs	1	1	1
Warmup Ratio	0.1	0.1	0.05
Min learning rate ratio	0.1	0.1	0.1
Learning Rate Scheduler	Cosine	Cosine	Cosine
Batch Size	512	512	256
Maximum Token Length	2048	2048	2048
Tra	aining Cost		
GPU Device	32×N	VIDIA A10)0-80G
Training Time	$\sim 11h$	$\sim \! 20 h$	${\sim}5h$

design prompt templates and manual annotation examples (in Appendix F.5) to guide and improve the quality of the annotations generated by LLM. Interestingly, we find that the powerful GPT-4 can even automatically identify invalid noisy labels in the labels with the customized prompt templates, helping us to clean up these invalid annotations.

E EXPERIMENTAL DETAILS

1864 E.1 TRAINING DETAILS

¹⁸⁶⁶ The training hyperparameters and cost are shown in Table 10.

Training Loss. Similar to LLMs (Brown et al., 2020; Touvron et al., 2023b) and MLLMs (Ge et al., 2023b; Jin et al., 2023) under the autoregressive discrete framework, as we discussed in Section 3.2, we train our MLLMs with an autoregressive objective to maximize the likelihood of predicting the next visual or textual token in a multimodal discrete token sequence as follows:

$$L = \sum_{i=1}^{|u|} \log P(u_i \mid u_1, \dots, u_{i-1})$$

1876 where *u* denotes the multimodal token sequence, and u_i denotes the *i*-th token in the sequence. We can divide the multimodal token sequence into two parts: visual tokens and textual tokens. In the MLLM training, we find that the loss of the visual tokens L_V is larger than the loss of the textual tokens L_T . To balance the loss of the visual and textual tokens, we introduce a loss scale ratio $\alpha = 0.1$ to scale the visual token loss. The final loss is defined as $L = \alpha \cdot L_V + L_T$.

Image-First and Text-First. Our pre-training data consists of image-caption pairs (IC) and image-text interleaved documents (MMGIC). We follow common practice (Alayrac et al., 2022) to randomly select the image-first or text-first data template for each data sample as shown in Section F.6. The difference between the two templates is the order of appearance of the image and text in the sample, which determines the condition of each token during autoregressive training. For the multimodal instruction data we collect from public datasets, we follow their original data formats and do not change the order of appearance of the image and text in the sample.

- **Training Parameters.** In all training stages, we only tune partial parameters of the LLM: the VL vocabulary, additional LoRA modules (Hu et al., 2021), norm layers and a language model head



Figure 11: A brief illustration of the hyperparameter *mask_prob* in the pre-training stage when the template is image-first.

Table 11: Evaluation benchmarks, eval splits and eval metrics for the pre-training stage.

Dataset	Task Description	Eval Split	Metric
COCO	Scene description	test	CIDEr
NoCaps	Scene description	test	CIDEr
MS-COCO VIST	Text-Conditional Image Synthesis Contextual image synthesis	val(30K) val	FID, CLIP-I, CLIP-7 FID, CLIP-I

1914

1907

1909

1898

1899

1900 1901

layer, to greatly improve efficiency. All data are pre-tokenized, and we don't load all visual modules during training.

Loss Calculation Strategy. Moreover, in SFT stage, we follow LLaVA (Liu et al., 2023c) to only 1915 calculate loss for the answer tokens. In the pre-training stage, we find that existing MLLMs typically 1916 calculate loss for all tokens or only for the textual tokens when visual tokens as the condition, and vise 1917 versa. We conduct experiments to investigate the impact of different loss calculation strategies on the 1918 performance of MLLMs. Take the image-first template as an example, we introduce a hyperparameter 1919 mask prob to control the probability of only computing the loss for the textual tokens in each sample. 1920 As shown in Figure 11, when $mask_prob = 0$, we calculate the loss for all tokens in each sample; 1921 when $mask_prob = 0.4$, we calculate the loss for the textual tokens only in 40% of the samples, and 1922 calculate the loss for all tokens in the rest 60% of the samples; when $mask_prob = 1$, we calculate 1923 the loss for the textual tokens only in all samples. mask prob can be seen as a sample-level mask ratio of modality tokens, which controls the difficulty of the pre-training task. Higher mask_prob 1924 means easier pre-training task since more tokens are masked out. We conduct experiments with 0.1 as 1925 the step size to investigate the impact of $mask_prob = [0, 1]$ on the performance of MLLMs on image 1926 captioning and image generation benchmarks. Interestingly, the situation is more complicated than 1927 we thought. For image captioning tasks, the performance increases with the increase of mask prob, 1928 but the performance suddenly drops at $mask_prob = 1.0$, even lower than the performance at 1929 $mask_prob = 0.0$. For image generation tasks, the performance decreases first and then increases 1930 with the increase of $mask_prob$ (the inflection point is $mask_prob = 0.5$, the performance is better 1931 than that at $mask_prob = 0.0$), and then decreases again (the inflection point is $mask_prob = 0.8$). 1932 To balance the performance of comprehension and generation tasks, we choose $mask_prob = 0.9$ as 1933 the final hyperparameter. Since this is not the focus of this paper, we do not conduct in-depth research 1934 and analysis on this, but we believe that this is an interesting phenomenon worthy of further research, 1935 and we encourage future researchers to conduct in-depth research on this.

1936

1937 E.2 EVALUATION DETAILS

Evaluation Benchmarks and Metrics. We show the evaluation benchmarks and their corresponding data splits and metrics in Table 11 and Table 12 for the pre-training and supervised fine-tuning stages, respectively.

1942

Evaluation Strategy. In this paper, we evaluate the performance of MLLMs on both multimodal comprehension and generation tasks, including image captioning, text-to-image generation, visual

1010	14010 121 2		Philo and Crai me	and and an
1946				
1947	Dataset	Task Description	Eval Split	Metric
1948	COCO	Scene description	test	CIDEr
1949	NoCaps	Scene description	test	CIDEr
1950	VQAv2	Scene understanding QA	test	VQA Acc
1951	GQA	Scene understanding QA	test	VQA Acc
1952	VizWiz	Scene understanding QA	test	VQA Acc
1052	POPE	Visual Hallucination	-	F1-score
1955	MME	Multimodal Comprehension	test	MME score
1954	MMBench	Multimodal Comprehension	dev	Acc
1955	ScienceQA-IMG	Multimodal Comprehension	test	Acc
1956	SEED-Bench-IMG	Multimodal Comprehension	test	Acc
1957	MS-COCO	Text-Conditional Image Synthesis	val(30K)	FID. CLIP-I. CLIP-T
1958	VIST	Contextual Image Synthesis	val	FID, CLIP-I
1959	DreamBench	Subject-driven Image Editing	From KOSMOS-G	DINO, CLIP-I, CLIP-T
1060				

Table 12:	Evaluation	benchmarks.	eval s	splits and	1 eval	metrics	after SFT.
14010 12.	Liuluution	o'chiennan no.	e rui i	opineo an	4 C . u	metres	arter or r.

1944 1945

1962 question answering, multi-choice benchmarks, and image editing. We follow standard evaluation 1963 strategies and dataset splits to ensure the comparability of our experimental results. We implement 1964 the generation of visual and textual tokens in the evaluation stage based on the vLLM library, and 1965 generate the final image through the Diffusers library. Notably, multi-choice benchmarks have different designs and implementations in task instruction templates and evaluation strategies. The 1966 impact of different task instruction templates on the model's performance is significant and cannot 1967 be ignored. Therefore, we manually design at least 8 task instruction templates for each dataset and 1968 select the best performance as the final result to ensure the stability of the experimental results. The 1969 evaluation strategy is more complicated, and some benchmarks, such as MMBench (Liu et al., 2023d), 1970 also introduce an external LLM (GPT-4 (Achiam et al., 2023)) as the choice extractor. To pursue 1971 simplicity and efficiency, we directly select the highest probability of the candidate option tokens, 1972 i.e., "Yes" and "No" for POPE (Li et al., 2023c), "A", "B", "C", "D", etc. for other benchmarks, as 1973 the final answer. We don't introduce any external LLMs in our experiments since it will increase the 1974 cost and leads to unstable results (Liu et al., 2023a) compared with the above strategy. Similarly, we 1975 don't evaluate on benchmarks that use GPT-4 as the default evaluator, since it will increase the cost 1976 and has strong bias on models that tuned with GPT-4 synthetic data (Lin et al., 2023).

1977 1978

F DETAILS OF MMGIC DATASET

1979 1980

In this section, we provide the details for constructing MMGIC dataset. The statistics of MMGIC dataset are shown in Table 13.

1983

1984 F.1 DATASET COLLECTION

We collect four large-scale object detection dataset, including Open Images (Kuznetsova et al., 2020),
Objects365 (Shao et al., 2019), V3Det (Wang et al., 2023a) and Visual Genome (Krishna et al., 2017)
for MMG1C.

1989 BigDetection (Cai et al., 2022) has found that Open Images and Objects365 are biased towards certain 1990 scale (small or large) of objects, while the improved object annotations for these two datasets in 1991 BigDetection are balanced across object scales. Therefore, we use the improved object annotations 1992 of Open Images and Objects365 from BigDetection. For Open Images, we merge the BigDetection 1993 object annotations with the attribute and relationship annotations of Open Images. Specifically, for each object region in BigDetection, we calculate the *IoU* with each Open Images object region in 1995 the same image, and select the Open Images object region with the largest *IoU* to pair with it. We only keep the relationship and attribute annotations in Open Images where subject object and the 1996 object object are both successfully paired. We do not use the LVIS (Gupta et al., 2019) subset in 1997 BigDetection since all images in LVIS come from COCO (Lin et al., 2014), which is our downstream

Table 13: Statistics of MMGIC dataset and four object detection datasets underlying it. The last 4 columns are based on the final tokenized dataset and the other columns are before tokenization. The 2000 # Concepts column denotes the number of unique category label-description pairs in each dataset. 2001 Note that Visual Genome and V3Det are repeated 3 times to keep data balance in the final dataset. 2002

2004 2005		Dataset	# Regions	# Images	# Concepts	# Visual Tokens	# Textual Tokens	# Used Regions	# Samples
2006	0	Open Images	7.16M	1.44M	642	0.70B	0.47B	6.21M	1.44M
2007	1	Objects365	15.71M	1.78M	600	1.31B	0.97B	12.58M	1.78M
2008	2	V3Det	1.18M	0.18M	12,976	0.26B	0.21B	2.30M	0.53M
2009	3	Visual Genome	1.50M	0.10M	20,830	0.28B	0.28B	2.77M	0.31M
2010	4	MMGIC	25.55M	3.48M	35,048	2.54B	1.92B	23.86M	4.06M

2013 evaluation dataset. For V3Det and Visual Genome, we use all their original annotations for further 2014 pre-processing.

2015 Since BigDetection dataset does not involve any changes to the positions of the objects in the Open 2016 Images and the objects in the same image are a subset of the objects in the Open Images, the matching 2017 process is correct and appropriate. We use *IoU* instead of exact matching because BigDetection uses 2018 coordinates to provide object positions, whereas Open Images use ratios, where using exact match 2019 may introduce errors. 2020

For V3Det, there are 60 images that were not successfully downloaded, and we simply discarded 2021 these images.

2023 2024

2025

2026

2027

2033

2034

2035

2041

1998

2003 2004

2011 2012

F.2 DATASET PRE-PROCESSING

Open Images, Objects365 and V3Det. There are 3 steps to pre-process these three datasets. Notice that we directly use the train and dev splits from the original datasets and don't use the test set.

- 2028 1. Clean labels: to remove the noise that may be introduced by long-tail distribution, we count the number of occurrences of each label in each dataset, and then select a threshold to remove 2029 low-frequency labels (typically noise). For V3Det, we remove object labels with frequency less 2030 than 3. For Open Images and Objects365, we don't remove any label since each label appears 2031 many times and has high quality. 2032
 - 2. Clean objects and object-related information: we remove objects that have illegal coordinates, exceed the image range or the corresponding label is already removed. For Open Images, we also remove the attribute and relationship annotations corresponding to these illegal objects.
- 3. Clean images: since the de-facto visual encoder (*i.e.*, Vision Transformer, ViT) requires a 2036 square input image, we resize all images to a fixed resolution of 224×224 pixels. Following 2037 standard practices (Zeng et al., 2023), to improve the data quality, we first remove images with 2038 short edges less than 224 pixels or aspect ratios greater than 3.0 or less than 0.33, to prevent the 2039 image from changing too much after resizing. We also remove the images without any objects. 2040
- **Visual Genome.** Visual Genome is a dataset trying to connect structured image concepts to 2042 language. Each image is annotated with: 1) image-level info, *i.e.*, image captions; 2) region-level info, 2043 including objects in the region and the region description; 3) object-level info, including object label 2044 and location, the attributes of objects, the relationships between objects, with all the labels mapped to 2045 WordNet synsets. We use the image-level info and object-level info for further pre-processing which 2046 require the following steps. 3 2047

²⁰⁴⁸

³We found that region-level annotations may not be suitable for constructing MMGIC since: 1.many regions 2049 only contains a single object; 2 the annotations of the objects in region may not be exactly in the region; 3 the 2050 region description may not be exactly matched with the region itself. Hence, We use v1.4 version of the dataset 2051 which includes image-level and object-level annotations except the attribute of objects. Since we found that the objects in v1.2 is the subset of that in v1.4, we use the attribute annotations in v1.2

1. **Clean images**: since all images need to be resized to 224×224 and input to the visual modules, we remove the images with the shorter side length less than 224, images with an aspect ratio greater than 3 or less than 0.33 to prevent the image from changing too much after resize.

- 2055
 2056
 2056
 2057
 2. Clean object annotations: we remove all objects that have more than one synset as we found that they are mostly noisy. We also remove objects that have illegal coordinates or exceed the image range.
- 2058 3. Construct, clean object label, clean object annotations again: in Visual Genome, the same 2059 object label may have different semantics distinguished by the synsets they are mapped to. To 2060 mitigate the noise brought by long-tail distribution, firstly, we count the frequency of each 2061 (object label, synset) pair if the mapped synset is not empty. Secondly, for those object labels 2062 that aren't mapped to any synset, we regard that it is mapped to the synset that has the highest 2063 frequency of the same object label in the first step (if there isn't any same object label, we regard that it is mapped to the empty synset) and add its frequency to the first step. Finally, we 2064 remove all the (object label, synset) pairs with frequency less than 2 and also remove all the 2065 corresponding objects. To get the final object label name, for the object labels which have more 2066 than two types of synsets, we manually check them with help of WordNet and if so, we add 2067 some keywords corresponding to the synset behind the object label name as the final name for better differentiation (e.g. (batter, batter.n.01) \rightarrow "batter (baseball player)", (batter, batter.n.02) 2069 \rightarrow "batter (semiliquid mixture)"). For other object labels, we directly use the object label name 2070 as the final name. 2071
 - 4. Clean images again: we remove all the images that don't have any object annotation.
- 5. Clean relationship and attribute annotations: firstly, we remove all the relationships where the corresponding object is removed by step 2 and step 3. Secondly, we count the frequency of each label and remove the relationship/attribute whose frequency is less than 5 or is empty. We don't add additional information in label name as in step 3 as we have checked manually that each label name has exactly one semantic.
- 8. Remove unreasonable labels and annotations: for each label, we ask ChatGPT if it is an reasonable label in Section F.5. For unreasonable labels, we remove them and the corresponding object, attribute, relationship annotations.
- 7. Dataset split: for images in Visual Genome with a coco id, we follow Karpathy's split: images in 'val' and 'restval' are used as the validation set, and images in 'test' are discarded. For images in Visual Genome with a flickr id but no coco id, we match them with images from flickr30k and discard the matched images. All other images are used as the training set.

2006 F.3 OBJECT ANNOTATION PRE-PROCESSING

2085

The fine-grained object annotations provided in these datasets include bounding box coordinates and category labels for each object. To accommodate varying aspect ratios of object bounding boxes 2089 and the requirement for a square input image, following KOSMOS-G (Pan et al., 2023), we crop 2090 a new larger square object region S_i that contains the original object region $R_i, R_i \subseteq S_i$, with 2091 their centers aligned as closely as possible. Since S_i may include not only the object region R_i 2092 but also the surrounding object regions, we design a simple but effective strategy to update the object label annotations of S_i by merging the object label annotations of R_i with the object label 2093 annotations of surrounding object regions, to improve the quality of object annotations. For each R_i 2094 and its corresponding new region S_i , we calculate the intersection-over-area (IoA) between S_i and 2095 all $R_{j(j\neq i)}$, *i.e.*, $IoA(S_i, R_j) = \frac{Area(S_i \cap R_j)}{Area(R_j)}$. If $IoA(S_i, R_j) \ge 0.8$, we consider R_j as a part of S_i 2096 and update the annotations of R_i to the annotations of S_i . This strategy can help to include more 2097 relevant object regions in the region S_i and improve the quality of object annotations. After that, we 2098 remove S_i with duplicate annotations and keep the unique ones. 2099

2100 Since each visual token in our framework corresponds to a 14×14 pixel region, we remove S_i 2101 with edge length less than 28 pixels or greater than 182 pixels to ensure that the object region is not 2102 too small or too large. We also remove all the images that don't have any bounding box and resize 2103 all images to 224×224 pixels. Finally, for objects in the image, MMGIC provides visual tokens 2104 of cropped regions and textual tokens of corresponding fine-grained category labels and location 2105 descriptions. The pre-processed object regions in the images are ready to be tokenized and fill in the 2106 template in Section F.6. 2106 Specifically, instead of transform bounding box coordinates into textual tokens in the text (Chen et al., 2107 2021; Liu et al., 2023c; Peng et al., 2023) or visual markers in the image (Yao et al., 2024; Shtedritski 2108 et al., 2023; Yang et al., 2023), for each object in the image, we directly provide visual tokens for the 2109 cropped region and textual tokens for the corresponding category labels. This not only upgrades our 2110 data from simple image-text pairs to more complex image-text interleaved data which has shown to be beneficial for vision-language learning (Lin et al., 2023; McKinzie et al., 2024), but also helps to 2111 locate and align concepts in the image and in the text by providing both the whole image and object 2112 regions. 2113

2115 F.4 CAPTION SYNTHESIS

Since only partial images are annotated with image captions in these datasets, following BLIP (Li et al., 2022a) and LAION-COCO (Schuhmann et al., 2023), we automatically synthesize captions for all images in the following steps:

- 1. Generate 10 candidate captions for each image with BLIP-2 (Li et al., 2023b);
 - 2. Rank candidate captions based on image-caption similarity scores calculated by CLIP (Radford et al., 2021);
- 2123 2124 2125

2128

2130

2120

2121

2122

2114

- 3. Filter too short (< 5 words) or too long (> 25 words) captions or captions with low image– caption similarity scores (< 0.25);
- 4. Select the Top-1 caption as the final caption if exists, otherwise repeat the above steps for 10 times. If no caption is selected, we select the caption with the highest similarity score.
- 2129 F.5 LABEL DESCRIPTION GENERATION

Label descriptions are corresponding concept descriptions of concrete concepts in the image, which 2131 convey understanding about a concept by visually observed details and relevant knowledge. Inspired 2132 by the success of previous works (Shen et al., 2022; Yao et al., 2022; Menon & Vondrick, 2023) that 2133 introduces label descriptions from WordNet (Miller, 1992), Wiktionary (Meyer & Gurevych, 2012) 2134 and LLMs (Brown et al., 2020) to help understand concepts, with the help of GPT-4 (Achiam et al., 2135 2023), we generate label descriptions for fine-grained category labels (object, object attribute, and 2136 relationship between objects) and manually check them. We design prompt templates for each type 2137 of category labels and carefully provided human-annotated examples. For V3Det, we directly use the 2138 label descriptions it provided. For Open Images, Objects 365 and Visual Genome, we generate label 2139 descriptions for object labels, attribute labels and relationship labels respectively.

Specifically, in system prompt, we first describe the task and add some tips that LLM might focus when generating descriptions. Then we give some examples. In user prompt, we instruct the LLM to generate the description. We also ask LLM to generate 'Invalid' first if the given label is noisy in Visual Genome These labels along with their corresponding annotations will be further removed. The full prompts are shown in Table 17, 18, 19, 20, 21, 22.

- 2146 F.6 TEMPLATE
- 2148 After completing all above steps, for each image sample, we have the following annotations:
- 2149 2150

2151

2152

2153

2154

2155

2159

2145

2147

- Coarse-grained concept annotations: captions and the localized narrative captions.⁴
- **Fine-grained concept annotations:** object labels, attribute labels and relationship labels along with their label descriptions.
 - **Region-level annotations**: the object regions of the image obtained in Section F.3 along with their locations and object labels.

We design templates to formulate the above annotations into data samples used in the pre-training stage. There are 2 types of templates: image-first template shown in Table 23 and text-first template shown in Table 24. Each template has 2 parts: image-annotation part and object-annotation part.

⁴Only part of images in Open Images have localized narrative captions.

2160 Image-first template. For image-annotation part or object-annotation part, we place the corresponding text after the image or region, aiming to enable the model to learn visual understanding.
 2162

2163 • **Image part**: we first place the image. Secondly, we place the coarse-grained concept annotation 2164 including caption and the localized narrative caption. Thirdly, we place image-level fine-grained concept annotations including object labels, attribute labels and relationship labels along with 2165 their descriptions. For object labels, we first place the label name and then place the label 2166 description, each object label name in the image appears exactly once. For attribute labels, 2167 we first place the attribute label and then place the associated objects, followed by the attribute 2168 label description. For relationship labels, we first place the relationship label and then place 2169 the associated subjects and objects, followed by the relationship label description.⁵ 2170

• **Object-annotation part**: we place the object regions of the image with their location descriptions and associated object labels after them.

Text-first template. Different from the image-first template, we place the corresponding text before
 the image or region, aiming to enable MLLMs to learn visual generation. Note that we place the
 captions close to the images, to help MLLMs better learn to utilize caption to generate image.

2178 F.7 PRE-TRAINING SAMPLE GENERATION

2171

2172 2173

2177

2179

2185

2186

2187

2213

Pre-tokenize images and regions. For images and regions, we directly use the visual tokenizer of LaVIT(Jin et al., 2023) to tokenize them into visual token sequences.

Fill in the template and tokenize data. There are 4 steps to generate a training sample. For V3Det and Visual Genome, we repeat the process 3 times to include as many regions as possible.

- 1. Select a template: we randomly select the image-first or text-first template with a probability of 0.5.
- 2. Fill in the template with textual annotations: For image-annotation part, we first fill the 2188 corresponding text. Then we remove the descriptions with a probability of 0.5 to prevent model 2189 overfitting which might be caused by description repetitions. For each object-annotation, 2190 we first fill in the region location. Specifically, the 224×224 square image is divided into a 2191 3×3 grid, with each cell named sequentially from top to bottom and left to right as Top Left, 2192 Top Middle, Top Right, Middle Left, Center, Middle Right, Bottom Left, Bottom Middle, and 2193 Bottom Right. The location of a region is designated by the name of the grid cell where its 2194 center is located. Then we fill in the object label with the annotations obtained in Section F.3. 2195 To get the final object-annotation part, we shuffle the regions and place as many regions as 2196 possible without exceeding the maximum tokenized length.
- 3. Tokenize the data and fill the template with visual annotations: firstly, we tokenize the data using the tokenizer of with [IMG], [/IMG] added as special tokens. Secondly, we replace the positions corresponding to [IMG] using the corresponding visual token sequence with [IMG] and [/IMG] inserted before and after it. Thirdly, we add <s> token and </s> token to the whole token sequence.
- 4. Dealing with samples with token length more than 2048: for those samples with object regions, we remove one region a time and go back to step 2 until the token sequence length is less than 2048. For those samples without regions but with descriptions, we go back to step 2 but fill in the template without descriptions. Other samples are discarded.

Statistics. Table 13 shows the statistics of the constructed MMGIC. The frequency distribution of object labels, attribute labels and relationship labels are shown in Figure 12 respectively, which illustrates that our label types are broadly distributed. Moreover, only a few labels appear with particularly low frequency. For object labels, attribute labels, and relationship labels, labels with frequency less than 5 accounted for 5%, 1.5%, and 0.3% of the total number of their respective label
types, respectively.

⁵We design the template in this way to reduce input length and duplication.



Figure 12: Label frequency distribution of objects, attributes and relationships in MMGIC.

DETAILS OF IC DATASET G

2260 2261

2259

2262 To strike a balance between increasing concept breadth, reducing data noise and improving efficiency, 2263 we collect several large-scale public image-caption datasets (Ordonez et al., 2011; Sharma et al., 2264 2018; Changpinyo et al., 2021; Schuhmann et al., 2022a;b; Sun et al., 2024), and follow LLaVA (Liu 2265 et al., 2023c) to first filter them by the frequency of noun-phrases extracted by spaCy from their given synthesized captions, and then automatically synthesize captions same as MMGIC. We name 2266 this dataset as IC, where 52M unique images are collected and selected, almost 15 times more than 2267 MMGIC. IC contains 3.19B visual tokens and 5.14B textual tokens.



Figure 13: Comparison of noun-phrase statistics before and after filtering IC (not including aesthetic part). The x-axis represents the unique noun-phrases ordered by frequency in the descending order. The total number of unique noun-phrases are reported in the legend.

2283 G.1 DATASET COLLECTION

2282

2289

2290

2291

2292

2293

2294

2295

2296

2297

2308

As we discussed in the last paragraph of Section D.2, LAION-5B dataset (and its subsets) cannot be downloaded from their official website. Therefore, we collect several large-scale public imagecaption datasets from other widely-used sources:

- BLIP-Captions: BLIP (Li et al., 2022a) provides image urls, web captions and synthesized COCO-style captions for CC3M (Sharma et al., 2018), CC12M (Changpinyo et al., 2021), SBU (Ordonez et al., 2011), LAION-400M (Schuhmann et al., 2021) (BLIP only uses and provides 115M of them). Finally, we collect valid ~ 116M images with both web and synthesized captions and denote them as BLIP-Captions.
- CapsFusion-120M: similarly, CapsFusion (Yu et al., 2023b) provides image urls, web captions and synthesized COCO-style captions from LAION-2B (Schuhmann et al., 2022a). Moreover, they also propose a LLM-based pipeline, denoted as CapsFusion, to merge two types of captions into a single caption, which is a well-structured sentence retaining the detailed real-world knowledge from web captions. We will discuss it later in the Section G.3.
- LAION-COCO-Aesthetic and JourneyDB: we follow LaVIT (Jin et al., 2023) to collect aesthetic image–caption data to further improve the diversity and aesthetic quality of image–caption data. LAION-COCO-Aesthetic is an unofficial dataset that contains 10% samples of the LAION-COCO (Schuhmann et al., 2023) dataset filtered by some text rules (remove url, special tokens, etc.), and image rules (image size > 384 × 384, aesthetic score > 4.75 and watermark probability
 cotic dataset that contains ~ 8M samples. JourneyDB is a large-scale generated aesthetic image–caption dataset that contains ~ 4M high-resolution Midjourney synthetic images and synthetic captions rewrote by GPT3.5 from real user prompts.

Finally, we collect ~ 229 M general image–caption samples and ~ 12 M aesthetic image–caption samples from the above open-source datasets.

2309 G.2 DATA FILTERING

2310 We follow LLaVA (Liu et al., 2023c) to downsample the above $\sim 229M$ general image-caption 2311 samples based on the frequency of noun-phrases extracted by spaCy from their given synthe-2312 sized captions, thereby reducing training costs while ensuring concept breadth (or concept cov-2313 erage). We denote the web captions, synthesized COCO-style captions and CapsFusion captions 2314 as caption_origin, caption_coco and caption_capsfusion respectively. We ex-2315 tract noun-phrases from caption_coco of BLIP-Captions and caption_capsfusion of 2316 CapsFusion-120M. Following LLaVA, we skip noun-phrases whose frequency is smaller than 20, as 2317 they are usually rare combinations of concepts and attributes that have already been covered by other 2318 captions. We then start from the noun-phrases with the lowest remaining frequency, add the captions 2319 that contain this noun-phrase to the candidate pool. If the frequency of the noun-phrase is larger than 50, we randomly choose a subset of size 50 out of all its captions. This results in $\sim 40M$ image-text 2320 pairs that can be successfully downloaded. The comparison of noun-phrase statistics before and after 2321 filtering IC is shown in Figure 13. The filtered dataset shows a good coverage of concepts whose

frequency is higher from 20, but with a smaller number of image-text pairs. Finally, we get $\sim 52M$ image-caption pairs as IC.

2325 G.3 CAPTION SYNTHESIS

As stated in Section F.4, we follow BLIP and LAION-COCO to synthesize captions for MMGIC. However, for IC, their images are more diverse and the synthesized COCO-style captions may lack in-depth real-world details. Following CapsFusion (Yu et al., 2023b), we synthesize caption_capsfusion based on caption_origin and caption_coco for each image, except for JourneyDB. Their GPT3.5-rewrote captions are already well-structured and contains detailed real-world knowledge, which can be directly seen as caption_capsfusion. In our preliminary experiments, we found that using both caption_coco and caption_capsfusion together is significantly better than using only one of them. The final data template (image-first as an example) is:

Image: [image]
Caption: [caption_coco]
Detailed caption: [caption_capsfusion]

H DETAILS OF SFT DATA

In this section, we provide the details of constructing our SFT data. The statistics of the collected dataset and final constructed dataset are shown in Table 14 and Table 15, respectively.

2345 2346 H.1 DATASET COLLECTION

The dataset we collected are shown in Table 14. Note that we playback 1M samples from MMGIC to avoid forgetting the knowledge in the pre-training stage. The preprocessing details are the same as Section F. We also sample 1M data from LAION-COCO-Aesthetic (LAION, 2024a) to keep the image captioning and text-to-image generation ability of the model, with half used as image captioning and half used as text-to-image generation respectively.

2353 H.2 SAMPLE GENERATION 2354

We use the template shown in Table 16 to format our collected dataset to training samples following according to the dataset task. For text-to-image generation and image editing, we follow the template provided by VL-GPT (Zhu et al., 2023a). For other task, we follow the template provided by LLaVA v1.5 (Liu et al., 2023b). We follow the system prompt used by VL-GPT, which is "You are a helpful, respectful and honest assistant. Always answer as helpfully as possible, while being safe. Your answers should not include any harmful, unethical, racist, sexist, toxic, dangerous, or illegal content. Please ensure that your responses are socially unbiased and positive in nature". The final statistics of the constructed dataset are shown in Table Table 15.

2362

2352

2335

2336

2337

2338 2339 2340

2341 2342

- 2363 2364
- 2365
- 2366
- 2367
- 2368
- 2369
- 2370
- 2371
- 2372
- 2373
- 2374
- 2375

Table 14: Statistics of the collected SFT data. VG denotes Visual Genome (Krishna et al., 2017).

Туре	Task	Datasets Involved	# Samples
	Image conversation	UniMM-Chat, Llavar, LLaVA	0.37M
	Open-ended VQA	VQAv2, GQA, OKVQA, OCRVQA	0.24M
Mutlimodal Understanding	Multi-chocie VQA	A-OKVQA	0.07M
_	Detailed Caption	ShareGPT4V, LaionGPT4V	0.11M
	Image Caption	LAION-COCO-Aesthetic, TextCaps	0.52M
Matting dat Comparties	Image Editing	Instructpix2pix, MagicBrush	0.32M
Muthmodal Generation	Text2Image	LAION-COCO-Aesthetic	0.50M
Text-only	Text-only	Alpaca, ShareGPT	0.09M
MMCIC Diavibaal	Image-first	Open Images, Objects365, V3Det, VG	0.50M
MMOIC Playback	Text-first	Open Images, Objects365, V3Det, VG	0.50M

Table 15: Statistics of the tokenized SFT data.

Туре	# Visual Tokens	# Textual Tokens	# Samples
Mutlimodal Understanding	115M	355.0M	1.31M
Mutlimodal Generation	104M	90.7M	0.82M
Text-only	0M	14.1M	0.09M
MMGIC Playback	602M	579.0M	1.00M
Total	821M	1038.8M	3.21M

Table 16: SFT data template. Each round starts with <s> and end ends with </s>. For each round, the instruction is the content between [INST] and [/INST], the other content is the response.
"Detail caption instruction" means one of our hand-crafted detailed caption instructions, e.g., "Explain the visual content of the image in great detail.", and "Analyze the image in a comprehensive and detailed manner.". "Editing instruction" is sample-specific, e.g., "Change the table to a dog.", and "Remove one potted plant".

Task	Template
Image Conversation	$\label{eq:ss} $$ $$ [INST] <>\n{SFT_SYSTEM_PROPMT}\n<>\n{IMG}\n{question} [/INST] {response} $$ $$ $$ $$ $$ $$ $$ $$ $$ $$ $$ $$ $$$
Open-ended VQA	<pre> <s> [INST] <<sys>>\n{SFT_SYSTEM_PROPMT}\n<</sys>>\n\n[IMG]\n{question}\nAnswer the question using a single word or phrase. [/INST] {response} </s></pre>
Multi-Choice	<pre><s>[INST] <<sys>>\n{SFT_SYSTEM_PROPMT}\n<</sys>>\n\n[IMG]\n{question}\nAnswer with the option's letter from the given choices directly. [/INST] {response} </s></pre>
Detailed Caption	<pre> <s> [INST] <<sys>>\n{SFT_SYSTEM_PROPMT}\n<</sys>>\n\n[IMG]"Detail caption in- struction"} [/INST] {response} </s></pre>
Image Caption	<pre><s>[INST] <<sys>>\n{SFT_SYSTEM_PROPMT}\n<</sys>>\n\n[IMG]\nProvide a one-sentence caption for the provided image. [/INST] {response} </s></pre>
Image Editing	<pre> <s> [INST] <<sys>>\n{SFT_SYSTEM_PROMPT}\n<</sys>>\n\n[IMG] {"Editing instruction"} /INST] Here is the edited image: [IMG] </s></pre>
Text2Image	<pre><s>[INST] <<sys>>\n{SFT_SYSTEM_PROMPT}\n<</sys>>\n\nCreate an image that visually represents the description: {Caption} [/INST] Here is the image: [IMG] </s></pre>
Text-only	$\label{eq:s} $$ $ [INST] <\n{SFT_SYSTEM_PROPMT}\n<\n{question} [/INST] {response} $$ $ $ $ $ $ $ $ $ $ $ $ $ $ $ $ $ $ $
MMGIC (Text-first)	<pre><s>[INST] <<sys>>\n{SFT_SYSTEM_PROPMT}\n<</sys>>\nCaption and fine-grained an- notations in the image-level part} [/INST] {Image in the image-level part} </s><s>[INST] {Location and object label in the related region} [/INST] {Region of the related region} </s></pre>
MMGIC (Image-first) <s> [INST] <<sys>>\n{SFT_SYSTEM_PROPMT}\n<</sys>>\n\{Image in the image-level part} [/INST] {Caption and fine-grained annotations in the image-level part} </s> (INST] {Region of the related region} [/INST] {Location and object label in the related region}

	System Prompt
	# Instructions for Object Category Description Generation
	You are a helpful respectful and honest assistant
	Now you are an expert in generating descriptions for category labels of objects in an image
	You are given some examples, each with their object category label and description.
	Your goal is to generate a description of a given object category label to help people better understa
	from both vision and language modalities.
	Informative, concise, accurate and clear descriptions are expected.
	Here are some useful tips for generating object category descriptions:
	1. Universality: Focus on features common to most instances of the object category.
	2. Multiple Semantics: Select the meaning of the most relevant and likely object category label in an in
	For instance, as an object category label, "bank" can be "a financial institution" or "a landform alongs
	3 Distinctive Features: Emphasize unique aspects differentiating the object from other similar objects
	4. Relevant Knowledge: Incorporate important concepts, historical, or cultural information that enrich
	understanding of the object but avoid excessive details.
	5. You can optionally focus on the following aspects when generating category descriptions:
	1. What are useful features for distinguishing the object of the given category label in an image?
	2. What does the object of the given category label in an image look like?
	5. what are the identifying characteristics of the object with the given category label in an image?
	4. what are the key visual indicators that help identify the object category label in an image?
	Here are some examples:
	"'
	Object Category Label: malayan tapir
	Cotagory Description: The Meleven tenir, a distinctive memmal found prodominantly in Southeast
	is known for its unique coloration: It boasts a unique appearance, featuring a black body with a w
front section. It resembles a large pig in shape, with a short, prehensile trunk, small eyes, a	front section. It resembles a large pig in shape, with a short, prehensile trunk, small eyes, and pointed
	Primarily nocturnal, tapirs are herbivorous and thrive in dense forests near water sources.
	"'
	Object Category Label: sheep
	Category Description: Sheep are medium-sized, quadrupedal mammals, typically covered in curly wo
	neece that varies from white to brown. Recognizable by their stout, flutty appearance, they have a distin-
	for their gentle demeanor and are primarily raised for wool, meat, and milk.
"' "'	"
	Object Category Label: compass
	Category Description: A compass is a navigational instrument typically used for orientation and direct
	finding. It consists of a magnetized needle that aligns with Earth's magnetic field, pointing toward
	north. The device usually has a circular scale marked with directions (North, South, East, West) and deg
	for precise navigation.
	u. Uran Dronant
	User Frompt # Object Category Description Generation
	Please directly generate an informative, concise, accurate and clear description for the given object cate
	label in about 50 words.
	"
	Object Category Label: {category label}

System Prompt
Instructions for Attribute Category Description Generation
You are a helpful, respectful, and honest assistant.
Now you are an expert in generating descriptions for attribute category labels of objects in an image.
You are given some examples, each example is in the following format:
"'
Attribute Category Label: the attribute category label required to generate a description.
Related Object Category Labels: the category labels of the objects which may be associated with the attribu-
category label, separated by commas.
Attribute Category Description: the generated description of the attribute category label
""
Your goal is to generate a description of a given attribute category label to help people better understand
from both vision and language modalities.
informative, concise, accurate and clear descriptions are expected.
Here are some useful tips for generating attribute category descriptions:
1. Universality: Focus on features common to most instances of the attribute category.
2. Multiple Semantics: Select the meaning of the most relevant and likely attribute category label in image. For instance, as an attribute category label "awake" is an adjective meaning "not asleen" rather the
a verb meaning "to stop sleeping or wake up from sleep".
3. Distinctive Features: Emphasize unique aspects differentiating the attribute from other similar attribut
4. Appearance, Sensory and Effect: Describe how the attribute typically appears or is perceived on
objects. This could include color, texture, size, sense, function, or any other aspect that is significant for attribute
5. You can optionally focus on the following aspects when generating attribute category descriptions:
1. How is the attribute category label usually observed or sensed?
2. What impact does this attribute have on the object, its perception or function?
3. What are the identifying characteristics of the given category label of an object in an image?
4. What are the key visual indicators that help identity the attribute category laber in an image:
Here are some examples:
"
Attribute Category Label: smiling
Related Object Category Labels: Man Woman Girl Roy Raby Child Animal
Actained Object Category Labers. Main, Wollian, Olli, Duy, Dauy, Clillu, Allillial,
Attribute Category Description: Smiling is a facial gesture where the mouth corners lift upwards, of
revealing the front teeth, and is usually associated with emotions like happiness, kindness, or amusement
generally signifies a person feeling joyful or content, making it a universally understood symbol of posit
««
User Prompt
Attribute Category Description Generation
Please directly generate an informative concise, accurate and clear description for the given attribute catego
label in about 50 words.
Attribute Category Label: {attribute category_label}

Attribute Category Description:

Sustan Drampt
Instructions for Relationship Category Description Generation
Ver an a helpful manageful and han at assistant
Now you are an expert in generating descriptions for relationship category labels of objects in an image.
You are given some examples, each example is in the following format:
Relationship Category Label: the relationship category label required to generate a description.
Palated Subject Object Pairs: the subject object category label pairs associated with the relationship category
label. Each pair in the form of [Subject Category Label, Object Category Label], different pairs are separate
by commas.
Relationship Category Description: the generated description of the relationship category label.
"'
Your goal is to generate a description of a given relationship label to help people better understand it fro
both vision and language modalities.
Informative, concise, accurate and clear descriptions are expected.
Here are some useful tips for generating attribute category descriptions:
1. Universality: Focus on features common to most instances of the relationship category.
2. Multiple Semantics: Select the meaning of the most relevant and likely relationship category label in a
image. For instance, as a relationship category label, "truck" is a verb meaning "convey by truck" rather th
a noun meaning "a large, heavy road vehicle used for carrying goods and materials".
relationships.
4. Nature of Relationship: Describe the nature of the relationship (e.g., spatial, action, functional, hierarchic
temporal, social) and how the subject and object interact or relate to each other.
5. You can optionally focus on the following aspects when generating attribute category descriptions:
2. What are the key characteristics or significance of the relationship?
3. What are common or typical scenarios in which this relationship is observed?
4. What are the key visual indicators that help identify the relationship category in an image?
mere are some examples:
Relationship Category Label: on
Related Subject-Object Pairs: [Bell pepper, Countertop], [Woman, Bicycle], [Tomato, Cutting board],
Relationship Category Description: The relationship lobal "on" indicates that the subject is positioned abo
and in contact with the surface of the object, without being suspended or supported by anything else. The
implies a direct interaction where an object rests upon or integrates into the surface, thereby becoming
prominent or integral part of its overall structure.
User Prompt
Relationship Category Description Generation
Please directly generate an informative, concise, accurate and clear description for the given relationsh category label in about 50 words
category laber in about 50 words.
"'
Relationship Category Label: {relationship category_label}

Relationship Category Description:

2592 2594 2595 Table 20: Prompts for the relationship object label description generation for Visual Genome. Note 2596 that there are 3 examples with 2 examples omitted. 2597 2598 System Prompt # Instructions for Object Category Description Generation 2600 2601 You are a helpful, respectful, and honest assistant. 2602 Now you are an expert in generating descriptions for category labels of objects in an image. 2603 You are given some examples, each with their object category label and description. 2604 Your goal is to generate a description of a given object category label to help people better understand it from both vision and language modalities. 2605 Informative, concise, accurate and clear descriptions are expected. 2606 2607 Attention, the given object category label may be invalid. In such cases, you should generate "Invalid." as its 2608 description and then explain why it is invalid after "Invalid.". 2609 Here are some useful tips for generating object category descriptions when the given object category label is 2610 valid: 2611 2612 1. Universality: Focus on features common to most instances of the object category. 2613 2. Multiple Semantics: Select the meaning of the most relevant and likely object category label in an image. For instance, as an object category label, "bank" can be "a financial institution" or "a landform alongside a 2614 river", and the former is more likely to be the meaning in an image. 2615 3. Distinctive Features: Emphasize unique aspects differentiating the object from other similar objects. 2616 4. Relevant Knowledge: Incorporate important concepts, historical, or cultural information that enrich the 2617 understanding of the object but avoid excessive details. 2618 5. You can optionally focus on the following aspects when generating category descriptions: 1. What are useful features for distinguishing the object of the given category label in an image? 2619 2. What does the object of the given category label in an image look like? 2620 3. What are the identifying characteristics of the object with the given category label in an image? 2621 4. What are the key visual indicators that help identify the object category label in an image? 2622 2623 Here are some examples: 2624 ... 2625 Object Category Label: malayan tapir 2626 2627 Category Description: The Malayan tapir, a distinctive mammal found predominantly in Southeast Asia, 2628 is known for its unique coloration: It boasts a unique appearance, featuring a black body with a white 2629 front section. It resembles a large pig in shape, with a short, prehensile trunk, small eyes, and pointed ears. Primarily nocturnal, tapirs are herbivorous and thrive in dense forests near water sources. 2630 2631 2632 2633 2634 User Prompt # Object Category Description Generation 2635 2636 Please directly generate an informative, concise, accurate and clear description for the given object category 2637 label in about 50 words 2638 2639 2640 Object Category Label: {category_label} 2641 Category Description: 2642 2643 2644 2645

Table 21: Prompts for the attribute label description generation for Visual Genome. The examples 2647 are omitted. 2648 2649 System Prompt 2650 # Instructions for Attribute Category Description Generation 2651 2652 You are a helpful, respectful, and honest assistant. 2653 Now you are an expert in generating descriptions for attribute category labels of objects in an image. Your goal is to generate a description of the given attribute category label to help people better understand it 2654 from both vision and language modalities. 2655 Informative, concise, accurate and clear descriptions are expected. 2656 2657 You are given some examples, each example is in the following format: 2658 ... 2659 2660 Attribute Category Label: the attribute category label required to generate a description. 2661 Related Object Category Labels: the category labels of the objects which may be associated with the attribute category label, separated by commas. 2663 Attribute Category Description: the generated description of the attribute category label. 2665 Attention, the given attribute category label and related object category labels may be noisy. Specifically: 2667 2668 1. The given related object category labels may be invalid or not be associated with the given attribute 2669 category label. You should disregard such noisy related object category labels. 2670 2. The given attribute category label may be invalid. You should generate "Invalid." as its description and then explain why it is invalid after "Invalid.". 2671 2672 Here are some useful tips for generating attribute category descriptions when the given attribute category 2673 label is valid: 2674 2675 1. Universality: Focus on features common to most instances of the attribute category. 2. Multiple Semantics: Select the meaning of the most relevant and likely attribute category label in an 2676 image. For instance, as an attribute category label, "awake" is an adjective meaning "not asleep" rather than 2677 a verb meaning "to stop sleeping or wake up from sleep". 2678 3. Distinctive Features: Emphasize unique aspects differentiating the attribute from other similar attributes. 2679 4. Appearance, Sensory and Effect: Describe how the attribute typically appears or is perceived on the objects. This could include color, texture, size, sense, function, or any other aspect that is significant for the 2680 attribute. 2681 5. You can optionally focus on the following aspects when generating attribute category descriptions: 1. How is the attribute category label usually observed or sensed? 2683 2. What impact does this attribute have on the object, its perception or function? 2684 3. What are the identifying characteristics of the given category label of an object in an image? 2685 4. What are the key visual indicators that help identify the attribute category label in an image? 2686 Here are some examples: 2687 2688 {examples} 2689 2690 User Prompt 2691 # Attribute Category Description Generation 2692 Please directly generate an informative, concise, accurate and clear description for the given attribute category 2694 label in about 50 words. 2695 ... 2696 Attribute Category Label: {attribute category_label} 2697 2698 Related Object Category Labels: {related object category labels}, ...

Attribute Category Description:

2700 2701 Table 22: Prompts for the relationship label description generation for Visual Genome. The examples 2702 are all omitted.

2703	
2704	System Prompt
2705	# Instructions for Relationship Category Description Generation
2706	
2707	You are a helpful, respectful, and honest assistant.
708	Your goal is to generate a description of a given relationship label to belp people better understand it from
700	both vision and language modalities.
0710	Informative, concise, accurate and clear descriptions are expected.
0711	
0710	You are given some examples, each example is in the following format:
712	"
717	Relationship Category Label: the relationship category label required to generate a description.
715	relationship category Zacon die rolationship category lacor required to generate a description
716	Related Subject-Object Pairs: the subject-object category label pairs associated with the relationship category
710	label. Each pair in the form of [Subject Category Label, Object Category Label], different pairs are separated
.717	by commas.
210	Relationship Category Description: the generated description of the relationship category label
2719	w
2720	
2/21	Attention, the given relationship category label and related subject-object pairs may be noisy. Specifically:
2722	1. The given related subject-object pairs may be invalid or not be associated with the given relationship
2723	category label. You should disregard such noisy related subject-object pairs.
724	2. The given relationship category label may be invalid. You should generate "Invalid." as its description and
725	then explain why it is invalid after invalid.
2726	Here are some useful tips for generating attribute category descriptions when the given relationship category
727	label is valid:
728	
729	1. Universality: Focus on features common to most instances of the relationship category.
2730	2. Multiple Semantics: Select the meaning of the most relevant and likely relationship category label in an image. For instance, as a relationship category label. "truck" is a varb meaning "convex by truck" rather than
2731	a noun meaning "a large, heavy road vehicle used for carrying goods and materials".
2732	3. Distinctive Features: Emphasize unique aspects differentiating the relationship from other similar
2733	relationships.
2734	4. Nature of Relationship: Describe the nature of the relationship (e.g., spatial, action, functional, hierarchical,
2735	temporal, social) and how the subject and object interact or relate to each other.
736	1. How is the relationship expressed or manifested in the image?
2737	2. What are the key characteristics or significance of the relationship?
738	3. What are common or typical scenarios in which this relationship is observed?
739	4. What are the key visual indicators that help identify the relationship category in an image?
740	
741	Here are some examples:
742	{examples}
2743	
744	User Prompt
745	# Relationship Category Description Generation
746	
2747	Please directly generate an informative, concise, accurate and clear description for the given relationship
748	category label in about 50 words.
0749	
2750	Relationship Category Label: {relationship category label}
751	Terminiship Curegory Europh (Terminiship Curegory_10001)
750	Related Subject-Object Pairs: {related subject-object pairs},
.132	
2753	Relationship Category Description:

2755 2756 2757 Table 23: Image-first template. The instructions for the related region part will be repeated if the 2758 image has more than one region. 2759 2760 image-first template with descriptions 2761 2762 # Detailed Analysis of Objects in the Image 2763 2764 Image: [IMG] 2765 Caption: [image caption] Localized narrative caption: [localized narrative caption] 2766 Objects and their descriptions: 2767 - [object label]: [object label des] 2768 - [object label]: [object label des] 2769 Attributes of objects and their descriptions: 2770 - [attribute label] [object label], [object label] : [attribute label des] 2771 - [attribute label] [object label], [object label] : [attribute label des] Relationships between objects and their descriptions: 2772 - [relationship label] [subject label]-[object label], [subject label]-[object label] : [relationship label des] 2773 - [relationship label] [subject label]-[object label], [subject label]-[object label] : [relationship label des] 2774 ## Overview of Selected Object Regions 2775 2776 ### Overview of a Selected Object Region 2777 Region: [IMG] 2778 Location of the selected region in the image: [location] 2779 Objects: 2780 - [object label] 2781 - [object label] 2782 2783 image-first template without descriptions 2784 # Detailed Analysis of Objects in the Image 2785 2786 Image: [IMG] 2787 Caption: [image caption] 2788 Localized narrative caption: [localized narrative caption] 2789 Objects: - [object label] 2790 - [object label] 2791 Attributes of objects: 2792 - [attribute label] [object label], [object label] 2793 - [attribute label] [object label], [object label] 2794 Relationships between objects: 2795 - [relationship label] [subject label]-[object label], [subject label]-[object label] - [relationship label] [subject label]-[object label], [subject label]-[object label] 2796 ## Overview of Selected Object Regions 2797 2798 ### Overview of a Selected Object Region 2799 Region: [IMG] 2801 Location of the selected region in the image: [location] Objects: 2802 - [object label] 2803 - [object label]

2805 2806 2807

2754

	ore 2.1. Test mot complute, the motivations for the related region part will be repeate
m	age has more than one region.
	text-first template with descriptions
	# Detailed Analysis of Objects in the Image
	Objects and their descriptions:
	- [object label]: [object label des]
	- [object label]: [object label des]
	Attributes of objects and their descriptions:
	- [attribute label] [object label], [object label] : [attribute label des]
	- [attribute label] [object label], [object label] : [attribute label des]
	Relationships between objects and their descriptions:
	- [relationship label] [subject label]-[object label], [subject label]-[object label] : [relationship label des
	- [relationship label] [subject label]-[object label], [subject label]-[object label] : [relationship label des
	Localized narrative caption: [localized narrative caption]
	Image: [TMG]
	## Overview of Selected Object Regions
	### Overview of a Selected Object Region
	Location of the selected region in the image: [location]
	Ubjects:
	- [ODJect label]
	- [OJECT TADE] Region: [IMG]
	text-first template without descriptions
	# Detailed Analysis of Objects in the Image
	Objects and their descriptions:
	- [object label]
	- [ODJect label]
	- [attribute label] [object label] [object label]
	- [attribute label] [object label], [object label]
	Relationships between objects:
	- [relationship label] [subject label]-[object label], [subject label]-[object label]
	- [relationship label] [subject label]-[object label], [subject label]-[object label]
	Caption: [image caption]
	Localized narrative caption: [localized narrative caption]
	Image: [IMG]
	## Overview of Selected Object Regions
	### Overview of a Selected Object Region
	Location of the selected region in the image: [location]
	Objects:
	- [object label]
	- [object label]
	Region: [IMG]