

Visual Prompting in Multimodal Large Language Models: A Survey

Junda Wu¹ Zhehao Zhang² Yu Xia¹ Xintong Li¹ Zhaoyang Xia³ Aaron Chang⁴
Tong Yu⁵ Sungchul Kim⁵ Ryan A. Rossi⁵ Ruiyi Zhang⁵ Subrata Mitra⁵
Dimitris N. Metaxas³ Lina Yao^{6,7} Jingbo Shang¹ Julian McAuley¹

¹UC San Diego ²Dartmouth College ³Rutgers University ⁴UC Los Angeles

⁵Adobe Research ⁶The University of New South Wales ⁷CSIRO's Data61

{juw069, yux078, xil240, jshang, jmcauley}@ucsd.edu zhehao.zhang.gr@dartmouth.edu

zx149@rutgers.edu aaronchang21@ucla.edu dnm@cs.rutgers.edu

{tyu, sukim, ryrossi, ruizhang, sumitra}@adobe.com lina.yao@data61.csiro.au

Abstract

Multimodal large language models (MLLMs) equip pre-trained large-language models (LLMs) with visual capabilities. While textual prompting in LLMs has been widely studied, visual prompting has emerged for more fine-grained and free-form visual instructions. This paper presents the first comprehensive survey on visual prompting methods in MLLMs, focusing on visual prompting, prompt generation, compositional reasoning, and prompt learning. We categorize existing visual prompts and discuss generative methods for automatic prompt annotations on the images. We also examine visual prompting methods that enable better alignment between visual encoders and backbone LLMs, concerning MLLM's visual grounding, object referring, and compositional reasoning abilities. In addition, we provide a summary of model training and in-context learning methods to improve MLLM's perception and understanding of visual prompts. This paper examines visual prompting methods developed in MLLMs and provides a vision of the future of these methods.

1 Introduction

Multimodal large language models (MLLMs) (Li et al., 2023b; Liu et al., 2024a), which augment pre-trained large language models (LLMs) with visual capabilities, enable visual understanding and reasoning on complex multimodal tasks (Zhou et al., 2024b; Jia et al., 2024). However, limited by using textual prompts to describe and specify visual elements (Lin et al., 2024a; Wu et al., 2024d), conventional prompting methods fall short of providing accurate visual grounding and referring to detailed visual information, which can cause visual hallucinations (Bai et al., 2024; Huang et al., 2024b) and language bias (Wu et al., 2024b; Qu et al., 2024).

Recently, visual prompting methods have emerged (Zhang et al., 2024a; Wu et al., 2024f)

as a new paradigm, complementing textual prompting and enabling more fine-grained and pixel-level instructions on multimodal input. Since visual prompting methods can take heterogeneous forms for various tasks and often operate at pixel-level granularity, general prompt templates might not apply to different images, making instance-level visual prompt generation necessary. Therefore, we provide a comprehensive categorization of current visual prompting methods (Section 2) and methods to generate (Section 3) such visual prompts.

Despite the success of visual prompting methods in augmenting MLLM's visual abilities, several works also suggest that MLLMs can be misaligned with visual prompts, due to the lack of heterogeneous visual prompting training data during the pre-training stage (Yan et al., 2024; Lin et al., 2024b). This misalignment can cause MLLMs to neglect or misinterpret certain visual prompts, causing hallucination problems. Therefore, we summarize existing efforts in aligning visual prompting with MLLM's perception and reasoning enabling more controllable compositional reasoning (Section 5). In addition, we examine existing pre-training, fine-tuning (Section 6), and in-context learning methods (Section 7) that fundamentally align MLLMs with multimodal augmented prompts.

Existing surveys on LLM prompting are limited to textual prompt design (Gu et al., 2023; Schulhoff et al., 2024) and in-context demonstrations (Xu et al., 2024b; Li, 2023), which lack literature coverage of pixel-level instructions and multimodal interactions. Visual prompting is also studied in computer vision. However, relevant surveys are limited to vision tasks with vision backbone models (Lei et al., 2024b; Zhang et al., 2024b), while multimodal perception and reasoning tasks involving MLLMs are absent. In addition, one recent survey on Segment Anything Models (SAM) (Zhang et al., 2023a) explores various applications of SAM in MLLMs. However, this work is limited to the

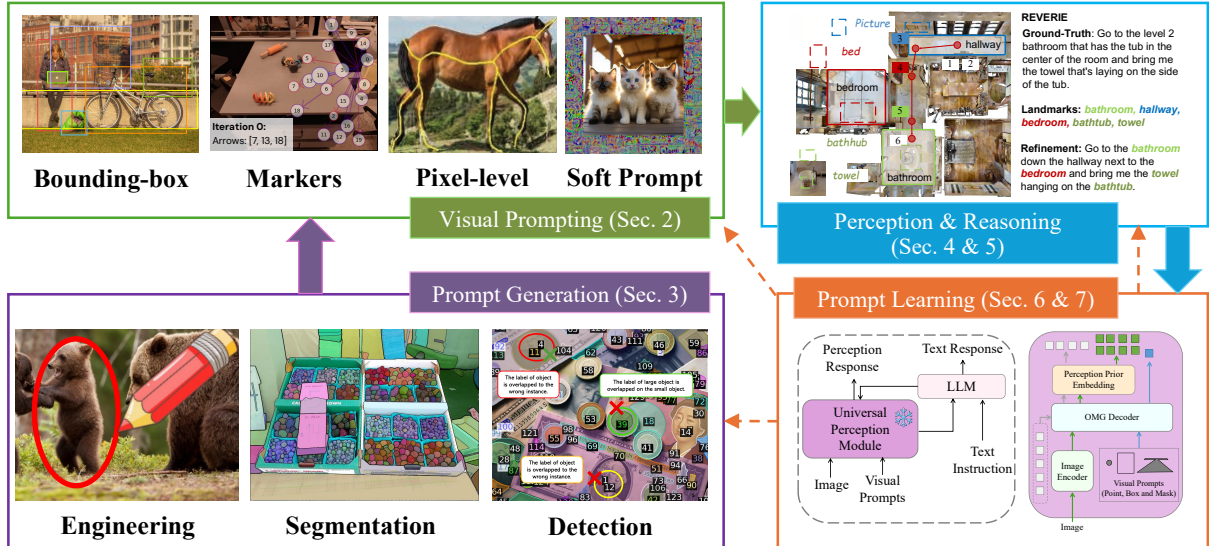


Figure 1: Taxonomy flow chart of visual prompting techniques. We illustrate in order of four stages of visual prompting including prompt generation, visual prompting, perception and reasoning, and prompt learning, where the solid arrows show the direction of each component’s information flow. We explain in detail various visual prompt generation techniques (Section 3), and how these generated prompts are used to prompt MLLMs (Section 2). Then we discuss the advanced perception and reasoning abilities achieved through visual prompting (Section 4 and 5). Finally, model pre-training, fine-tuning, instruction tuning, and in-context learning further update previous model components, which are illustrated by the dashed arrows (Section 6 and 7).

SAM model and lacks comprehensive studies on diverse visual prompting methods. In this paper, we present the first comprehensive survey on visual prompting in MLLMs to address these gaps and extend the current understanding of visual prompt generation, multimodal prompting, perception and reasoning, and prompt learning. We illustrate the taxonomy of our survey in Figure 1 and summarize our contributions as follows:

- We provide a comprehensive categorization of visual prompting and prompt generation methods in MLLMs.
- We explain the integration of visual prompts into MLLM’s perception and reasoning for more controllable compositional reasoning, which helps to prevent hallucination and language bias issues.
- We summarize MLLM alignment methods with visual prompts, covering model training and in-context learning, addressing issues of misinterpretation, and proposing strategies for more controllable compositional reasoning.

2 Visual Prompt Categorization

Visual prompts are essential tools in MLLMs, guiding models in interpreting and processing visual data. These prompts (Wu et al., 2024f) can take

various forms, such as bounding boxes, markers, pixel-level prompts, and soft prompts. They provide additional information to enhance the models’ visual perception capabilities. By manipulating images and videos with different techniques, visual prompts improve model performance in complex understanding and reasoning tasks.

2.1 Bounding-Box

Bounding boxes are used to demarcate objects or regions within an image, enabling MLLMs to extract visual features (Lin et al., 2024a). These features help the model understand the image content and correlate it with the corresponding text, thereby enhancing fine-grained and grounded image understanding. Previous work, such as Shikra (Chen et al., 2023b) and VTPrompt (Jiang et al., 2024), quantize bounding boxes to represent key objects numerically, modeling both input and output positions. Other approaches modify bounding boxes for specific tasks: A3VLM (Huang et al., 2024a) uses 3D bounding boxes to locate actionable parts of an image, CityLLaVA (Duan et al., 2024) scales up the bounding box, and TextCoT (Luan et al., 2024) extends the shorter sides of the bounding box to match the longer side, ensuring it encompasses the entire region of interest. In addition, CRG (Wan et al., 2024) masks out specific regions with black pixels to reduce priors, providing a way to correct predic-

tions without additional training. Groma (Ma et al., 2024a) and InstructDET (Dang et al., 2023) encode user-specified regions (i.e., bounding boxes) into visual tokens, enhancing the localization ability of MLLMs by directly integrating them into user instructions. Another framework (Lin et al., 2024b) further enhances the localization capabilities of MLLMs by integrating contextual embeddings from external knowledge within bounding boxes, serving as visual prompts to boost the fine-grained cognitive abilities of various MLLMs.

2.2 Markers

Similar to bounding boxes, visual markers are specific elements within visual data (such as images or videos) used to highlight, identify, or draw attention to particular features or regions. They are often employed to indicate particular parts of an image that are relevant to the task. Prior work (Shtedritski et al., 2023) has demonstrated that models trained on web-scale data can focus on specific visual markers, such as red circles, to highlight desired regions instead of cropping the image around them. AutoAD-Zero(Xie et al., 2024) introduced a two-stage, training-free approach that incorporates character information by "circling" characters in the frame and color-coding each identity. More recently, Set-of-Mark (SoM) prompting(Yang et al., 2023) overlays visual markers directly onto images to help models generate answers grounded in specific image regions. ViP-LLaVA(Cai et al., 2024) expands on this by incorporating arbitrary visual cues like scribbles and arrows, using fine-tuned models to recognize these markers. Liao et al. (2024) also leverage the SoM technique to introduce feedback, converting it into text or visual marks to improve semantic grounding. SoM-LLaVA (Yan et al., 2024) proposes a method to enhance SoM’s tag association by listing items individually and comprehensively describing all tagged items within an image. Other methods, such as ToL (Fan et al., 2024b) and OWG (Tziafas and Kasaei, 2024), link each segment in the frame with a unique ID, while Pivot (Nasiriany et al., 2024) projects a 3D location into image space and draws a visual marker at this projected location to refer to spatial concepts in the output space.

2.3 Pixel-level

Previous approaches relied on coarse markers like colorful boxes or circles, which introduced ambiguity in accurately highlighting objects. To address

this issue, pixel-level prompts (Ma et al., 2024b) use individual pixels in images or videos, enhancing the semantic localization capability of MLLMs. Methods such as FGVP (Yang et al., 2024a), EVP (Liu et al., 2023b), DOrA (Wu et al., 2024e), and CoLLaVO (Lee et al., 2024) employ pixel-level prompts to convey semantic information for precise object localization. OMG-LLaVA (Zhang et al., 2024e) and VisionLLM (Wang et al., 2024b) tokenize images into pixel-centric visual tokens, aligning visual tasks with language instructions. Techniques such as Image Inpainting (Bar et al., 2022) decode visual tokens into pixels, while ControlMLLM (Wu et al., 2024d) models rich semantic relations between pixels and text prompts. There are also coordinate prompt methods, such as SCAF-FOLD (Lei et al., 2024a) and AO-Planner (Chen et al., 2024a), which convert input images into coordinates using metrics, enhancing spatial understanding and reasoning abilities in MLLMs.

2.4 Soft Visual Prompt

Soft visual prompts, learned in the pixel space and applied directly to the image, allow models to adapt more effectively to specific downstream tasks. In particular, TVP (Zhang et al., 2024g), BlackVIP (Oh et al., 2023), and VPGTrans (Zhang et al., 2024a) add pixel-level prompts to images, either by surrounding the image with universal prompts or designing prompts matching the image’s shape. In Learned Prompt (Rezaei et al., 2024), WVPrompt (Ren et al., 2024), and ILM-VP (Chen et al., 2023a), task-relevant perturbation patterns are injected into the pixel space to modify the input sample. Additionally, ImageBrush (Yang et al., 2024b) enhances semantic understanding by extracting tokenized features from images.

3 Visual Prompt Generation

Different from textual prompts, visual prompts are typically position-aware and instance-specific, involving particular visual objects, relationships, and contexts. Current approaches use visual prompt generation methods and models to improve the accuracy and comprehension of visual prompts by MLLMs, which generate visual prompts, such as segmentation, detection, and image inpainting, for individual images and videos. Additionally, toolchains of visual prompt methods are employed to enable multi-step visual reasoning and planning. To create universally applicable visual prompts

learnable pixel values have also been developed.

3.1 Prompt Engineering

Understanding human-engineered visual prompts can be important in practical use cases, where visual prompts are especially efficient for expressing one’s intention or attention to the current visual evidence. Early exploration (Shtedritski et al., 2023) discovers that drawing a simple red circle around an object can direct a model’s attention to that region. In addition, to enrich detailed visual evidence, MIVPG (Zhong et al., 2024) leverages instance correlations within images or patches.

ViP (Cai et al., 2024) introduces a novel multi-modal model capable of decoding free-form visual prompts, allowing users to intuitively mark images with natural cues. This approach does not require complex region encodings and achieves state-of-the-art performance on region-specific comprehension tasks. In addition, ViP-Bench (Cai et al., 2024) is also proposed to evaluate MLLM’s perception of such naturally engineered visual prompts. In domain-specific CityLLaVA (Duan et al., 2024) framework, engineered visual prompts are collected and tailored for urban scenarios, which further augments the fine-tuned MLLM.

3.2 Visual Segmentation

Segmentation methods such as OpenSeeD (Zhang et al., 2023b), SAM (Kirillov et al., 2023), and SegFormer (Xie et al., 2021), are used to delineate and identify specific regions, objects, or structures within images, thus enabling the models to focus on relevant visual information more accurately. With pre-trained segmentation models, external visual knowledge can be transferred and integrated into MLLM’s prompt. Yang et al. (2024a) explore a fine-grained visual prompting method by pixel-level annotations annotated from image inpainting (Bar et al., 2022) method. Lin et al. (2024b) propose an instruction tuning method to directly incorporate fine-grained segmentation knowledge in the spatial embedding map as visual prompts, which enhances the model’s context-awareness of the visual scene. VAP (Chen et al., 2024a) develops a visual affordance prompting method that grounds visual elements by SAM (Kirillov et al., 2023) in navigation tasks. DOrA (Wu et al., 2024e) further introduces 3D spatial and contextual information to improve 3D visual grounding tasks.

Fine-grained segmentation information also augments MLLM’s visual perception and reasoning

abilities. OMG-LLaVA (Zhang et al., 2024e) integrates multi-level visual prompts that enable MLLM’s coarse-to-fine visual perception to more comprehensive visual understanding. Liu et al. (2023b) propose to enhance the model’s ability to understand and process low-level structural elements within images. He et al. (2024) further incorporate such visual prompts into MLLM fine-tuning to augment the model’s capacity in fine-grained visual perception. CoLLaVO (Lee et al., 2024) proposes a crayon prompting which further augments with panoptic segmentation method with image in-painting color maps to better discriminate multi-objects within the image.

3.3 Object Detection

Object detection models like SoM (Yang et al., 2023), RCNN (Girshick, 2015), and Omni3D (Brazil et al., 2023) provide precise object identification and localization in the visual context, which assists MLLM’s visual grounding abilities and guides MLLM’s attention on semantically meaningful contents. SoM-LLaVA developed by Yan et al. (2024) uses numeric tags to align visual objects with textual descriptions. Object tags enable the model to list and describe these objects accurately, which enhances visual reasoning and visual instruction following capabilities. InstructDET (Dang et al., 2023) incorporates generalized instructions into the training process, diversifying object detection by enabling the model to understand and follow various referring instructions. This enhances the model’s flexibility in understanding user intentions and instructions in different task contexts. Wan et al. (2024) propose to improve the grounding of vision-language models by contrastive region guidance. By guiding the model’s attention to relevant regions, MLLM can more accurately associate visual regions with corresponding textual instructions. Cho et al. (2024) extend vision-language models to understand 3D environments, by improving spatial awareness and the understanding of object interactions in three-dimensional spaces.

3.4 Visual Prompt Toolchain

To enable more complex multimodal understanding by multi-step or interactive reasoning, several methods aggregate various visual prompting methods as toolchains (Wu et al., 2024f) to be called by the MLLM and assist individual reasoning sub-tasks. Zhou et al. (2024b) propose an image-of-thought method that can automatically determine

each reasoning step’s visual information extraction method and implement it as visual prompts, which prompt MLLM to follow a certain reasoning path and enable step-by-step multimodal reasoning. Tzifafas and Kasaei (2024) focus on adapting vision-language models for open-world grasping tasks by incorporating a list of visual prompting methods including open-end segmentation and object grounding to enable open-world grasping tasks. To enable more transferable and generalizable visual prompts, Sheng et al. (2024) create a more unified in-context learning method by integrating various contextual visual prompts into a unified representational space. MineDreamer (Zhou et al., 2024a) further develops a versatile visual prompt generation method for imaginary visual scenes, which are consistent with current decision-making intention and visually express the next-step goal.

3.5 Learnable and Soft Visual Prompt

Learnable or soft visual prompts are employed to adapt the visual encoder in MLLMs, enabling more controlled and versatile use of visual prompts that are aligned with downstream tasks. Such techniques are used in multimodal instruction tuning with visual instructions. Rezaei et al. (2024) investigates how visual prompts can be learned to guide the attention mechanisms in ViT. Li et al. (2023a) fine-tune MLLMs to follow zero-shot demonstrative instructions using learnable visual prompts. Chen et al. (2023a) focus on better mapping visual inputs to corresponding labels through learned prompts. For some specific and domain-oriented problems, (Ren et al., 2024) develop a learnable visual prompting method as image watermarking to identify the image’s copyright and ownership.

At the same time, learnable visual prompts can also be transferable across MLLMs and downstream tasks. VPGTrans (Zhang et al., 2024a) proposes a transferable visual prompt generator, which adapts the pre-trained source MLLM to target MLLM with low cost in training data points and computation. Memory-space visual prompt (Jie et al., 2024) injects learnable prompts in the key and value layers in the vision-transformer architecture, which enables efficient vision-language fine-tuning. Wu et al. (2023) also injects soft visual tokens as visual compositional operations, which are learned to better compose multimodal information with few-shot examples. The black-box visual prompting method (Oh et al., 2023) focuses on robust transfer learning, where the visual prompts

help models adapt to new tasks and domains without direct access to model parameters.

4 Visual Perception

4.1 Visual Grounding and Referring

Recent visual prompting works have significantly improved MLLM’s visual grounding and referring abilities. Some works emphasize the importance of iterative feedback and multimodal interaction to refine semantic grounding, while others explore object-centric perception and the comprehension of visual relations. To improve MLLM’s regional grounding and object detection abilities, SoM-LLaVA (Yan et al., 2024) employs the Set-of-Mark (SoM) model to tag all the objects in the image and ask the model to list all the items. Instruct-DET (Dang et al., 2023) and VTPrompt (Jiang et al., 2024) further enable multimodal grounding, which extracts object entities from the text and these objects’ regional bounding boxes.

With a fine-grained visual grounding encoder, several works further use visual cues to guide MLLM’s attention to relevant regions within images and achieve better regional referring abilities. CRG (Wan et al., 2024) uses contrastive regional guidance to direct the model’s attention to specific areas of interest within an image, without model finetuning. RelationVLM (Huang et al., 2024c) leverages visual prompts to enhance MLLM’s understanding and reasoning about objects’ spatial relations. Shikra (Chen et al., 2023b) further applies to visual dialogue systems, where MLLM responds to referential cues within a dialogue, enabling more precise and context-aware interactions. In addition, several works aim to provide a comprehensive framework that incorporates various visual prompting methods in different granularity levels, to enable more fine-grained and flexible multimodal interactions, including free-form visual prompt inputs (Lin et al., 2024a) and feedback mechanisms (Liao et al., 2024) on visual prompts.

4.2 Multi-image and Video Understanding

To improve the models’ understanding of complex visual relationships and ensure that they can accurately reference and describe objects across diverse multi-image inputs, several works propose visual prompts in multi-image inputs and novel evaluation benchmarks to test their effectiveness. Fan et al. (2024c) present a novel benchmark dataset with multipanel images to test MLLM’s abilities in

distinguishing objects across panels and navigating between different visual elements. Pan et al. (2024) leverage morph-token auto-encoding to enhance the model’s capacity to process visual grounding across multiple images. Li et al. (2023a) fine-tune MLLMs to follow in-context demonstrative instructions across multiple images. In addition, AIM (Gao et al., 2024) proposes to dynamically adapt its grounding and referring abilities to accommodate new visual contexts across several images.

Several methods are also developed to allow MLLMs to identify specific regions of interest, improving their ability to handle complex and dynamic video content. OmAgent (Zhang et al., 2024c) develops a visual prompting method to enable task division in video understanding, by annotating a series of visual features. RACCooN (Yoon et al., 2024) uses visual prompts to guide MLLMs in identifying the target regions in the video for manipulation. Wu et al. (2024c) ground objects across videos, enabling the model to comprehend and refer to objects in dynamic scenes.

4.3 3D Visual Understanding

Recent works use visual prompting for better 3D visual understanding. Li et al. (2024) constructs an extensive dataset comprising instruction-responses pairs for 3D scenes and introduces 3DMIT for efficient prompt tuning while eliminating the alignment stage between 3D scenes and languages. DOrA (Wu et al., 2024e) proposes a novel 3D visual grounding framework with Order-Aware referring. This method leverages LLM to infer ordered object series that used to guide the progressive feature refinement process.

Cho et al. (2024) constructs a large-scale dataset named LV3D and introduces a new MLLM CubeLLM pre-trained on the proposed dataset. Zhang et al. (2024d) proposes Agent3D-Zero, which introduces novel visual prompts by employing bird’s-eye view images and selecting viewpoints to unleash the MLLM’s ability to observe 3D scenes. 3DAP (Liu et al., 2023a) develops a novel visual prompting method that creates a 3D coordinate system and additional annotation to empower GPT-4V to complete 3D spatial tasks.

5 Compositional Reasoning

This section discusses how visual prompting enhances compositional and multimodal learning in MLLMs, enabling improvements in tasks like vi-

sual planning, reasoning, and action generation. We examine how visual prompts facilitate complex step-by-step reasoning, decision-making, and control over visual generation models, expanding their capabilities across diverse tasks. We also review several frontier applications (Appendix 9), which can be under-explored and lack sufficient solutions.

5.1 Visual Planning

Recent works demonstrate that visual prompting improves visual planning tasks. Zhou et al. (2024b) proposes an Image-of-Thought(IoT) prompting method that compels MLLMs to automatically design visual and textual steps and leverages external image processing tools to generate a multi-model rationale series, which is used to assist MLLMs with complex visual reasoning tasks through a step-by-step process. OWG (Tziafas and Kasaei, 2024) combines segmentation and grasp synthesis models, which unlocks the grounded world understanding through segmentation, grasp planning, and ranking. Zhou et al. (2024a) introduces the Chain-of-Imagination (CoI) method and creates an embodied agent in Minecraft named MineDreamer. This method envisions the step-by-step process of executing instructions with the help of an LLM-enhanced diffusion model that translates imaginations into precise visual prompts to support the accurate generation of the agent’s actions. BEVInstructor (Fan et al., 2024a) incorporates Bird’s Eye View representations as visual prompts into MLLMs for navigation instruction generation. AO-Planner (Chen et al., 2024a) achieves affordances-oriented motion planning and action decision-making with a VAP approach and a high-level PathAgent.

5.2 Chain-of-thought

To enable more complex image reasoning, recent works incorporate visual prompting with Chain-of-Thought methods. Luan et al. (2024) proposes a novel Chain-of-Thought framework for text-rich image understanding, named TextCoT. This method consists three stages including image overview for global information, coarse localization for estimating the section that encompasses the answer and fine-grained observation for furnishing precise answers. Wu et al. (2024f) proposes DetToolChain to unlock the potential of MLLMs in object detection task. This method involves using a "detection prompting toolkit," which includes visual processing and detection reasoning prompts,

combined with a multimodal detection Chain-of-Thought method to reason the sequential implementation of the detection prompts.

6 Model Training

This section presents key approaches to align multimodal large language models (MLLMs) using visual prompting techniques, including pre-training, fine-tuning, and instruction tuning, which aim to unify multi-modal prompts and improve cross-task transferability. In addition to model training techniques, we also summarize evaluation datasets (Appendix 8), which inspire future work to develop more powerful visual prompting methods.

6.1 Pre-training

To improve MLLM’s ability on more fine-grained vision perception or reasoning tasks, a line of works focuses on designing better pre-training objectives including visual prompts. PSALM (Zhang et al., 2024h) extends the capabilities of MLLM to address various image segmentation tasks by incorporating a mask decoder and a flexible input schema. This approach unifies multiple segmentation tasks within a single model, supporting generic, referring, interactive, and open-vocabulary segmentation, while demonstrating strong performance on both in-domain and out-of-domain pixel-level segmentation tasks. OMG-LLaVA (Zhang et al., 2024e) proposes a unified framework that bridges image-level, object-level, and pixel-level reasoning and understanding in a single model that combines a universal segmentation method as the visual encoder with an LLM, enabling flexible user interaction through various visual and text prompts. VisionLLM v2 (Wu et al., 2024a) introduces an end-to-end generalist MLLM that unifies visual perception, understanding, and generation within a single framework. The model employs a novel "super link" technique to connect the central LLM with task-specific decoders, enabling flexible information transmission and end-to-end optimization across hundreds of vision and vision-language tasks. UrbanVLP (Hao et al., 2024) proposes a vision-language pretraining framework for urban region profiling that integrates multi-granularity information from both satellite (macro-level) and street-view (micro-level) imagery, overcoming previous limitations. This method also incorporates an automatic text generation and calibration mechanism to produce high-quality textual descriptions

of urban areas, enhancing interpretability.

6.2 Fine-tuning

Zhang et al. (2024g) propose Transferable Visual Prompting (TVP), a method to improve the transferability of soft visual prompts which are a small amount of learnable parameters across different MLLMs for downstream tasks. Lin et al. (2024b) integrate fine-grained external knowledge such as OCR and segmentation into multimodal MLLMs through visual prompts, which embed fine-grained knowledge information directly into a spatial embedding map. CoLLaVO (Lee et al., 2024) enhances MLLMs’ object-level image understanding through a visual prompt called Crayon Prompt, which is derived from panoptic color maps generated by a panoptic segmentation model. CityLLaVA (Duan et al., 2024) introduces an efficient fine-tuning framework for MLLM designed for urban scenarios which incorporates visual prompt engineering techniques, including bounding box-guided, view selection, and global-local joint views. ViP-LLaVA (Cai et al., 2024) is enabled to understand arbitrary visual prompts, which is trained by directly overlaying visual markers onto images. ImageBrush (Yang et al., 2024b) introduces a framework for exemplar-based image manipulation that learns visual in-context instructions without language prompts.

Explicit Visual Prompting (EVP) (Liu et al., 2023b) proposes a unified approach for low-level structure segmentation tasks with a frozen pre-trained vision transformer backbone and introduces task-specific soft prompts derived from frozen patch embeddings and high-frequency image components. BlackVIP (Oh et al., 2023) adapts large pre-trained models with a Coordinator to generate soft visual prompts and SPSA-GC for efficient gradient estimation, enabling robust few-shot adaptation across diverse domains. Iterative Label Mapping-based Visual Prompting (ILM-VP) (Chen et al., 2023a) improves the accuracy and interpretability of soft visual prompting by jointly optimizing input patterns and label mapping through bi-level optimization. MemVP (Jie et al., 2024) efficiently combines pre-trained vision encoders and language models for vision-language tasks by injecting visual information directly into the feed-forward network weights of MLLMs, treating them as additional factual knowledge. VPG-C (Li et al., 2023a) enhances visual prompting in MLLMs by completing missing visual details to better compre-

hend demonstrative instructions with interleaved multimodal context. It extends traditional Visual Prompt Generators by using LLM-guided, context-aware visual feature extraction to create more comprehensive visual prompts.

6.3 Instruction Tuning

Instruction tuning has proved to effectively improve the overall ability of both text-only LLMs and MLLMs such as instruction following and structured output (Ouyang et al., 2022; Wang et al., 2022; Liu et al., 2024a). For MLLMs with a focus on visual prompts, AnyRef (He et al., 2024) introduces a unified referring representation that enables the MLLM to handle diverse input modalities and visual prompts (text, bounding boxes, images, audio) through instruction tuning. This model uses special tokens and prompts to format multi-modal inputs, allowing it to process various referring formats consistently. A refocusing mechanism enhances mask embeddings by incorporating grounded textual embeddings, improving segmentation accuracy. AnyRef combines vision and audio encoders with an LLM, using projection layers to align different modalities in the language space. The model is instruction-tuned end-to-end with a combination of text loss and mask loss, enabling it to generate both textual descriptions and pixel-level segmentation in response to multi-modal prompts.

7 In-context and Few-shot Learning

Beyond methods that optimize performance using single data points as input, some works focus on enhancing in-context learning (ICL) with visual prompts. Image-of-Thought (IoT) prompting (Zhou et al., 2024b) is a train-free approach to enhance MLLMs for visual question-answering tasks by integrating discrete image processing actions. IoT enables MLLMs to automatically design and extract visual rationales step-by-step, combining them with textual rationales to improve both accuracy and interpretability. CRG (Wan et al., 2024) is a training-free method that improves visual grounding in MLLMs by contrasting model outputs with and without specific image regions masked which guides models to focus on relevant image areas. AIM (Gao et al., 2024) enables any MLLM to perform efficient ICL by aggregating image information from demonstrations into the latent space of corresponding textual labels which reduces memory costs by discarding visual tokens after aggrega-

tion, approximating multimodal ICL prompts to contain only a single query image. I2L(Wang et al., 2024a) combines demonstrations, visual cues, and reasoning into a single image to enhance multimodal models’ performance on complex tasks through ICL. I2L-Hybrid extends this by automatically selecting between I2L and other in-context learning methods for each task instance.

Few-shot learning through visual prompts can also improve the capabilities of MLLMs with minimum computational cost and better data efficiency. CoMM (Chen et al., 2024b) proposes a high-quality coherent interleaved image-text dataset designed to enhance the generation capabilities of MLLMs and investigate their in-context learning ability. M2oEGPT (Sheng et al., 2024) propose an ICL framework by using multimodal quantization and unified embedding to enable joint learning of multimodal data in a general token embedding space, combining an autoregressive transformer with a Mixture of Experts (MoEs) for stable multi-task co-training. Partial2Global (Xu et al., 2024a) select optimal in-context examples in visual ICL through a transformer-based list-wise ranker to compare multiple alternative samples and a consistency-aware ranking aggregator to achieve globally consistent ranking. Hossain et al. (2024) introduces learnable visual prompts for both base and novel classes on semantic segmentation, along with a novel-to-base causal attention mechanism that allows novel prompts to be contextualized by base prompts without degrading base class performance. Emu2 (Sun et al., 2024) is MLLM trained to predict the next element in diverse multimodal sequences. Its unified architecture enables strong multimodal in-context learning abilities, allowing it to quickly adapt to new tasks with just a few examples.

8 Evaluation

This section explores and compares the current MLLM visual prompting training datasets and benchmarks, as visualized in Section 7.1. The three main categories for the visual prompting techniques are Semantic Prompting (SP), Textual Prompting (TP), and GP (Generative Prompting).

The datasets and benchmarks that fall into the Semantic Prompting (SP) utilize high-level descriptions to help the model understand the semantic relationships present in the data. Some examples include creating bounding boxes (Huang et al., 2024a; Wu et al., 2024c), labeling regions of interest (Li

Reference	SP	TP	GP	Image	Video	Audio	Manual	Automatic
MDVP-Bench (Lin et al., 2024a)	✓	✓	✓	✓			✓	✓
A3VLM (Huang et al., 2024a)	✓			✓				✓
VLM Feedback (Li et al., 2023c)	✓			✓	✓		✓	✓
GPT-4V Challenger (Fu et al., 2023)	✓		✓	✓				
EarthMarker (Zhang et al., 2024f)	✓	✓	✓	✓			✓	
RACCoN (Yoon et al., 2024)			✓		✓		✓	✓
Safety of MLLMs (Liu et al., 2024b)			✓	✓	✓	✓		
GLEE (Wu et al., 2024c)	✓	✓		✓	✓			
AutoAD-Zero (Xie et al., 2024)	✓		✓	✓		✓		✓
MultipanelVQA (Fan et al., 2024c)			✓	✓				
MM-Vid (Lin et al., 2023)	✓	✓	✓	✓	✓	✓	✓	
Groma (Ma et al., 2024a)	✓	✓		✓			✓	✓

Table 1: We compare different benchmarks and training datasets, and they are each grouped into three different criteria—Semantic Prompting (SP), Textual Prompting (TP), and Generative Prompting (GP). Then, based on the different modalities, they can be classified if they contain pixel-level images (Image), video encoding and decoding (Video), and if they are supplemented by an audio transcript (Audio). Finally, the last categorization determines if the specified method visual prompting is done manually (Manual), automated (Automatic), or a combination of both.

et al., 2023c), and tagging objects (Lin et al., 2024a; Zhang et al., 2024f). Another general method is Textual Prompting (TP) where either user or LLM generated text is supplemented into the model input that relates the visual aspects in the image. Image and video descriptions can be generated and used as a visual prompt (Lin et al., 2023), drawing relationships and descriptions on the image itself in order to add location-specific analysis (Lin et al., 2024a; Wu et al., 2024c), and embedding localization into image tokenization (Ma et al., 2024a). Given the extensive effort required for manual visual prompting in MLLMs, some techniques have adopted automatic generation methods to streamline the visual prompting process using Generative Prompting (GP). Automatic modality conversion uses LLMs to generate text from images/videos and vice versa for users to easily modify and cater the prompts (Yoon et al., 2024). Audio descriptions are generated and then summarized by an LLM [(Xie et al., 2024), [(Lin et al., 2023)]]]. Similarly, generation is used to create difficult and unique benchmarks to assess the capabilities and weaknesses of specific models [(Fan et al., 2024c)].

The final taxonomy system distinguishes between those visual prompting techniques that are done manually between those that are done automatically. The manual techniques offer precision and customization, but in turn sacrifice time and efficiency. They are suitable for tasks that are smaller scale and require detail. Automatic techniques provide scalability and productivity—they work well with large scale tasks that do not require

fine-grained accuracy. Some techniques apply a combination of these (Lin et al., 2024a; Li et al., 2023c; Yoon et al., 2024; Ma et al., 2024a) and those that do not have either checked were either training datasets or surveys themselves (Liu et al., 2024b; Fan et al., 2024c; Fu et al., 2023).

9 Frontier Applications

9.1 Jailbreaking & Safety

While visual prompting enables fine-grained instructions to MLLMs for better response generation, it can also be intentionally designed to expose critical safety issues of MLLMs (Liu et al., 2024b; Ni et al., 2024). Several works have explored jailbreaking of MLLMs with visual prompts. Instead of feeding harmful textual instructions directly, Gong et al. (2023) converts them into images through typography and feeds them to MLLMs as visual prompts. The results show that even if the underlying LLM has been aligned for safety, visual prompting opens a new jailbreaking surface generating harmful responses.

To further expose the safety problems of MLLMs for red-teaming, multimodal jailbreaking prompts combining both textual and visual instructions are also studied. Ying et al. (2024) first embed harmful perturbation in the visual prompt and then optimize the textual prompt through LLM reasoning on the harmful intent in the image. Meanwhile, Liu et al. (2024c) utilize a red-team MLLM and a red-team LLM guided by reinforcement learning to automatically generate visual and textual jailbreak-

king prompts respectively. Their results suggest that multimodal prompts could lead to stronger attack on MLLMs that fuse multimodal input features. Furthermore, Gu et al. (2024) observe a more severe safety issue of infectious jailbreak in multi-agent MLLM environments. With an adversarial image simply jailbreaking one agent and without any further intervention, almost all agents will start exhibit harmful behaviors in an exponential infection rate during multi-agent interaction.

9.2 Hallucination

The more fine-grained visual contexts provided with visual prompting are also useful for multimodal hallucination mitigation. To address the issue that MLLMs’ textual outputs are often not grounded in the reference images, Favero et al. (2024) propose a mutual-information decoding strategy to amplify the influence of visual prompts on model generation. To reduce MLLMs’ object hallucination and enhance fine-grained understanding in object-oriented perception tasks, Jiang et al. (2024) develop a prompting strategy jointly utilizing visual and textual prompts. A specialized detection model is employed to highlight relevant visual objects and visual prompts based on the key concepts extracted from textual prompts. While previous works mostly focus on single-object hallucination, Chen et al. (2024c) utilize visual referring prompts to evaluate multi-object hallucination of MLLMs. The results show MLLMs tend to experience more hallucinations when tasked with focusing multiple objects at the same time and authors suggest probing objects individually in visual prompts to enhance performances.

9.3 Debiasing

Despite the impressive capabilities of MLLMs, the biases and robustness of them remain a crucial challenge where models tend to utilize spurious correlations between input and target variables for predictions leading to potential social biases on certain topics, e.g., gender and racial biases (Ye et al., 2024). As visual prompting enables more fine-grained understanding of visual objects and relationships, it serves as a promising solution to mitigate potential biases in MLLMs’ generations by grounding the outputs with essential visual information and thus avoiding spurious correlations of non-essential inputs. It may also enhance the causal understanding of MLLMs between objects from the same modality and across different modalities for

generating more robust and grounded responses.

9.4 Visual Generation

Visual generation models, especially text-to-image diffusion models (Rombach et al., 2022), are becoming popular. Considering large-scale pre-trained diffusion models as MLLMs broadly, visual prompting plays an important role in controlling the generation and enable diffusion models for unseen visual tasks. Zhang et al. (2023c), Mou et al. (2024) propose ControlNet and T2I Adapter, which take various visual prompts for spatial control in image generation. In this survey, we discuss works that focus on visual prompting instead of controllable generation (Cao et al., 2024) in general. Prompt Diffusion (Wang et al., 2023) proposes a diffusion-based generative model that takes a novel vision-language prompts and outputs the target images, which unlocks the ability of in-context generation after fine-tuned on six visual tasks. ImageBrush (Yang et al., 2024b) proposes to achieve adaptive image manipulation under the instruction of a pair of exemplar demonstrations in order to address the issue of language ambiguity in image editing task. MPerceiver (Ai et al., 2024) introduces a multi-modal prompt learning approach using generative priors of diffusion models to enhance the all-in-one image restoration. Chen et al. (2024d) proposes VP3D, which leverages rich knowledge in 2D visual prompts to improve text-to-3D generation quality and trigger a new task of stylized text-to-3D generation. PromptCharm (Wang et al., 2024c) proposes an interaction system that supports text-to-image creation through multi-modal prompting and image refinement, which suggests the necessity of visual prompting for better image creation.

10 Conclusion

In this survey, we provide the first comprehensive review of visual prompting methods in MLLMs. We categorized various visual prompting techniques and discussed their generation processes, examining their integration into MLLMs for enhanced visual reasoning and perception. Our work also examines existing training and in-context learning methods in MLLMs with visual prompting. We inspire future directions that leverage visual prompts for better MLLM compositional reasoning.

11 Limitations

While our survey offers a comprehensive overview, it may be limited by the rapidly evolving nature of the field and potential gaps in the available literature. Future work should focus on expanding the scope of visual prompts and refining alignment techniques to further enhance MLLM capabilities.

References

- Yuang Ai, Huaibo Huang, Xiaoqiang Zhou, Jiexiang Wang, and Ran He. 2024. Multimodal prompt perceiver: Empower adaptiveness generalizability and fidelity for all-in-one image restoration. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 25432–25444.
- Zechen Bai, Pichao Wang, Tianjun Xiao, Tong He, Zongbo Han, Zheng Zhang, and Mike Zheng Shou. 2024. Hallucination of multimodal large language models: A survey. *arXiv preprint arXiv:2404.18930*.
- Amir Bar, Yossi Gandelsman, Trevor Darrell, Amir Globerson, and Alexei Efros. 2022. Visual prompting via image inpainting. *Advances in Neural Information Processing Systems*, 35:25005–25017.
- Garrick Brazil, Abhinav Kumar, Julian Straub, Nikhila Ravi, Justin Johnson, and Georgia Gkioxari. 2023. Omni3d: A large benchmark and model for 3d object detection in the wild. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 13154–13164.
- Mu Cai, Haotian Liu, Siva Karthik Mustikovela, Gregory P Meyer, Yuning Chai, Dennis Park, and Yong Jae Lee. 2024. Vip-llava: Making large multimodal models understand arbitrary visual prompts. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12914–12923.
- Pu Cao, Feng Zhou, Qing Song, and Lu Yang. 2024. Controllable generation with text-to-image diffusion models: A survey. *arXiv preprint arXiv:2403.04279*.
- Aochuan Chen, Yuguang Yao, Pin-Yu Chen, Yihua Zhang, and Sijia Liu. 2023a. Understanding and improving visual prompting: A label-mapping perspective. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 19133–19143.
- Jiaqi Chen, Bingqian Lin, Xinmin Liu, Xiaodan Liang, and Kwan-Yee K Wong. 2024a. Affordances-oriented planning using foundation models for continuous vision-language navigation. *arXiv preprint arXiv:2407.05890*.
- Keqin Chen, Zhao Zhang, Weili Zeng, Richong Zhang, Feng Zhu, and Rui Zhao. 2023b. Shikra: Unleashing multimodal llm’s referential dialogue magic. *arXiv preprint arXiv:2306.15195*.
- Wei Chen, Lin Li, Yongqi Yang, Bin Wen, Fan Yang, Tingting Gao, Yu Wu, and Long Chen. 2024b. Comm: A coherent interleaved image-text dataset for multimodal understanding and generation. *arXiv preprint arXiv:2406.10462*.
- Xuwei Chen, Ziqiao Ma, Xuejun Zhang, Sihan Xu, Shengyi Qian, Jianing Yang, David F Fouhey, and Joyce Chai. 2024c. Multi-object hallucination in vision-language models. *arXiv preprint arXiv:2407.06192*.
- Yang Chen, Yingwei Pan, Haibo Yang, Ting Yao, and Tao Mei. 2024d. Vp3d: Unleashing 2d visual prompt for text-to-3d generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4896–4905.
- Jang Hyun Cho, Boris Ivanovic, Yulong Cao, Edward Schmerling, Yue Wang, Xinshuo Weng, Boyi Li, Yurong You, Philipp Krähenbühl, Yan Wang, et al. 2024. Language-image models with 3d understanding. *arXiv preprint arXiv:2405.03685*.
- Ronghao Dang, Jianguan Feng, Haodong Zhang, Chongjian Ge, Lin Song, Lijun Gong, Chengju Liu, Qijun Chen, Feng Zhu, Rui Zhao, et al. 2023. Instructdet: Diversifying referring object detection with generalized instructions. *arXiv preprint arXiv:2310.05136*.
- Zhizhao Duan, Hao Cheng, Duo Xu, Xi Wu, Xiangxie Zhang, Xi Ye, and Zhen Xie. 2024. Cityllava: Efficient fine-tuning for vlms in city scenario. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7180–7189.
- Sheng Fan, Rui Liu, Wenguan Wang, and Yi Yang. 2024a. Navigation instruction generation with bev perception and large language models. *arXiv preprint arXiv:2407.15087*.
- Yue Fan, Lei Ding, Ching-Chen Kuo, Shan Jiang, Yang Zhao, Xinze Guan, Jie Yang, Yi Zhang, and Xin Eric Wang. 2024b. Read anywhere pointed: Layout-aware gui screen reading with tree-of-lens grounding. *arXiv preprint arXiv:2406.19263*.
- Yue Fan, Jing Gu, Kaiwen Zhou, Qianqi Yan, Shan Jiang, Ching-Chen Kuo, Xinze Guan, and Xin Eric Wang. 2024c. Muffin or chihuahua? challenging large vision-language models with multipanel vqa. *arXiv preprint arXiv:2401.15847*.
- Alessandro Favero, Luca Zancato, Matthew Trager, Siddharth Choudhary, Pramuditha Perera, Alessandro Achille, Ashwin Swaminathan, and Stefano Soatto. 2024. Multi-modal hallucination control by visual information grounding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14303–14312.
- Chaoyou Fu, Renrui Zhang, Haojia Lin, Zihan Wang, Timin Gao, Yongdong Luo, Yubo Huang, Zhengye Zhang, Longtian Qiu, Gaoxiang Ye, et al. 2023. A challenger to gpt-4v? early explorations of gemini in visual expertise. *arXiv preprint arXiv:2312.12436*.

- Jun Gao, Qian Qiao, Ziqiang Cao, Zili Wang, and Wenjie Li. 2024. Aim: Let any multi-modal large language models embrace efficient in-context learning. *arXiv preprint arXiv:2406.07588*.
- Ross Girshick. 2015. Fast r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 1440–1448.
- Yichen Gong, DeLong Ran, Jinyuan Liu, Conglei Wang, Tianshuo Cong, Anyu Wang, Sisi Duan, and Xiaoyun Wang. 2023. Figstep: Jailbreaking large vision-language models via typographic visual prompts. *arXiv preprint arXiv:2311.05608*.
- Jindong Gu, Zhen Han, Shuo Chen, Ahmad Beirami, Bailan He, Gengyuan Zhang, Ruotong Liao, Yao Qin, Volker Tresp, and Philip Torr. 2023. A systematic survey of prompt engineering on vision-language foundation models. *arXiv preprint arXiv:2307.12980*.
- Xiangming Gu, Xiaosen Zheng, Tianyu Pang, Chao Du, Qian Liu, Ye Wang, Jing Jiang, and Min Lin. 2024. Agent smith: A single image can jailbreak one million multimodal llm agents exponentially fast. *arXiv preprint arXiv:2402.08567*.
- Xixuan Hao, Wei Chen, Yibo Yan, Siru Zhong, Kun Wang, Qingsong Wen, and Yuxuan Liang. 2024. Urbanvlp: A multi-granularity vision-language pre-trained foundation model for urban indicator prediction. *arXiv preprint arXiv:2403.16831*.
- Junwen He, Yifan Wang, Lijun Wang, Huchuan Lu, Jun-Yan He, Jin-Peng Lan, Bin Luo, and Xuansong Xie. 2024. Multi-modal instruction tuned llms with fine-grained visual perception. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13980–13990.
- Mir Rayat Imtiaz Hossain, Mennatullah Siam, Leonid Sigal, and James J Little. 2024. Visual prompting for generalized few-shot segmentation: A multi-scale approach. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 23470–23480.
- Siyuan Huang, Haonan Chang, Yuhan Liu, Yimeng Zhu, Hao Dong, Peng Gao, Abdeslam Boularias, and Hongsheng Li. 2024a. A3vlm: Actionable articulation-aware vision language model. *arXiv preprint arXiv:2406.07549*.
- Wen Huang, Hongbin Liu, Minxin Guo, and Neil Zhenqiang Gong. 2024b. Visual hallucinations of multimodal large language models. *arXiv preprint arXiv:2402.14683*.
- Zhipeng Huang, Zhizheng Zhang, Zheng-Jun Zha, Yan Lu, and Baining Guo. 2024c. Relationvlm: Making large vision-language models understand visual relations. *arXiv preprint arXiv:2403.12801*.
- Mengzhao Jia, Zhihan Zhang, Wenhao Yu, Fangkai Jiao, and Meng Jiang. 2024. Describe-then-reason: Improving multimodal mathematical reasoning through visual comprehension training. *arXiv preprint arXiv:2404.14604*.
- Songtao Jiang, Yan Zhang, Chenyi Zhou, Yeying Jin, Yang Feng, Jian Wu, and Zuozhu Liu. 2024. Joint visual and text prompting for improved object-centric perception with multimodal large language models. *arXiv preprint arXiv:2404.04514*.
- Shibo Jie, Yehui Tang, Ning Ding, Zhi-Hong Deng, Kai Han, and Yunhe Wang. 2024. Memory-space visual prompting for efficient vision-language fine-tuning. *arXiv preprint arXiv:2405.05615*.
- Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. 2023. Segment anything. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4015–4026.
- Byung-Kwan Lee, Beomchan Park, Chae Won Kim, and Yong Man Ro. 2024. Collavo: Crayon large language and vision model. *arXiv preprint arXiv:2402.11248*.
- Xuanyu Lei, Zonghan Yang, Xinrui Chen, Peng Li, and Yang Liu. 2024a. Scaffolding coordinates to promote vision-language coordination in large multi-modal models. *arXiv preprint arXiv:2402.12058*.
- Yiming Lei, Jingqi Li, Zilong Li, Yuan Cao, and Hongming Shan. 2024b. Prompt learning in computer vision: a survey. *Frontiers of Information Technology & Electronic Engineering*, 25(1):42–63.
- Juncheng Li, Kaihang Pan, Zhiqi Ge, Minghe Gao, Wei Ji, Wenqiao Zhang, Tat-Seng Chua, Siliang Tang, Hanwang Zhang, and Yueting Zhuang. 2023a. Fine-tuning multimodal llms to follow zero-shot demonstrative instructions. In *The Twelfth International Conference on Learning Representations*.
- Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. 2023b. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In *International conference on machine learning*, pages 19730–19742. PMLR.
- Yinheng Li. 2023. A practical survey on zero-shot prompt design for in-context learning. *arXiv preprint arXiv:2309.13205*.
- Zeju Li, Chao Zhang, Xiaoyan Wang, Ruilong Ren, Yifan Xu, Ruifei Ma, and Xiangde Liu. 2024. 3dmit: 3d multi-modal instruction tuning for scene understanding. *arXiv preprint arXiv:2401.03201*.
- Zongjie Li, Chaozheng Wang, Chaowei Liu, Pingchuan Ma, Daoyuan Wu, Shuai Wang, and Cuiyun Gao. 2023c. Vrptest: Evaluating visual referring prompting in large multimodal models. *arXiv preprint arXiv:2312.04087*.
- Yuan-Hong Liao, Rafid Mahmood, Sanja Fidler, and David Acuna. 2024. Can feedback enhance semantic grounding in large vision-language models? *arXiv preprint arXiv:2404.06510*.

- Kevin Lin, Faisal Ahmed, Linjie Li, Chung-Ching Lin, Ehsan Azarnasab, Zhengyuan Yang, Jianfeng Wang, Lin Liang, Zicheng Liu, Yumao Lu, et al. 2023. Mmvid: Advancing video understanding with gpt-4v (ision). *arXiv preprint arXiv:2310.19773*.
- Weifeng Lin, Xinyu Wei, Ruichuan An, Peng Gao, Bocheng Zou, Yulin Luo, Siyuan Huang, Shanghang Zhang, and Hongsheng Li. 2024a. Draw-and-understand: Leveraging visual prompts to enable mllms to comprehend what you want. *arXiv preprint arXiv:2403.20271*.
- Yuanze Lin, Yunsheng Li, Dongdong Chen, Weijian Xu, Ronald Clark, Philip Torr, and Lu Yuan. 2024b. Rethinking visual prompting for multimodal large language models with external knowledge. *arXiv preprint arXiv:2407.04681*.
- Dingning Liu, Xiaomeng Dong, Renrui Zhang, Xu Luo, Peng Gao, Xiaoshui Huang, Yongshun Gong, and Zhihui Wang. 2023a. 3daxiesprompts: Unleashing the 3d spatial task capabilities of gpt-4v. *arXiv preprint arXiv:2312.09738*.
- Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2024a. Visual instruction tuning. *Advances in neural information processing systems*, 36.
- Weihuang Liu, Xi Shen, Chi-Man Pun, and Xiaodong Cun. 2023b. Explicit visual prompting for low-level structure segmentations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 19434–19445.
- Xin Liu, Yichen Zhu, Yunshi Lan, Chao Yang, and Yu Qiao. 2024b. Safety of multimodal large language models on images and text. *arXiv preprint arXiv:2402.00357*.
- Yi Liu, Chengjun Cai, Xiaoli Zhang, Xingliang Yuan, and Cong Wang. 2024c. Arondight: Red teaming large vision language models with auto-generated multi-modal jailbreak prompts. *arXiv preprint arXiv:2407.15050*.
- Bozhi Luan, Hao Feng, Hong Chen, Yonghui Wang, Wengang Zhou, and Houqiang Li. 2024. Textcot: Zoom in for enhanced multimodal text-rich image understanding. *arXiv preprint arXiv:2404.09797*.
- Chuofan Ma, Yi Jiang, Jiannan Wu, Zehuan Yuan, and Xiaojuan Qi. 2024a. Groma: Localized visual tokenization for grounding multimodal large language models. *arXiv preprint arXiv:2404.13013*.
- Huan Ma, Yan Zhu, Changqing Zhang, Peilin Zhao, Baoyuan Wu, Long-Kai Huang, Qinghua Hu, and Bingzhe Wu. 2024b. Invariant test-time adaptation for vision-language model generalization. *arXiv preprint arXiv:2403.00376*.
- Chong Mou, Xintao Wang, Liangbin Xie, Yanze Wu, Jian Zhang, Zhongang Qi, and Ying Shan. 2024. T2i-adapter: Learning adapters to dig out more controllable ability for text-to-image diffusion models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 4296–4304.
- Soroush Nasiriany, Fei Xia, Wenhao Yu, Ted Xiao, Jacky Liang, Ishita Dasgupta, Annie Xie, Danny Driess, Ayzaan Wahid, Zhuo Xu, et al. 2024. Pivot: Iterative visual prompting elicits actionable knowledge for vlms. *arXiv preprint arXiv:2402.07872*.
- Minheng Ni, Yeli Shen, Lei Zhang, and Wangmeng Zuo. 2024. Responsible visual editing. *arXiv preprint arXiv:2404.05580*.
- Changdae Oh, Hyeji Hwang, Hee-young Lee, YongTaek Lim, Geunyoung Jung, Jiyoung Jung, Hosik Choi, and Kyungwoo Song. 2023. Blackvip: Black-box visual prompting for robust transfer learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 24224–24235.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. 2022. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35:27730–27744.
- Kaihang Pan, Siliang Tang, Juncheng Li, Zhaoyu Fan, Wei Chow, Shuicheng Yan, Tat-Seng Chua, Yueting Zhuang, and Hanwang Zhang. 2024. Auto-encoding morph-tokens for multimodal llm. *arXiv preprint arXiv:2405.01926*.
- Leigang Qu, Haochuan Li, Tan Wang, Wenjie Wang, Yongqi Li, Liqiang Nie, and Tat-Seng Chua. 2024. Unified text-to-image generation and retrieval. *arXiv preprint arXiv:2406.05814*.
- Huali Ren, Anli Yan, Chong-zhi Gao, Hongyang Yan, Zhenxin Zhang, and Jin Li. 2024. Are you copying my prompt? protecting the copyright of vision prompt for vpaas via watermark. *arXiv preprint arXiv:2405.15161*.
- Razieh Rezaei, Masoud Jalili Sabet, Jindong Gu, Daniel Rueckert, Philip Torr, and Ashkan Khakzar. 2024. Learning visual prompts for guiding the attention of vision transformers. *arXiv preprint arXiv:2406.03303*.
- Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. 2022. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695.
- Sander Schulhoff, Michael Ilie, Nishant Balepur, Konstantine Kahadze, Amanda Liu, Chenglei Si, Yin-heng Li, Aayush Gupta, HyoJung Han, Sevien Schulhoff, et al. 2024. The prompt report: A systematic survey of prompting techniques. *arXiv preprint arXiv:2406.06608*.

- Dianmo Sheng, Dongdong Chen, Zhentao Tan, Qiankun Liu, Qi Chu, Jianmin Bao, Tao Gong, Bin Liu, Shengwei Xu, and Nenghai Yu. 2024. Towards more unified in-context visual understanding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13362–13372.
- Aleksandar Shtedritski, Christian Rupprecht, and Andrea Vedaldi. 2023. What does clip know about a red circle? visual prompt engineering for vlms. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 11987–11997.
- Quan Sun, Yufeng Cui, Xiaosong Zhang, Fan Zhang, Qiyang Yu, Yueze Wang, Yongming Rao, Jingjing Liu, Tiejun Huang, and Xinlong Wang. 2024. Generative multimodal models are in-context learners. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14398–14409.
- Georgios Tzifas and Hamidreza Kasaei. 2024. Towards open-world grasping with large vision-language models. *arXiv preprint arXiv:2406.18722*.
- David Wan, Jaemin Cho, Elias Stengel-Eskin, and Mohit Bansal. 2024. Contrastive region guidance: Improving grounding in vision-language models without training. *arXiv preprint arXiv:2403.02325*.
- Lei Wang, Wanyu Xu, Zhiqiang Hu, Yihuai Lan, Shan Dong, Hao Wang, Roy Ka-Wei Lee, and Ee-Peng Lim. 2024a. All in an aggregated image for in-image learning.
- Wenhai Wang, Zhe Chen, Xiaokang Chen, Jiannan Wu, Xizhou Zhu, Gang Zeng, Ping Luo, Tong Lu, Jie Zhou, Yu Qiao, et al. 2024b. Visionllm: Large language model is also an open-ended decoder for vision-centric tasks. *Advances in Neural Information Processing Systems*, 36.
- Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A Smith, Daniel Khashabi, and Hananeh Hajishirzi. 2022. Self-instruct: Aligning language models with self-generated instructions. *arXiv preprint arXiv:2212.10560*.
- Zhendong Wang, Yifan Jiang, Yadong Lu, Pengcheng He, Weizhu Chen, Zhangyang Wang, Mingyuan Zhou, et al. 2023. In-context learning unlocked for diffusion models. *Advances in Neural Information Processing Systems*, 36:8542–8562.
- Zhijie Wang, Yuheng Huang, Da Song, Lei Ma, and Tianyi Zhang. 2024c. Promptcharm: Text-to-image generation through multi-modal prompting and refinement. In *Proceedings of the CHI Conference on Human Factors in Computing Systems*, pages 1–21.
- Jiannan Wu, Muyan Zhong, Sen Xing, Zeqiang Lai, Zhaoyang Liu, Wenhai Wang, Zhe Chen, Xizhou Zhu, Lewei Lu, Tong Lu, et al. 2024a. Visionllm v2: An end-to-end generalist multimodal large language model for hundreds of vision-language tasks. *arXiv preprint arXiv:2406.08394*.
- Junda Wu, Xintong Li, Tong Yu, Yu Wang, Xiang Chen, Jiuxiang Gu, Lina Yao, Jingbo Shang, and Julian McAuley. 2024b. Commit: Coordinated instruction tuning for multimodal large language models. *arXiv preprint arXiv:2407.20454*.
- Junda Wu, Rui Wang, Handong Zhao, Ruiyi Zhang, Chaochao Lu, Shuai Li, and Ricardo Henao. 2023. Few-shot composition learning for image retrieval with prompt tuning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pages 4729–4737.
- Junfeng Wu, Yi Jiang, Qihao Liu, Zehuan Yuan, Xiang Bai, and Song Bai. 2024c. General object foundation model for images and videos at scale. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3783–3795.
- Mingrui Wu, Xinyue Cai, Jiayi Ji, Jiale Li, Oucheng Huang, Gen Luo, Hao Fei, Xiaoshuai Sun, and Rongrong Ji. 2024d. Controlmllm: Training-free visual prompt learning for multimodal large language models. *arXiv preprint arXiv:2407.21534*.
- Tung-Yu Wu, Sheng-Yu Huang, and Yu-Chiang Frank Wang. 2024e. Dora: 3d visual grounding with order-aware referring. *arXiv preprint arXiv:2403.16539*.
- Yixuan Wu, Yizhou Wang, Shixiang Tang, Wenhao Wu, Tong He, Wanli Ouyang, Jian Wu, and Philip Torr. 2024f. Dettolchain: A new prompting paradigm to unleash detection ability of mllm. *arXiv preprint arXiv:2403.12488*.
- Enze Xie, Wenhai Wang, Zhiding Yu, Anima Anandkumar, Jose M Alvarez, and Ping Luo. 2021. Segformer: Simple and efficient design for semantic segmentation with transformers. *Advances in neural information processing systems*, 34:12077–12090.
- Junyu Xie, Tengda Han, Max Bain, Arsha Nagrai, Gül Varol, Weidi Xie, and Andrew Zisserman. 2024. Autoad-zero: A training-free framework for zero-shot audio description. *arXiv preprint arXiv:2407.15850*.
- Chengming Xu, Chen Liu, Yikai Wang, and Yanwei Fu. 2024a. Towards global optimal visual in-context learning prompt selection. *arXiv preprint arXiv:2405.15279*.
- Xin Xu, Yue Liu, Panupong Pasupat, Mehran Kazemi, et al. 2024b. In-context learning with retrieved demonstrations for language models: A survey. *arXiv preprint arXiv:2401.11624*.
- An Yan, Zhengyuan Yang, Junda Wu, Wanrong Zhu, Jianwei Yang, Linjie Li, Kevin Lin, Jianfeng Wang, Julian McAuley, Jianfeng Gao, et al. 2024. List items one by one: A new data source and learning paradigm for multimodal llms. *arXiv preprint arXiv:2404.16375*.

- Jianwei Yang, Hao Zhang, Feng Li, Xueyan Zou, Chunyuan Li, and Jianfeng Gao. 2023. Set-of-mark prompting unleashes extraordinary visual grounding in gpt-4v. *arXiv preprint arXiv:2310.11441*.
- Lingfeng Yang, Yueze Wang, Xiang Li, Xinlong Wang, and Jian Yang. 2024a. Fine-grained visual prompting. *Advances in Neural Information Processing Systems*, 36.
- Yifan Yang, Houwen Peng, Yifei Shen, Yuqing Yang, Han Hu, Lili Qiu, Hideki Koike, et al. 2024b. Imagebrush: Learning visual in-context instructions for exemplar-based image manipulation. *Advances in Neural Information Processing Systems*, 36.
- Wenqian Ye, Guangtao Zheng, Yunsheng Ma, Xu Cao, Bolin Lai, James M Rehg, and Aidong Zhang. 2024. Mm-spubench: Towards better understanding of spurious biases in multimodal llms. *arXiv preprint arXiv:2406.17126*.
- Zonghao Ying, Aishan Liu, Tianyuan Zhang, Zhengmin Yu, Siyuan Liang, Xianglong Liu, and Dacheng Tao. 2024. Jailbreak vision language models via bi-modal adversarial prompt. *arXiv preprint arXiv:2406.04031*.
- Jaehong Yoon, Shoubin Yu, and Mohit Bansal. 2024. Raccoon: Remove, add, and change video content with auto-generated narratives. *arXiv preprint arXiv:2405.18406*.
- Ao Zhang, Hao Fei, Yuan Yao, Wei Ji, Li Li, Zhiyuan Liu, and Tat-Seng Chua. 2024a. Vpgrans: Transfer visual prompt generator across llms. *Advances in Neural Information Processing Systems*, 36.
- Chaoning Zhang, Fachrina Dewi Puspitasari, Sheng Zheng, Chenghao Li, Yu Qiao, Taegoo Kang, Xinru Shan, Chenshuang Zhang, Caiyan Qin, Francois Rameau, et al. 2023a. A survey on segment anything model (sam): Vision foundation model meets prompt engineering. *arXiv preprint arXiv:2306.06211*.
- Hao Zhang, Feng Li, Xueyan Zou, Shilong Liu, Chunyuan Li, Jianwei Yang, and Lei Zhang. 2023b. A simple framework for open-vocabulary segmentation and detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1020–1031.
- Jingyi Zhang, Jiaying Huang, Sheng Jin, and Shijian Lu. 2024b. Vision-language models for vision tasks: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*.
- Lu Zhang, Tiancheng Zhao, Heting Ying, Yibo Ma, and Kyusong Lee. 2024c. Omagent: A multi-modal agent framework for complex video understanding with task divide-and-conquer. *arXiv preprint arXiv:2406.16620*.
- Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. 2023c. Adding conditional control to text-to-image diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3836–3847.
- Sha Zhang, Di Huang, Jiajun Deng, Shixiang Tang, Wanli Ouyang, Tong He, and Yanyong Zhang. 2024d. Agent3d-zero: An agent for zero-shot 3d understanding. *arXiv preprint arXiv:2403.11835*.
- Tao Zhang, Xiangtai Li, Hao Fei, Haobo Yuan, Shengqiong Wu, Shunping Ji, Chen Change Loy, and Shuicheng Yan. 2024e. Omg-llava: Bridging image-level, object-level, pixel-level reasoning and understanding. *arXiv preprint arXiv:2406.19389*.
- Wei Zhang, Miaoxin Cai, Tong Zhang, Yin Zhuang, and Xuerui Mao. 2024f. Earthmarker: A visual prompt learning framework for region-level and point-level remote sensing imagery comprehension. *arXiv preprint arXiv:2407.13596*.
- Yichi Zhang, Yinpeng Dong, Siyuan Zhang, Tianzan Min, Hang Su, and Jun Zhu. 2024g. Exploring the transferability of visual prompting for multimodal large language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 26562–26572.
- Zheng Zhang, Yeyao Ma, Enming Zhang, and Xiang Bai. 2024h. Psalm: Pixelwise segmentation with large multi-modal model. *arXiv preprint arXiv:2403.14598*.
- Wenliang Zhong, Wenyi Wu, Qi Li, Rob Barton, Boxin Du, Shioulin Sam, Karim Bouyarmene, Ismail Tutar, and Junzhou Huang. 2024. Enhancing multimodal large language models with multi-instance visual prompt generator for visual representation enrichment. *arXiv preprint arXiv:2406.02987*.
- Enshen Zhou, Yiran Qin, Zhenfei Yin, Yuzhou Huang, Ruimao Zhang, Lu Sheng, Yu Qiao, and Jing Shao. 2024a. Minedreamer: Learning to follow instructions via chain-of-imagination for simulated-world control. *arXiv preprint arXiv:2403.12037*.
- Qiji Zhou, Ruochen Zhou, Zike Hu, Panzhong Lu, Siyang Gao, and Yue Zhang. 2024b. Image-of-thought prompting for visual reasoning refinement in multimodal large language models. *arXiv preprint arXiv:2405.13872*.