Multimodal Chain-of-Thought Reasoning: A Comprehensive Survey

Yaoting Wang¹, Shengqiong Wu¹, Yuechen Zhang², Shuicheng Yan¹, Ziwei Liu³, Jiebo Luo⁴, Hao Fei^{1*}

¹NUS, ²CUHK, ³NTU, ⁴UR

Survey Project: https://github.com/yaotingwangofficial/Awesome-MCoT

Abstract

By extending the advantage of chain-of-thought (CoT) reasoning in human-like step-by-step processes to multimodal contexts, multimodal CoT (MCoT) reasoning has recently garnered significant research attention, especially in the integration with multimodal large language models (MLLMs). Existing MCoT studies design various methodologies and innovative reasoning paradigms to address the unique challenges of image, video, speech, audio, 3D, and structured data across different modalities, achieving extensive success in applications such as robotics, healthcare, autonomous driving, and multimodal generation. However, MCoT still presents distinct challenges and opportunities that require further focus to ensure consistent thriving in this field, where unfortunately an up-to-date review of this domain is lacking. To bridge this gap, we present the first systematic survey of MCoT reasoning, elucidating the relevant foundational concepts and definitions. We offer a comprehensive taxonomy and an in-depth analysis of current methodologies from diverse perspectives across various application scenarios. Furthermore, we provide insights into existing challenges and future research directions, aiming to foster innovation toward multimodal AGI.

Keywords— Multimodal Reasoning, Chain-of-Thought, Multimodal Large Language Models



^{*}Corresponding Author. (haofei37@nus.edu.sg)

Contents

1	Intr	oduction	3				
	1.1	Contributions	4				
	1.2	Survey Organization	4				
2	Background and Preliminary						
	2.1	From CoT to MCoT	6				
	2.2	Thought Paradigm	7				
	2.3	Multimodal LLMs	8				
3	MC	oT Reasoning Under Various Modalities	9				
	3.1	MCoT Reasoning over Image	9				
	3.2	MCoT Reasoning over Video	10				
	3.3	MCoT Reasoning over 3D	11				
	3.4	MCoT Reasoning over Audio and Speech	11				
	3.5	MCoT Reasoning over Table and Chart	12				
	3.6	Cross-modal CoT Reasoning	12				
4	Methodologies in MCoT Reasoning						
	4.1	From Rationale Construction Perspective	13				
	4.2	From Structural Reasoning Perspective	13				
	4.3	From Information Enhancing Perspective	15				
	4.4	From Objective Granularity Perspective	15				
	4.5	From Multimodal Rationale Perspective	16				
	4.6	From Test-Time Scaling Perspective	17				
5	App	lications with MCoT Reasoning	18				
	5.1	Embodied AI	18				
	5.2	Agentic System	19				
	5.3	Autonomous Driving	19				
	5.4	Medical and Healthcare	19				
	5.5	Social and Human	20				
	5.6	Multimodal Generation	20				
6	MC	oT Datasets and Benchmarks	20				
	6.1	Datasets for MLLMs Finetuning with Rationale	20				
	6.2	Benchmarks for Downstream Capability Assessment	22				
7	Lim	itations, Challenges and Future Directions	22				
8	Con	clusion	25				

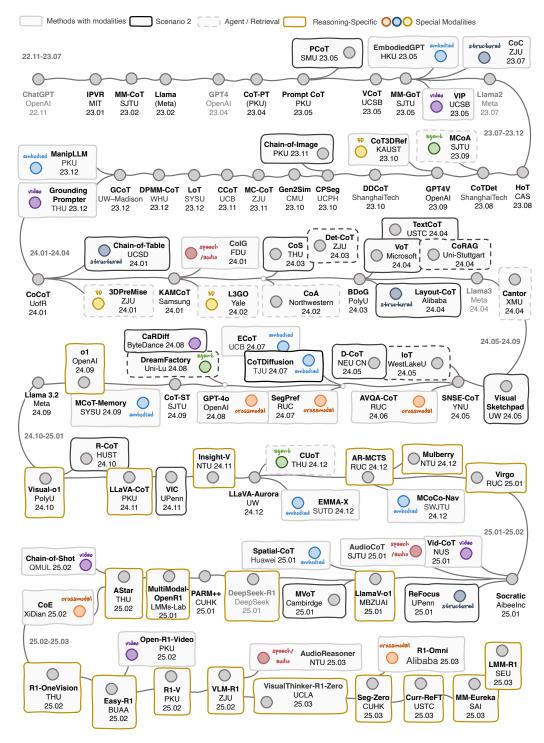


Figure 1: Developing timeline of Multimodal Chain-of-Thought (MCoT) reasoning. Models with names in gray are text-only LLMs. For clarity, the models in the figure are assumed to include the image modality by default, unless specified with special modalities indicated by colored circles.

1 Introduction

The emergence of large language models (LLMs) [1–7] has ushered in an unprecedented era in the artificial intelligence (AI) field. It has been long recognized the necessity of aligning with the inherently multimodal nature of real-world environments, and correspondingly, the AI field evolves from

LLMs to multimodal LLMs (MLLMs) [8–18], integrating diverse modalities into language intelligence. Achieving human-level intelligence requires transcending basic perceptual capabilities to attain sophisticated cognitive reasoning—a hallmark of human cognition that enables iterative reasoning through contextual understanding and self-correction. Inspired by this observation, in-context learning (ICL) techniques have empowered LLMs to demonstrate stepwise reasoning—commonly known as chain-of-thought (CoT) reasoning mechanisms [19–24]. This technique enables models to break down problems into a series of intermediate steps, enhancing both transparency in decision-making and performance on intricate reasoning tasks. The remarkable success of CoT reasoning on a wide range of downstream complex tasks has driven its widespread adoption across academia and industry. Especially the recent advancements implicitly integrating this capability into cutting-edge systems like OpenAI's o1/o3 [25] and DeepSeek R1 [26] has garnered widespread attention.

The integration of CoT reasoning into multimodal contexts has then catalyzed transformative progress in AI, giving rise to Multimodal Chain-of-Thought (MCoT) reasoning [27, 28]. The MCoT topic has generated a spectrum of innovative outcomes due to both the CoT attributes and the heterogeneous nature of cross-modal data interactions. On one hand, the original CoT framework has evolved into advanced reasoning architectures incorporating hierarchical thought structures, from linear sequences [19] to graph-based representations [23]. On the other hand, unlike the unimodal text setting, diverse modalities such as visual, auditory, and spatiotemporal data demand specialized processing strategies—visual reasoning requires precise perception and analysis of static scenes, object relationship, while video understanding necessitates robust temporal dynamics modeling. These requirements have spurred the development of array of sophisticated MCoT methodologies that adapt reasoning processes to modality-specific characteristics, such as Multimodal-CoT [29], MVoT [30], Video-of-Thought [31], Audio-CoT [32], Cot3DRef [33] and PARM++ [34]. The demonstrated effectiveness of MCoT has also led to its successful application in critical domains such as autonomous driving [35–38], embodied AI [39–41], robotics [42–45] and healthcare [46–50], positioning it as a foundational technology for achieving multimodal AGI.

In recent years, research on MCoT has attracted growing attention. Figure 1 presents a comprehensive timeline of key milestones in this emerging field. Despite its promising potential in enhancing multimodal reasoning, MCoT also poses significant challenges and leaves several critical questions unanswered—for example, determining the most effective strategies for leveraging varied multimodal context, designing CoT processes that truly enhance MLLMs' reasoning capabilities, and implementing implicit reasoning within these models. Notably, the absence of comprehensive surveys hinders knowledge consolidation in this emerging field. To bridge this critical gap, this paper provides the first systematic overview of MCoT reasoning, providing a structured analysis of technological development, methodologies, practical applications, and future directions. We hope this survey will serve as an authoritative reference, spurring further innovation and progress in this rapidly evolving domain.

1.1 Contributions

- **First Survey:** This paper represents the first survey dedicated to an inaugural thorough review of MCoT reasoning.
- **Comprehensive Taxonomy:** We propose a meticulous taxonomy (cf. Figure 2) that categorizes the diverse approaches in MCoT research.
- Frontiers and Future Directions: We discuss emerging challenges and outline promising avenues for future research.
- **Resource Sharing:** We compile and make publicly available all relevant resources to support and accelerate progress within the research community.

1.2 Survey Organization

The remainder of this survey is organized as follows. We begin by introducing the fundamental concepts and background knowledge related to MCoT ($\S 2$). We then review the state-of-the-art research in MCoT across different modalities ($\S 3$). Next, we provide a taxonomy and consolidate the mainstream methods in MCoT under various perspectives ($\S 4$). Following this, we summarize the extensive downstream applications of MCoT ($\S 5$). Subsequently, we present an overview of datasets and benchmarks from multiple perspectives ($\S 6$). Finally, we discuss the challenges and future directions in this field ($\S 7$).

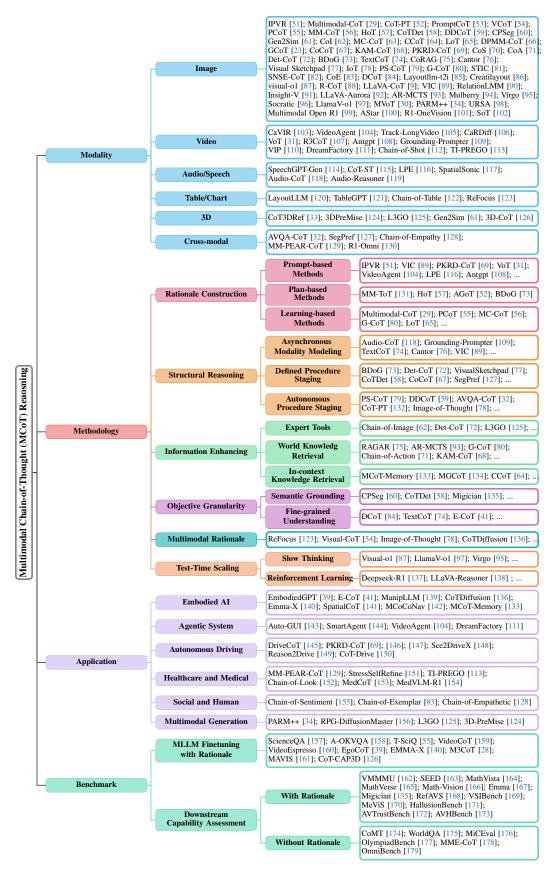


Figure 2: Taxonomy of MCoT reasoning.

Terms	Abbrev.	Description
In-context Learning	ICL	Prompting LLMs with task-specific examples without additional explicit training.
Chain-of-Thought	СоТ	Prompting LLMs to reason step-by-step or breaks complex problems into logical steps.
Multimodal CoT	MCoT	Extends CoT to reason with multimodalities, e.g., audio, image.
Cross-modal CoT		Reasoning with two or more multimodalities, e.g., audio-visual.
Thought		A single reasoning step in CoT.
Rationale		Built upon multiple thoughts to support the final answer.

Table 1: Interpretation of MCoT-related terms.

2 Background and Preliminary

Recent advancements in the scale of model pretraining have driven a significant shift in the application paradigm of language models, transitioning from the conventional "pretrain-then-finetune" approach to a more adaptive "pretrain-then-prompt" framework [180–184]. Within this evolving landscape, researchers have explored innovative techniques to enhance the reasoning capabilities of LLMs for complex tasks, notably ICL [180, 185, 186] and CoT¹ reasoning [19]. The essence of ICL lies in supplying task-relevant examples or demonstrations within the prompt, enabling LLMs to better interpret user intent and generate outputs aligned with expectations. This method leverages contextual guidance to steer the model toward task-appropriate responses. In contrast, CoT reasoning emulates human problem-solving by decomposing complex tasks into a sequence of manageable sub-tasks, systematically constructing solutions. The intermediate reasoning steps or trajectories, termed the rationale, elucidate the logical progression underlying the model's conclusions. Building on this foundation, MCoT reasoning extends the CoT paradigm by incorporating diverse data modalities, such as images, videos, and audio. This augmentation broadens the scope of multi-step reasoning, enhancing its applicability to increasingly intricate scenarios.

2.1 From CoT to MCoT

We provide explanations of the terms related to MCoT in Table 1. To formalize the MCoT framework, we begin by defining \mathcal{P} , \mathcal{S} , \mathcal{Q} , \mathcal{A} and \mathcal{R} to represent the prompt, instruction, query, answer, and rationale, respectively. Each of these elements is represented as a sequence of language tokens, with length denoted as $|\cdot|$. We also use lowercase letter to denote individual tokens, e.g., a_i refers to the i-th token of the answer \mathcal{A} . Next, we define a standard ICL process that integrates few-shot demonstration pairs, which can be expressed as follows:

$$\mathcal{P}_{ICL} = \{ \mathcal{S}, (x_1, y_1), \dots, (x_n, y_n) \}, \tag{1}$$

where \mathcal{P}_{ICL} represents the prompt for ICL, consisting of an instruction \mathcal{S} along with n demonstration pairs of questions x and their corresponding answers y. Then, the probability of generating an answer sequence \mathcal{A} given the prompt \mathcal{P}_{ICL} and a query \mathcal{Q} is mathematically defined as:

$$p(\mathcal{A} \mid \mathcal{P}_{ICL}, \mathcal{Q}) = \prod_{i=1}^{|\mathcal{A}|} \mathcal{F}(a_i \mid \mathcal{P}_{ICL}, \mathcal{Q}, a_{< i}),$$
 (2)

where \mathcal{F} denotes the probabilistic language model. Note that when n=0, the process simplifies to the standard zero-shot prompting scenario.

Then, we can define the vanilla CoT as:

$$\mathcal{P}_{CoT} = \{ \mathcal{S}, (x_1, e_1, y_1), \dots, (x_n, e_n, y_n) \},$$
(3)

where \mathcal{P}_{CoT} denotes the prompt used for CoT reasoning, and e_i represents the example rational. Next, we define the joint probability of generating an answer \mathcal{A} and rationale \mathcal{R} given the input

¹For clarity and consistency, we use "CoT" to denote multi-step reasoning technologies and "topology" to describe distinct thought structures (*e.g.*, chain or graph topologies). Specific approaches are assigned unique identifiers, such as vanilla CoT.

prompt \mathcal{P}_{CoT} and a query \mathcal{Q} :

$$p(\mathcal{A}, \mathcal{R} \mid \mathcal{P}_{CoT}, \mathcal{Q}) = p(\mathcal{A} \mid \mathcal{P}_{CoT}, \mathcal{Q}, \mathcal{R}) \cdot p(\mathcal{R} \mid \mathcal{P}_{CoT}, \mathcal{Q}), \tag{4}$$

where the right side represents two conditional probabilities of generating the answer A and rationale R, which can be defined as:

$$p(\mathcal{R} \mid \mathcal{P}_{CoT}, \mathcal{Q}) = \prod_{i=1}^{|\mathcal{R}|} \mathcal{F}(r_i \mid \mathcal{P}_{CoT}, \mathcal{Q}, r_{< i}),$$
 (5)

$$p(\mathcal{A} \mid \mathcal{P}_{CoT}, \mathcal{Q}, \mathcal{R}) = \prod_{i=1}^{|\mathcal{A}|} \mathcal{F}(a_i \mid \mathcal{P}_{CoT}, \mathcal{Q}, a_{< i}).$$
 (6)

In contrast to the ICL approach, as shown in Equation (2), the CoT framework necessitates the generation of a rationale \mathcal{R} prior to arriving at the answer \mathcal{A} , as reflected in Equations (5) and (6).

When considering MCoT, it is crucial to highlight that, unlike CoT, MCoT incorporates multimodal information into the components $\mathcal{P}, \mathcal{Q}, \mathcal{A}$, and \mathcal{R} . However, it is not necessary for all these components to simultaneously encompass multimodal information. That is, given language-based input \mathcal{T} and language-excluded multimodal context \mathcal{M} , we have $\exists \vartheta \in \{\mathcal{P}, \mathcal{Q}, \mathcal{A}, \mathcal{R}\} : \mathcal{M}(\vartheta)$. Hence we categorize MCoT into two distinct scenarios based on the composition of rationale: one relying exclusively on language-based information and another incorporating multimodal signals beyond linguistic content:

► *Scenario-1:* MCoT with text-only thought to tackle multimodal input and output:

$$\mathcal{R} \in \mathcal{L}$$
.

▶ Scenario-2: MCoT with multimodal thought to tackle unimodal or multimodal scenes:

$$\mathcal{R} \in \left\{\mathcal{M}, \mathcal{M} \oplus \mathcal{L}\right\}.$$

Scenario-1 aims to address tasks that involve multimodal information in either the input or output while utilizing a rationale composed solely of language. In contrast, Scenario-2 highlights the integration of given, retrieved or generated multimodal information within the rationale itself.

2.2 Thought Paradigm

Since the introduction of vanilla CoT [19], various paradigms have emerged to enhance multimodal and multi-step reasoning. Based on the construction of thought generation during reasoning, the community categorizes the reasoning structures [187] or topologies [188] into chain, tree, and graph types, as depicted in Figure 3. In these topologies, thoughts are represented as nodes, with edges indicating dependencies between them [188]. Chain topologies [19, 189, 190] facilitate linear and sequential thought generation, progressively converging toward the final answer. However, chain topologies lack the capacity for in-depth exploration of individual thoughts during reasoning.

In contrast, tree topologies [191, 192] enable exploration and backtracking within the reasoning process. At each node (*i.e.*, thought) in a tree topology, a thought generator produces multiple child nodes, as illustrated on the left side of Figure 3.C. These child nodes are then evaluated by a state evaluator, which assigns scores to them. These scores can be derived from the LLM itself or based on specific rules. Search algorithms, such as breadth-first search (BFS) or depth-first search (DFS), then guide the tree's expansion.

Graph topologies [23] also allow for the generation of multiple child nodes from a single parent. However, they introduce cycles and N-to-1 connections, meaning that a single node can have multiple parent nodes. This facilitates aggregation among multiple nodes, as illustrated by the blue arrows in Figure 3. Hypergraph topologies [57] extend graph topologies by employing hyperedges, which connect more than two thoughts. This structure inherently supports joint reasoning by integrating information from diverse modalities. Furthermore, self-consistency [193] can be seamlessly incorporated into various reasoning methods. For instance, using the chain topology as a baseline

(Figure 3.A), multiple chain-based reasoning processes can be executed in parallel, with the final answer determined by majority voting to ensure consistency across several rationales. Overall, the evolution of reasoning topologies reflects a progression from linear dependencies to branching exploration, aggregation with refinement, and higher-order associations.

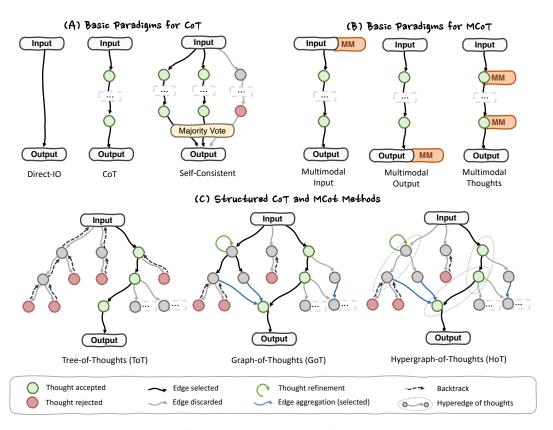


Figure 3: Different thought paradigms of CoT and MCoT.

2.3 Multimodal LLMs

The release of models such as GPT-4V [181], Gemini 2.0 [7], and Claude3 [194] has demonstrated remarkable capabilities in multimodal understanding, sparking significant interest in MLLMs within the research community. Initial investigations into MLLMs focused on developing robust language models capable of interpreting multimodal content and generating textual responses. In the domain of image-text understanding, notable progress has been achieved with Visual Large Language Models (VLLMs) such as BLIP2 [195], OpenFlamingo [196], MiniGPT-4 [197], and LLaVA [13]. Concurrently, advancements in video-text understanding have emerged, with significant contributions from VideoChat [198] and Video-ChatGPT [17]. Audio and speech comprehension have also garnered attention, exemplified by models like Qwen-Audio [199, 16] and LLaSM [200]. A noteworthy development is VideoLLaMA [18], which leverages Qformer [195] to enable both audio and video understanding. In simple terms, mainstream MLLMs typically follow a consistent model architecture by processing multimodal embeddings or tokens into the decoder structure and generating contextually relevant outputs in an autoregressive manner, as shown in the left of Figure 4.

Parallel to these works about multimodal understanding, research also explored multimodal content generation. In image generation, models such as Kosmos-2 [201], GILL [202], Emu [203], and MiniGPT-5 [204] have achieved breakthroughs. Audio generation has seen advancements with SpeechGPT [205, 206] and AudioPaLM [207], while video generation research, including CogVideo [208], VideoPoet [209], Video-Lavit [210], and StreamingT2V [211], has laid the groundwork for multimodal content creation. The recent introduction of GPT-4o [212], capable of both understanding and generating images and audio, has shifted attention toward "any-to-any" paradigm models.

Common Architectures of MLLMs

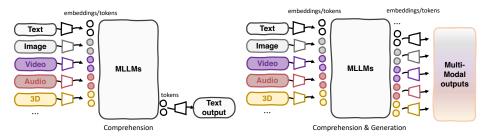


Figure 4: Common architectures for comprehension-only and comprehension-generation MLLMs.

Prior works, such as NExT-GPT [14] advances this objective for the first time by integrating multimodal adapters with various diffusion models. AnyGPT [213] utilizes multimodal discrete tokens to facilitate the generation of diverse multimodal content. Subsequently, Mini-Omni2 [214, 215] introduces a command-based interruption mechanism, enhancing user interaction and aligning further with GPT-4o's capabilities. Compared to MLLMs that only support comprehension, as shown in Figure 4, MLLMs that integrate both comprehension and generation either utilize an autoregressive approach to generate multimodal tokens [213], or connect decoders of varying modalities to decode multimodal embeddings [14].

Most recently, the release of the reasoning-focused OpenAI o1 [216] model has drawn interest in enhancing reasoning capabilities through deliberate, extended processing and test-time scaling. Models such as Mulberry [94], AStar [100], and LlamaV-o1 [97], by adopting long-MCoT reasoning strategies, have demonstrated robust performance in multimodal reasoning, further advancing the field of multimodal understanding.

3 MCoT Reasoning Under Various Modalities

MCoT extends the reasoning capabilities of LLMs/MLLMs to tackle complex tasks across diverse modalities, *e.g.*, images, videos, audio, 3D, tables/charts and beyond, by employing chained thought reasoning. As demonstrated in Figure 5, integrating MCoT with these modalities facilitates the implementation of numerous fundamental and significant applications. This section systematically reviews MCoT research across these modalities, emphasizing key advancements and their contributions to the development of multimodal reasoning.

3.1 MCoT Reasoning over Image

The prevalence of image data and associated tasks has driven the extensive application of MCoT mostly in Visual Question Answering (VQA). Early implementations, such as IPVR [51] and Multimodal-CoT [29], establish the foundational MCoT framework by generating intermediate rationales before final predictions. Subsequent advancements have further refined this paradigm: MC-CoT [56] integrates self-consistency [193] with MCoT, employing word-level majority voting during training to enhance the quality of generated rationales. SoT [102] leverages a router model to dynamically select reasoning paradigms (*i.e.*, conceptual chaining, chunked symbolism, and expert lexicons) inspired by human cognitive strategies to enhance reasoning efficiency. CoCoT [67] improves multi-image comprehension in MLLMs through similarity and difference analysis across inputs, while RelationLMM [90] explicitly addresses object relationship modeling via task decomposition. HoT [57] extends the Graph-of-Thought framework by introducing hyperedges to connect multiple reasoning nodes, thereby enhancing multimodal reasoning capabilities.

Structured reasoning mechanisms have been proposed to enhance controllability and interpretability. DDCoT [59] and Socratic Questioning [96] employ staged reasoning processes to systematically refine multimodal outcomes. Interaction methodologies between text and vision modalities also critically influence rationale generation. Chain-of-Spot [70], TextCoT [74], and DCoT [84] prioritize region-of-interest analysis to improve contextual understanding. RAGAR [75] and Cantor [76] integrate automated processes with low-level image attributes to strengthen reasoning, whereas



Figure 5: Examples of MCoT applications in various modalities and tasks.

KAM-CoT [68] and PKRD-CoT [69] incorporate external knowledge bases, which are further augmented by graph-based techniques described in [134] and [73]. The reliance of MCoT on annotated reasoning data has spurred research into automated data augmentation. G-CoT [80], STIC [81], PS-CoT [79], SNSE-CoT [82], Chain-of-Exemplar [83], and R-CoT [88] address this limitation by innovating methods to automate and enhance training data generation. In addition, static image feature extraction often results in inconsistencies when handling complex reasoning demands. To mitigate this, DPMM-CoT [66] and LLavA-AURORA [92] regenerate image features from latent space. Beyond text-based rationales, recent approaches leverage multimodal rationales for comprehensive reasoning, for instance, Visual-CoT [54], Chain-of-Image [62], VisualSketchpad [77], MVoT [30], and Visualization-of-Thought [217] effectively process multimodal scenes with multimodal thoughts.

MCoT's applicability also extends beyond VQA to specialized domains. For fine-grained instance-level tasks, CoTDet [58], Det-Cot [72], and CPSeg [60] demonstrate notable advancements. In image generation, PromptCoT [53] focuses refining the input prompts, PARM++ [34] optimizes reward mechanisms, and LayoutLLM-T2I [85] and CreatiLayout [86] employ text-based layout construction prior to synthesis, significantly improving output quality.

3.2 MCoT Reasoning over Video

Video understanding also relies on essential reasoning capabilities, as beyond the image understanding that require process static visual content and spatial relationships, videos present challenges of temporal dynamics, particularly in the case of long videos. As a basic usage, CaVIR [103] enhances intent question answering, which demands contextual and commonsense comprehension, by imple-

menting a zero-shot MCoT approach. Similarly, VideoAgent [104] and HM-Prompt [105] employ zero-shot MCoT to improve long-video reasoning and reduce hallucinations. AntGPT [108] extends few-shot MCoT to action categorization in egocentric videos. In generative tasks, DreamFactory [111] applies few-shot MCoT to produce consistent key frames for long-video synthesis.

For complex video comprehension, Video-of-Thought [31] proposes a comprehensive five-stage framework: task and target identification, object tracking, action analysis, ranked question answering, and answer verification. This structured approach ensures thorough interpretation of video content. Likewise, CaRDiff [106] decomposes intricate video tasks into sub-components—caption generation, saliency inference, and bounding box production—to guide diffusion processes for salient object mask creation. R3CoT [107] introduces a three-stage model (*i.e.*, refining, retrieving, reasoning) tailored for video rumor detection, while Grounding-Prompter [109] integrates global and local perception for temporal sentence grounding, localizing video moments based on linguistic queries. Efficiency in long-video analysis is addressed by frameworks such as VIP [110], which prioritizes key frames and extracts critical features (*e.g.*, focus, action, emotion, objects, background) to evaluate reasoning through attribute prediction for intermediate and future frames. Chain-of-Shot [112] further optimizes frame sampling by employing binary video summaries during training and assessing frame-task relevance for efficient inference.

MCoT's utility also spans specialized domains, such as medical video analysis [113, 151] and affective computing [218, 219, 130]. Collectively, these advancements underscore MCoT's role in decomposing complex video tasks, enhancing reasoning accuracy, and improving computational efficiency across diverse applications, marking a significant step forward in long-video understanding.

3.3 MCoT Reasoning over 3D

Reasoning in 3D scenes entails significant challenges due to the integration of complex, high-dimensional data, including shape, spatial relationships, and physical properties. Traditional approaches reliant on costing manual annotations and inflexible rules, prompting the adoption of MCoT to decompose intricate tasks into manageable and structured processes.

Several frameworks exemplify MCoT's efficacy in 3D generation. 3D-PreMise [124] employs MCoT to direct LLMs in generating 3D shapes and programming parameters, streamlining object synthesis. Similarly, L3GO [125] introduces Chain-of-3D-Thought, enabling 3D image generation through iterative trial-and-error and tool invocation within simulated environments, enhancing adaptability and precision. Gen2Sim [61] leverages MCoT to advance robot skill learning by generating 3D assets as MCoT inputs, subsequently prompting LLMs to produce task descriptions and reward functions. This approach reduces human intervention while facilitating scalable and diverse task acquisition in simulations.

When meeting language instructions, CoT3DRef [33] breaks down complex 3D grounding into interpretable steps. Meanwhile, 3D-CoT [126] improves 3D vision-language alignment by integrating MCoT with a dataset that includes structural reasoning annotations, covering aspects such as shape recognition, functional inference, and causal reasoning. Together, these advancements highlight the crucial role of MCoT for efficiently addressing complex and compositional 3D tasks.

3.4 MCoT Reasoning over Audio and Speech

MCoT has been effectively extended to step-wise and manageable speech and audio processing, mitigating the gap between the waveform signal and language semantics. CoT-ST [115] breaks the speech translation into discrete stages of speech recognition and subsequent translation. Xie et al. [116] integrate automatic speech recognition and emotion detection prior to empathetic dialogue generation. Audio-CoT [118] incorporates vanilla CoT into audio understanding and reasoning tasks. Furthermore, Audio-Reasoner [119] achieves the first long-MCoT reasoning by integrating a four-step structured reasoning framework (*e.g.*, Planning, Captioning, Reasoning, Summary). For the generative task, SpatialSonic [117] employs vanilla MCoT to derive pertinent attributes and captions, supporting the creation of a spatial audio generation. SpeechGPT-Gen [114] further introduces the Chain-of-Information-Generation approach, which systematically models semantic and perceptual information in sequential steps to facilitate natural speech generation. These developments highlight the adaptability and efficacy of MCoT in enhancing speech and audio processing, fostering more natural and contextually responsive outcomes.

Rationale Construction Input MM Input MM Input MM Input MM Prompt-based (i.e., tree, graph) Output Output

Figure 6: MCoT reasoning methods under different *rationale construction* perspectives.

3.5 MCoT Reasoning over Table and Chart

LLMs demonstrate proficiency in document comprehension but face challenges with structured data, such as tables and charts, due to their intricate layouts and implicit patterns. LayoutLLM [120] advances document understanding by integrating layout-aware pretraining at the document, region, and segment levels, employing vanilla MCoT to enhance processing. Similarly, Dai et al. [220] incorporate scene graphs to interpret charts, utilizing vanilla MCoT to mitigate hallucination in LLMs' responses. While these efforts provide coarse-level analysis, recent methodologies address such limitations by decomposing tasks into sequential, executable operations. TableGPT [121] introduces the Chain-of-Command approach, leveraging command sets (e.g., SelectCondition and GroupBy) to systematically process tabular questions. Wang et al. [122] propose Chain-of-Table, enabling LLMs to dynamically generate requisite operations and arguments, reconstructing tables to retain only pertinent information. In contrast, ReFocus [123] emulates human attention by producing visual thoughts through editing operations, such as adding highlights or masking regions in tables, thereby improving comprehension. These advancements collectively illustrate the efficacy of MCoT in navigating the complexities of structured data.

3.6 Cross-modal CoT Reasoning

CoT-like reasoning also excels in integrating multiple modalities beyond text and a single additional modality, *i.e.*, cross-modal CoT reasoning. AVQA-CoT [32] decomposes complex queries into simpler sub-questions, addressing audio-visual question answering (AVQA) sequentially through LLMs and pre-trained models. Similarly, SegPref [127] employs VLLMs to detect potential sounding objects within visual scenes, subsequently combining text rationales with mask decoders for audio-visual segmentation (AVS), thereby reducing over-reliance on visual features. In parallel, Chain-of-Empathy [128] leverages vanilla MCoT alongside psychotherapeutic principles to enhance LLMs' reasoning about human emotions, promoting empathetic responses. Likewise, MM-PEAR-CoT [129] applies vanilla MCoT to analyze linguistic emotion, integrating it with audio and video inputs to improve multimodal emotion recognition and mitigate hallucinations. R1-Omni [130] presents the first application of Reinforcement Learning with Verifiable Reward (RLVR) to an Omni-multimodal LLM in the context of emotion recognition. Although cross-modal CoT reasoning predominantly relies on text-based rationales, these advancements in MCoT demonstrate superior performance across diverse downstream tasks.

4 Methodologies in MCoT Reasoning

To thoroughly investigate the robust reasoning capabilities of MCoT in multimodal contexts, the research community has developed a wide array of methods and strategies centered on MCoT. To provide a systematic and comprehensive analysis, we categorize these approaches from multiple perspectives: rationale construction, structural reasoning, information enhancing, objective granularity, multimodal rationale, and test-time scaling.

4.1 From Rationale Construction Perspective

This part summarizes the methodologies employed in constructing rationales for MCoT reasoning. Unlike traditional direct input-output approaches, which prioritize final answers, CoT and MCoT emphasize deriving the correct answer through the reasoning process. Hence, MCoT reasoning methodologies primarily concern the construction of rationales and can be categorized into three distinct types: prompt-based, plan-based, and learning-based methods, as shown in Figure 6.

Prompt-based Methods. Prompt-based MCoT reasoning employs carefully designed prompts, including instructions or in-context demonstrations, to guide models in generating rationales during inference, typically in zero-shot or few-shot settings. For instance, the simplest instruction is "think step-by-step to understand the given text and image inputs" serves as a zero-shot prompt [110] to elicit a rationale for addressing multimodal problems. However, most MCoT approaches specify explicit steps to ensure the reasoning follows specific guidance [31, 51, 62, 69, 89, 104, 116]. In addition, expert tools are often integrated to deepen insights into detailed information [72, 76] or to incorporate multimodal data into verbal reasoning [62, 77, 217], particularly in image and video understanding. In few-shot scenarios, prompts may include explicit reasoning examples to further steer the reasoning process [108, 109, 113]. This prompt-based methodology demonstrates notable flexibility, making it advantageous for scenarios where computational resources are constrained or swift responses are essential.

Plan-based Methods. Plan-based MCoT reasoning enables models to dynamically explore and refine thoughts during the reasoning process. MM-ToT [131] utilizes GPT-4 [181] and Stable Diffusion [221] to generate multimodal outputs, applying DFS and BFS to select optimal outputs based on a 0.0–1.0 metric scale. HoT [57] produces interconnected thoughts from multimodal inputs, encapsulated within a single hyperedge. In contrast, Aggregation Graph-of-Thought (AGoT) [52] constructs a reasoning aggregation graph, integrating multiple reasoning facets at each reasoning step and subsequently incorporating visual data. Blueprint Debate on Graph (BDoG) [73] adopts a distinctive approach, forgoing search algorithms in favor of three agents—an affirmative debater, a negative debater, and a moderator. Through iterative debate, these agents address multimodal questions, with the moderator synthesizing a final answer, implicitly forming a graph-of-thought that explores and aggregates diverse thoughts. PARM++ [34] trains an image generation model with chained verification steps, such as potential assessment, to filter out bad outputs during the image generation. In summary, unlike prompt-based methods with their linear, example-driven inference, plan-based MCoT variants enable models to traverse multiple reasoning pathways, enhancing adaptability and problem-solving depth.

Learning-based Methods. Learning-based MCoT reasoning embeds rationale construction within the training or fine-tuning process, requiring models to explicitly learn reasoning skills along-side multimodal inputs. Multimodal-CoT [29] pioneers this approach by fine-tuning models with reasoning data containing rationales, fostering inherent reasoning capabilities. PCoT [55] refines this paradigm for rationale generation, while MC-CoT [56] incorporates multimodal consistency and majority voting during training to enhance reasoning in smaller models. G-CoT [80] employs ChatGPT to produce reasoning data, activating reasoning potential transferable to autonomous driving via fine-tuning. LoT [65] boosts creativity through fine-tuning with leap-of-thought data, and PromptCoT [53] enhances prompt generation for image synthesis through targeted fine-tuning. In summary, Learning-based methods focus on embedding reasoning patterns during training. However, following the release of OpenAI o1 [216] in late 2024, interest has surged in augmenting long-CoT reasoning with scaled test-time computation [91, 137, 222, 223] to tackle complex reasoning tasks, which we will further discuss in the Section 4.6.

4.2 From Structural Reasoning Perspective

Beyond simplistic rationale generation based on next token prediction, recent studies have proposed structural reasoning frameworks to enhance the controllability and interpretability of the rationale generation process. Figure 7 demonstrates the structured formats categorized into three types: asynchronous modality modeling, defined procedure staging, and autonomous procedure staging.

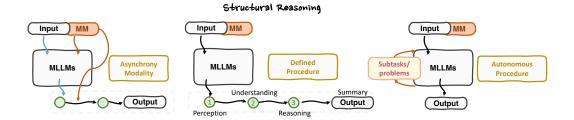


Figure 7: MCoT methods under different *structural reasoning* perspectives.

Asynchronous Modality Modeling. Early research explorations in MCoT generate rationales directly from multimodal contexts, as exemplified by Multimodal-CoT [29], Audio-CoT [118] and Grounding-Prompter [109]. However, neuroscience research by Wu et al. [224] suggest that recognition and reasoning operate in distinct cognitive modules, adhering to a "description then decision" strategy. This insight has motivated asynchronous approaches to modality processing. For instance, IPVR [51] introduced a three-stage "see, think, confirm" framework for VQA, decoupling perception from reasoning. Visualization-of-Thought [217] simulates mental imagery by generating 2D grid-based text representations to guide search and navigation tasks. Similarly, TextCoT [74] employs a two-phase process: first summarizing the image context, then generating responses grounded in visual inputs. Cantor [76] separates perception and decision-making stages, where the perception phase extracts low-level attributes (*e.g.*, objects, colors, shapes) from images or textual descriptions, while the decision phase integrates these features for accurate problem-solving. In contrast, VIC [89] decomposes tasks into text-based sub-steps before incorporating visual inputs to derive final rationales. These methods align with neuroscience by isolating perceptual encoding from high-level reasoning, thereby enhancing interpretability and alignment with human cognitive processes.

Defined Procedure Staging. Several studies explicitly define structured reasoning stages to enhance process interpretability. BDoG [73] employs a fixed debate-summarization pipeline with specialized agents, while Det-CoT [72] formalizes VQA reasoning into templated instruction parsing, subtask decomposition, execution, and verification. VisualSketchpad [77] structures rationales into "Thought, Action, Observation" phases, whereas CoTDet [58] implements object detection through object listing, affordance analysis, and visual feature summarization. Socratic Questioning [96] decomposes VQA into self-guided subquestion generation, detailed captioning, and summarization. The Grounding-Prompter [109] performs global understanding, noise assessment, and partition understanding before the final decision-making. For multi-image comprehension, CoCoT [67] systematically compares similarities and differences across inputs. LLaVA-CoT [225] achieves long-MCoT reasoning via summary, captioning, analysis, and conclusion phases, and Audio-Reasoner [119] also implements in the same manner.

Structured staging also facilitates dataset construction and downstream applications development. URSA [98] generates mathematical reasoning datasets via rationale distillation and trajectory rewriting. Paralleled by Chain-of-Sentiment [155], Chain-of-Exemplar [83] and Chain-of-Empathetic [128] in educational and affective computing domains. SmartAgent [144] builds personal assistants through GUI navigation, reasoning, and recommendation stages. CoT-ST [115] combines speech recognition and machine translation for speech translation. SegPref [127] robustly locates the sounding objects in the visual space by leveraging global understanding, sounding object filtering, and noise information removal. In generative tasks, PARM [34] generates images with clarified judgment, potential assessment, and best-of-N selection, while SpeechGPT-Gen [114] synthesizes speech from a perceptual to a semantic perspective.

Autonomous Procedure Staging. Recent studies have explored autonomous procedural staging, enabling LLMs to self-determine the sequence of reasoning steps. PS-CoT [79] allows LLMs to autonomously generate task-solving plans before rationale generation, while DDCoT [59] and AVQA-CoT [32] decompose problems into subquestions for iterative resolution. CoT-PT [132] employs hierarchical reasoning from abstract to concrete concepts (e.g., object \rightarrow animal \rightarrow dog). Image-of-Thought [78] automatically segments VQA tasks into subtasks with corresponding image ma-

Input MM Input MM Input MM Relation Modeling (e.g., counter) MLLMs Expert Tools MLLMs World Knowledge MLLMs Internal Knowledge Output Output

Figure 8: MCoT reasoning under perspectives with *information enhancing*.

nipulation actions. Insight-V [91] dynamically determines the focus of each reasoning step and autonomously decides whether to proceed or summarize intermediate results. Chain-of-Table [122] generates stepwise queries to modify table structures (*e.g.*, adding a "country" header), synthesizes operation arguments, and optimizes data storage for efficient answer derivation. In embodied intelligence tasks, E-CoT [41] and Emma-X [140] enable LLMs to infer executable subtask sequences.

4.3 From Information Enhancing Perspective

Enhancing multimodal inputs facilitates comprehensive reasoning through the integration of expert tools and internal or external knowledge.

Using Expert Tools. Recent studies leverage specialized tools to enhance multimodal reasoning through structured visual or geometric operations. For mathematical and geometric tasks, approaches such as Chain-of-Image [62] and VisualSketchpad [77] generate auxiliary visualizations via expert tools or codes. Similarly, Det-CoT [72], Cantor [76], and Image-of-Thought [78] employ image manipulation tools (*e.g.*, zoom-in, ruler markers) to improve fine-grained visual analysis. In parallel, L3GO [125] and 3D-Premise [124] integrate 3D generation tools to support spatial reasoning workflows. These methodologies underscore the growing emphasis on integrating domain-specific toolkits to augment both interpretability and precision in multimodal reasoning tasks.

Using World Knowledge Retrieval. Recent studies augment reasoning processes by integrating external knowledge sources. Approaches such as RAGAR [75], AR-MCTS [93], and Chain-of-Action [71] leverage retrieval-augmented generation (RAG) to incorporate domain-specific or commonsense knowledge during inference. G-CoT [80] distills task-relevant commonsense information from ChatGPT, while CoTDet [58] retrieves object affordances to provide context for detection tasks. KAM-CoT [68] jointly reasons over images, textual data, and structured knowledge graphs to enhance multimodal comprehension. These methods demonstrate the critical role of knowledge-aware architectures in bridging perceptual inputs with conceptual understanding.

Leveraging In-context Knowledge Retrieval. Beyond external knowledge augmentation, several studies improve reasoning by retrieving and organizing information directly from input content or generated rationales from LLMs/MLLMs themselves. DCoT [84] focuses on prioritizing image regions of interest during inference. In contrast, MCoT-Memory [133], MGCoT [134], Video-of-Thought [31], CCoT [64], and BDoG [73] implicitly retrieve in-context knowledge by modeling relationships between objects or concepts through scene graph representations. Similarly, CoT3DRef [33] generates target anchors in grounding reference sentences, effectively functioning as simplified scene graphs. Together, these approaches demonstrate the effectiveness of structured in-context knowledge extraction in enhancing reasoning fidelity.

4.4 From Objective Granularity Perspective

Research methodologies often align with the granularity of objective, as illustrated in Figure 9. While most question-answering tasks emphasize overview information, *i.e.*, coarse understanding, some fine-grained tasks such as grounding place greater importance on individual instances, *i.e.*, semantic grounding and fine-grained understanding.

Objective Granularity

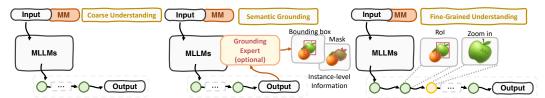


Figure 9: MCoT reasoning under the perspectives of various *objective granularities*.

Multimodal Rationale

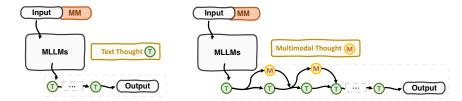


Figure 10: MCoT reasoning with multimodal rationale.

Coarse Understanding Level. As the most widely explored information processing level, coarse understanding is commonly used in tasks such as VQA and AQA, exemplified by methods like Multimodal-CoT [29] and Audio-CoT [118]. These approaches aim to achieve an overview of the given multimodal information without focusing on the details.

Semantic Grounding Level. Semantic grounding tasks are addressed through specialized reasoning paradigms. CPSeg [60], CoTDet [58], and Migician [135] refine grounding references via LLMs to enhance the alignment between textual prompts and target visual instances, improving precision in downstream mask decoders or bounding box proposers. Similarly, SegPref [127] utilizes Visual Large Language Models (VLLMs) to infer potential sounding objects from global scene information, then incorporates audio information to locate the sounding objects within the visual space.

Fine-grained Understanding Level. Fine-grained understanding also necessitates capturing detailed information within multimodal contexts. DCoT [84], TextCoT [74], and Chain-of-Spot [70] first focus on regions of interest within images, followed by the identification of fine-grained information from these retrieved areas. E-CoT [41] retrieves bounding boxes of target objects and identifies gripper positions for robotic tasks, facilitating embodied interactions.

4.5 From Multimodal Rationale Perspective

As introduced in the preliminary Section §2, reasoning processes can adopt either text-only or multimodal rationales. The predominant focus has been on text-centric approaches, exemplified by Multimodal-CoT [29], PCoT [55], MC-CoT [63], LLaVA-CoT [225], and Grounding-Prompter [109]. These approaches predominantly employ textual representation to encode multimodal information, enabling seamless integration with the reasoning mechanisms of LLMs or MLLMs.

Emerging methods, however, explore multimodal rationale construction inspired by human cognitive processes. For instance, ReFocus [123] simulates visual attention by overlaying highlighted regions on tabular data, while Visual-CoT [54] addresses logical gaps in sequential image reasoning by generating intermediate imaginary states. Chain-of-Image [62] emulates human tool-assisted reasoning by generating auxiliary diagrams for mathematical or geometric problem-solving, and Image-of-Thought [78] combines textual and visual retrieval to dynamically localize objects which should be referenced to answer a fine-grained question. Visualization-of-Thought [217] constructs 2D grids to represent human mental images for solving spatial reasoning tasks, whereas MVoT [30] further visualizes each reasoning step. Similarly, CoTDiffusion [136] leverages diffusion models to

Test-time Scaling Input MM Input MM SFT with Long SFT with Long MCoT RL Reinforcement MIIMs MIIMs MIIMs MIIMs MIIMs Internal slow thinking Output Output External slow thinking

Figure 11: MCoT reasoning with *test-time scaling* strategies. RL can help improve reasoning quality, or active long-CoT reasoning ability without annotated long-CoT training data. SFT is optional.

decompose robotic manipulation tasks into coherent visual subgoal plans, bridging abstract reasoning and physical execution. This progression from text-centric to multimodal rationales reflects an increasing emphasis on emulating human-like cognitive mechanisms.

4.6 From Test-Time Scaling Perspective

The reasoning response of LLMs can be categorized into direct responses and CoT responses, analogous to the two distinct reasoning systems present in human cognition [226]. System 1, characterized by rapid, heuristic-driven decision-making, contrasts with System 2, which employs deliberate, logical reasoning to yield more accurate and less biased outcomes [227]. Snell et al. [222] further substantiated the efficacy of "slow thinking" in LLMs that, based on System-2, scaling test-time computation optimally during inference may outperform scaling model parameters in efficiency. The release of OpenAI's o1 [216] has further fueled interest in large-scale reasoning models that combine both internal and external slow-thinking mechanisms [228–231], offering potential solutions to complex challenges, particularly in fields like mathematics and coding, while Deepseek-R1 [137] showcases that reinforcement learning (RL) alone can awaken the long-CoT reasoning ability.

Slow-Thinking-based Models. Internal slow-thinking enhances reasoning depth and quality through training or finetuning. Conversely, external slow-thinking improves reasoning with iterative sampling and refining solutions during inference. As a pioneering work, Qwen-QwQ [232] undergoes supervised fine-tuning (SFT) with 7,000 long CoT samples, unlocking the long-CoT reasoning capability. Macro-o1 [233] also integrates internal slow-thinking via CoT fine-tuning but further employs Monte Carlo Tree Search (MCTS), a heuristic search algorithm, to activate the external slow-thinking capability.

Building on these foundations, research has extended into large multimodal reasoning models. Visual-o1 [87] addresses ambiguous instructions using a multimodal and multi-turn CoT framework, while LLaVA-CoT [225] and Audio-Reasoner [119] achieve long-McoT reasoning with SFT and structural reasoning. LlamaV-o1 [97] active the long-McoT reasoning ability through a curriculum learning approach. Virgo [95] constructs a multimodal slow-thinking system by fine-tuning an MLLM with a compact set of long-form textual data. AR-MCTS [93] dynamically retrieves multimodal insights during MCTS expansion to enrich sampling diversity and reliability. AStar [100] distills reasoning patterns from only 500 data samples via MCTS to steer reasoning, whereas Mulberry [94] enhances tree search with collective learning from multiple MLLMs, leveraging negative paths to generate reflective data for improved self-reflection. RedStar [234] demonstrates that a few thousand samples suffice to activate long-CoT capabilities, with efficacy scaling alongside model size, and can enhance the generalization capability of MLLMs. These developments underscore the transformative potential of slow-thinking paradigms in advancing multimodal reasoning capabilities.

Reinforcement Learning-based Models. RL has demonstrated significant efficacy in advancing the reasoning capabilities of LLMs. Deepseek-R1 [137] illustrates this by activating long-CoT reasoning using RL alone, with further improvements surpassing OpenAI o1 in some aspects via SFT cold starts and iterative self-improvement, sparking interest in RL-driven models like Open-R1 [235] and TinyZero [236]. In multimodal contexts, early efforts like LLaVA-Reasoner [138] and Insight-V [91] refined reasoning capabilities by employing long-MCoT data fine-tuning and direct

Model	Foundational LLMs	Modality	Learning	Cold Start	Algorithm	Aha-moment
Deepseek-R1-Zero [137]	Deepseek-V3	T	RL	X	GRPO	~
Deepseek-R1 [137]	Deepseek-V3	T	SFT+RL	/	GRPO	-
LLaVA-Reasoner [138]	LLaMA3-LLaVA-NEXT-8B	T,I	SFT+RL	V	DPO	-
Insight-V [91]	LLaMA3-LLaVA-NEXT-8B	T,I	SFT+RL	V	DPO	-
Multimodal-Open-R1 [99]	Qwen2-VL-7B-Instruct	T,I	RL	Х	GRPO	Х
R1-OneVision [101]	Qwen2.5-VL-7B-Instruct	T,I	SFT	-	-	-
R1-V [237]	Qwen2.5-VL	T,I	RL	Х	GPRO	×
VLM-R1 [238]	Qwen2.5-VL	T,I	RL	X	GPRO	X
LMM-R1 [239]	Qwen2.5-VL-Instruct-3B	T,I	RL	X	PPO	×
Curr-ReFT [244]	Qwen2.5-VL-3B	T,I	RL+SFT	X	GPRO	-
Seg-Zero [245]	Qwen2.5-VL-3B + SAM2	T,I	RL	Х	GPRO	×
MM-Eureka [246]	InternVL2.5-Instruct-8B	T,I	SFT+RL	V	RLOO	-
MM-Eureka-Zero [246]	InternVL2.5-Pretrained-38B	T,I	RL	Х	RLOO	~
VisualThinker-R1-Zero [247]	Qwen2-VL-2B	T,I	RL	X	GPRO	✓
Easy-R1 [240]	Qwen2.5-VL	T,I	RL	Х	GRPO	-
Open-R1-Video [243]	Qwen2-VL-7B	T,I,V	RL	X	GRPO	X
R1-Omni [130]	HumanOmni-0.5B	T,I,V,A	SFT+RL	✓	GRPO	-
VisRL [248]	Qwen2.5-VL-7B	T,I	SFT+RL	✓	DPO	-
R1-VL [249]	Qwen2-VL-7B	T,I	RL	X	StepGRPO	-

Table 2: Multimodal reasoning models utilizing reinforcement learning. Deepseek-R1 serves as a text-only LLM for comparison.

preference optimization (DPO) guided by human feedback. Then advancing based on the Deepseek-R1, Multimodal-Open-R1 [99] integrates the GPRO framework [137] to develop an R1-like MLLM, while R1-V [237] verifies that RL can enhance the generalization in visual reasoning tasks. The effectiveness of RL in visual reasoning is also verified by concurrent works such as R1-OneVision [101], VLM-R1 [238], LMM-R1 [239], and Easy-R1 [240]. In addition, compared to the aforementioned outcome reward models (ORM), progress reward models (PRM) such as MSTaR [241] and VisualPRM [242] assess and provide feedback at each reasoning step, further enhancing the self-consistency and self-evolving capabilities of MLLMs.

This RL-based reasoning paradigm further extends to video reasoning tasks with Open-R1-Video [243], detection tasks with Curr-ReFT [244], segmentation tasks with Seg-Zero [245], and multimodal emotion recognition tasks incorporating audio in R1-Omni [130]. Moreover, RL fosters the emergence of "aha-moments" that enables the reflection and backtracking during the reasoning, which is first identified by Deepseek-R1 in text-only scenarios. MM-Eureka [246] and VisualThinker-R1-Zero [247] further reproduce the phenomena successfully in visual reasoning. Table 2 concludes the techniques used by MLLMs with RL for better long-MCoT reasoning. In summary, RL unlocks complex reasoning and "aha-moment" without SFT, demonstrating its potential to enhance model capabilities through iterative self-improvement and rule-based approaches, ultimately paving the way for more advanced and autonomous multimodal reasoning systems.

5 Applications with MCoT Reasoning

MCoT's robust capability to decompose complex tasks into manageable subtasks has facilitated its application across diverse domains, including embodied systems, agents, autonomous driving, healthcare innovations, and multimodal generation frameworks. Each domain exemplifies how MCoT reasoning enhances task decomposition, decision-making, and generalization, thereby offering significant insights into its transformative potential in addressing real-world AI challenges.

5.1 Embodied AI

Recent advancements in embodied AI have significantly enhanced robotic capabilities across planning, manipulation, and navigation. EmbodiedGPT [39] and E-CoT [41] utilize MCoT reasoning to segment tasks into actionable subgoals. Notably, EmbodiedGPT introduces the EgoCoT dataset

for vision-language pre-training, while E-CoT focuses on the sequential execution of textual commands. ManipLLM [139] enhances manipulation via fine-tuned MLLMs for object-centric tasks, while CoTDiffusion [136] employs diffusion-generated visual subgoals to achieve precision in long-horizon activities. In spatial reasoning, Emma-X [140] integrates grounded planning and predictive movement, while SpatialCoT [141] uses coordinate alignment for complex spatial reasoning. In navigation, MCoCoNav [142] optimizes multi-robot coordination with a global semantic map and score-based collaboration, whereas MCoT-Memory [133] improves long-horizon planning by incorporating memory retrieval and scene graph updates, retaining high-confidence experiences for robust decision-making. Collectively, these studies underscore a trend toward integrating multimodal data and chained reasoning for adaptable, generalizable embodied systems.

5.2 Agentic System

Advancements in AI-driven agent systems have expanded autonomous interaction and content generation capabilities. Auto-GUI [143] employs Multimodal Chain-of-Action (MCoA) to manipulate graphical interfaces directly, enhancing efficiency without reliance on external tools or APIs. Similarly, SmartAgent [144] integrates GUI navigation with Chain-of-User-Thought (CoUT) reasoning to provide personalized recommendations for embodied agents. In video understanding, VideoAgent [104] leverages LLMs with a reflective three-step decision process for accurate interpretation of long-form content. Complementing these, DreamFactory [111] pioneers long-video generation through a multi-agent framework, ensuring scene consistency via keyframe iteration and MCoT reasoning. These studies collectively illustrate the pivotal role of chaining mechanisms and agent collaboration in addressing complex real-world AI challenges.

In addition, a significant paradigm shift in AI agent systems has emerged recently, integrating "perception-reasoning" with "planning-execution." The advent of Manus [250] exemplifies this transition, igniting interest in tool-use AI agents like OpenManus [251]. Leveraging LLMs for natural language understanding and generation, Manus iteratively refines solutions through goal-directed self-reflection. As a tool-use agent, it incorporates diverse functionalities, such as web search, data querying, and code execution, employing a chain-of-tools approach to address complex multimodal tasks in the real world. Future advancements in tool-use agents are expected to build upon foundational models, enhancing long-MCoT reasoning and integrating varied multimodal interfaces and tools. This developmental trajectory suggests a progression toward AI agents with increasingly human-like capabilities.

5.3 Autonomous Driving

Recent advancements in autonomous driving have increasingly leveraged MLLMs and MCoT reasoning to improve decision-making and adaptability. DriveCoT [145] integrates MCoT into end-to-end driving systems, supported by a tailored dataset, while PKRD-CoT [69] employs zero-shot MCoT prompting to address perception, knowledge, reasoning, and decision-making in dynamic environments. Human interaction is emphasized by Ma et al. [147] and Cui et al. [146], who incorporate LLMs to interpret feedback and verbal instructions effectively. Sce2DriveX [148] enhances end-to-end control through multimodal scene understanding and demonstrates robust generalization. Furthermore, Reason2Drive [149] provides 600K+ video-text pairs to explore interpretable reasoning, augmenting LLMs with object-level perception to strengthen planning capabilities. Collectively, these efforts indicate a transition toward human-like reasoning, enhanced interactivity, and improved generalization within autonomous driving systems.

5.4 Medical and Healthcare

The innovative applications of AI in healthcare have harnessed chained reasoning to enhance various medical tasks. StressSelfRefine [151] detects stress in videos through a psychology-inspired "Describe, Assess, Highlight" process, refined via DPO to enhance accuracy. TI-PREGO [113] integrates ICL with Automatic Chain-of-Thought (ACoT) to identify procedural errors in egocentric videos, leveraging action sequences and logical reasoning. Chain-of-Look [152] tackles surgical triplet recognition in endoscopic videos by breaking down tasks into video reasoning stages using vision-language prompts. Meanwhile, MedCoT [153] improves medical visual question answering through a hierarchical expert system that culminates in a Mixture-of-Experts diagnosis. Furthermore, MedVLM-R1 [154] employs RL with only 600 medical VQA samples, aiming to enhance

the medical reasoning capabilities of vision-language models. Collectively, these efforts illustrate the effectiveness of MCoT reasoning in enhancing interpretability and precision across a range of medical AI applications.

5.5 Social and Human

MCoT has been effectively extended to the humanities and social sciences, leveraging their capacity for complex task decomposition. For instance, Chain-of-Empathetic [128] employs MCoT to facilitate empathetic dialogue generation. MM-PEAR-CoT framework [129] enhances multimodal sentiment analysis through a structured Preliminaries-Question-Answer-Reason approach, generating textual rationales prior to late-stage multimodal fusion. Advancing the affective computing domain, Chain-of-Sentiment [155] refines sentiment analysis in conversational contexts, with complementary contributions from concurrent studies [218, 219]. Furthermore, Chain-of-Exemplar [83] expands MCoT into the field of education, X-Reflect [252] extends MCoT into multimodal recommendation systems, while Yu and Luo [253] employs zero-shot MCoT for demographic inference. These advancements highlight the potential of MCoT to tackle complex challenges in human-centric and social scientific contexts.

5.6 Multimodal Generation

Recent advancements in AI-driven image and 3D generation highlight diverse multimodal synthesis strategies. PARM and PARM++ [34] employ iterative step-by-step reasoning with clarity and potential assessments, augmented by reflection mechanisms, to produce high-quality images. GoT [254] constructs a generation chain of thought by introducing a clear language reasoning process that analyzes semantic relationships and spatial arrangements before generating and editing images. RPG-DiffusionMaster [156] utilizes MLLMs for text-to-image diffusion, breaking down prompts into detailed subregions for coherent outputs, while L3GO [125] leverages language agents with a Chain-of-3D-Thoughts approach to create unconventional 3D objects, outperforming traditional diffusion models on out-of-distribution descriptions. Additionally, 3D-PreMise [124] integrates LLMs with program synthesis to generate parametric 3D shapes, yielding promising industrial outcomes when guided by explicit reasoning examples. These studies collectively underscore the potential of reasoning-augmented AI to overcome the limitations of data-driven generation, enabling precise and innovative multimodal outputs.

6 MCoT Datasets and Benchmarks

MCoT reasoning necessitates specialized datasets and benchmarks to support both model finetuning and performance evaluation. Incorporating with Table 3, this section surveys the landscape of MCoT-related resources, categorized into two key areas: datasets designed for finetuning MLLMs with reasoning rationales, and benchmarks developed to assess downstream capabilities, with or without accompanying rationales. These resources collectively address the diverse needs of training and evaluating MLLMs across multiple domains, modalities, and reasoning complexities.

6.1 Datasets for MLLMs Finetuning with Rationale

Several studies explore to activate the MCoT reasoning capabilities of MLLMs through specific datasets. ScienceQA [157] presents multimodal science questions paired with annotated answers, lectures, and explanations, illustrating how language models can leverage MCoT to enhance multihop reasoning. A-OKVQA [158] offers essential data for VQA by providing MLLMs with extensive commonsense and world knowledge. Building on ScienceQA, T-SciQ [55] enriches reasoning rationale through advanced LLMs. VideoCoT [159] supplies reasoning data for step-by-step video question answering (VideoQA), although its reasoning is limited to textual explanations. In contrast, VideoEspresso [160] provides VideoQA pairs that maintain spatial and temporal coherence along with multimodal reasoning annotations. Additionally, the MAVIS [161] dataset propels the training and application of MLLMs in mathematics by automating the generation of mathematical visual data, thus offering a wealth of aligned vision-language pairs and reasoning rationales. Ego-CoT [39] and EMMA-X [140] propose an egocentric dataset for the training model to execute the embodied tasks with sub-goal tasks. M3CoT [28] further advances VLLMs reasoning with multi-

Datasets	Year	Task	Domain	Modality	Format	Samples		
Training with rationale								
ScienceQA [157]	2022	VQA	Science	T, I	MC	21K		
A-OKVQA [158]	2022	VQA	Common	T, I	MC	25K		
EgoCoT [39]	2023	VideoQA	Common	T, V	Open	200M		
VideoCoT [159]	2024	VideoQA	Human Action	T, V	Open	22K		
VideoEspresso [160]	2024	VideoQA	Common	T, V	Open	202,164		
EMMA-X [140]	2024	Robot Manipulation	Indoor	T, V	Robot Actions	60K		
M3CoT [28]	2024	VQA	Science, Math, Common	T, I	MC	11.4K		
MAVIS [161]	2024	ScienceQA	Math	T, I	MC and Open	834K		
LLaVA-CoT-100k [225]	2024	VQA	Common, Science	T, I	MC and Open	834K		
MAmmoTH-VL [255]	2024	Diverse	Diverse	T, I	MC and Open	12M		
Mulberry-260k [94]	2024	Diverse	Diverse	T, I	MC and Open	260K		
MM-Verify [256]	2025	MathQA	Math	T, I	MC and Open	59,772		
VisualPRM400K [242]	2025	ScienceQA	Math, Science	T, I	MC and Open	400K		
R1-OneVision [101]	2025	Diverse	Diverse	T, I	MC and Open	155K		
		Evaluation v	vithout rationale					
MMMU [162]	2023	VQA	Arts, Science	T, I	MC and Open	11.5K		
SEED [163]	2023	VQA	Common	T, I	MC	19K		
MathVista [164]	2023	ScienceQA	Math	T, I	MC and Open	6,141		
MathVerse [165]	2024	ScienceQA	Math	T, I	MC and Open	15K		
Math-Vision [166]	2024	ScienceQA	Math	T, I	MC and Open	3040		
OSWorld [257]	2024	Agent	Real Comp. Env.	T,I	Agent Actions	369		
AgentClinic [258]	2024	MedicalQA	Medical	T,I	Open	335		
MeViS [170]	2023	Referring VOS	Common	T, V	Dense Mask	2K		
VSIBench [169]	2024	VideoQA	Indoor	T, V	MC and Open	5K		
HallusionBench [171]	2024	VQA	Common	T, I	Yes-No	1,129		
AV-Odyssey [259]	2024	AVQA	Common	T, V, A	MC	4,555		
AVHBench [173]	2024	AVQA	Common	T, V, A	Open	5,816		
RefAVS-Bench [168]	2024	Referring AVS	Common	T, V, A	Dense Mask	4,770		
MMAU [260]	2024	AQA	Common	T, A	MC	10K		
AVTrustBench [172]	2025	AVQA	Common	T, V, A	MC and Open	600K		
MIG-Bench [135]	2025	Multi-image Grounding	Common	T, I	BBox	5.89K		
MedAgentsBench [261]	2025	MedicalQA	Medical	T, I	MC and Open	862		
		Evaluation	with rationale					
CoMT [174]	2024	VQA	Common	T, I	MC	3,853		
OmniBench [179]	2024	VideoQA	Common	T, I, A	MC	1,142		
WorldQA [175]	2024	VideoQA	Common	T, V, A	Open	1,007		
MiCEval [176]	2024	VQA	Common	T, I	Open	643		
OlympiadBench [177]	2024	ScienceQA	Maths, Physics	T, I	Open	8,476		
MME-CoT [178]	2025	VQA	Science, Math, Common	T, I	MC and Open	1,130		
EMMA [167]	2025	VQA	Science	T, I	MC and Open	2,788		
VisualProcessBench [242]	2025	ScienceQA	Math, Science	T, I	MC and Open	2,866		

Table 3: Datasets and Benchmarks for MCoT Training and Evaluation. "MC" and "Open" refer to multiple-choice and open-ended answer formats, while "T", "I", "V", and "A" represent Text, Image, Video, and Audio, respectively.

domain and multi-hop reasoning samples. Meanwhile, MAmmoTH-VL-Instruct [255] constructs 12 million long-MCoT reasoning data across 118 datasets and 10 categories, advancing the long-MCoT reasoning capabilities of MLLMs. In addition, datasets such as LLaVA-CoT-100k [225], Mulberry-260k [94], MM-Verify [256] and VisualPRM400K [242] are commonly proposed by their corresponding reasoning models to active the long-MCoT reasoning abilities. These datasets are commonly linked with learning-based rationale construction as we mentioned in Section 4.1.

6.2 Benchmarks for Downstream Capability Assessment

A variety of benchmarks have been developed to evaluate downstream capabilities, particularly in the domains of commonsense and scientific reasoning. As illustrated in Table 4, we present a performance comparison of MLLMs from various institutions across four benchmarks: MMMU [162], MathVista [164], Math-Vision [166], and EMMA [167]. While MMMU and EMMA focus on multidisciplinary scenarios, MathVista and Math-Vision primarily assess mathematical reasoning.

Datasets without Rationale. Several multimodal evaluation benchmarks have been widely adopted to assess the performance of MLLMs. Although these benchmarks do not provide rationales, their diversity and challenges suggest that, with the help of MCoT, MLLMs could further enhance their performance on these benchmarks. MMMU [162] involves visual-language questions across six core disciplines, aiming to measure the three fundamental abilities of perception, knowledge, and reasoning in LMMs. SEED [163] further introduces the video modality to assess the understanding and generation capabilities of MLLMs. MathVista [164], MathVerse [165] and Math-Vision [166] specifically focuses on visual perception and reasoning in the mathematical domain. Emma [167] introduces reasoning challenges that cannot be solved by independent reasoning within each modality in solo, providing a comprehensive evaluation of multimodal reasoning abilities in MLLMs.

In addition to the general and mathematical domains, MCoT demonstrates its effectiveness in various downstream tasks due to its step-by-step reasoning capabilities. Migician [135] provides an evaluation of multi-image grounding and shows the effectiveness of MCoT in such tasks. RefAVS [168] introduces grounding into the audiovisual context, incorporating temporal and spatial information, which can be further addressed through MCoT. VSIBench [169] offers an evaluation of spatial reasoning abilities in MLLMs. MeViS [170] provides an evaluation of video segmentation with motion expressions. In particular, through the step-by-step reasoning of MCoT, hallucination phenomena that arise in MLLMs are expected to be further addressed. HallusionBench [171] evaluates hallucination phenomena in VLLMs, while AVTrustBench [172] and AVHBench [173] assess hallucinations in audiovisual contexts, which can be further mitigated through MCoT reasoning. OSWorld [257] and AgentClinic [258] provide benchmarks for assessing agent capabilities in multimodal scenarios, which can be enhanced through reasoning.

Datasets with Rationale. With the emergence of OpenAI of [216] and Deepseek-rf [137], interest in scaling test-time computation and slow-thinking has steadily grown, leading to the development of evaluation benchmarks designed to assess the quality of rationales generated by MLLMs during reasoning. CoMT [174] is introduced to address the limitations of traditional multimodal benchmarks that only reasoning with language, by requiring both multimodal input and output, aiming to better mimic human-like reasoning and explore complex visual operations. WorldQA [175] challenges MLLMs to answer questions using language, vision, and audio, while incorporating long-MCoT reasoning and world knowledge. MiCEval [176] is crafted to evaluate the accuracy of reasoning chains by meticulously assessing the quality of both the descriptive component and each individual reasoning step. OlympiadBench [177] features 8000+ bilingual Olympiad-level mathematics and physics problems with annotated rationales to effectively evaluate the advanced capabilities of MLLMs. MME-CoT [178] provides a systematic evaluation of MCoT reasoning, revealing critical insights, including how reflection mechanisms enhance reasoning quality and how MCoT prompting can negatively impact perception-intensive tasks, potentially due to overthinking. OmniBench [179] is the first comprehensive evaluation benchmark involving text, vision, and audio.

7 Limitations, Challenges and Future Directions

Despite the increasing attention and research efforts devoted to MCoT, several critical aspects remain unresolved and underexplored, which might be the key bottlenecks for achieving human-level multimodal AGI. Below, we summarize some challenges to shed light on future works.

Computational Sustainability and Slow-thinking Paradox. Despite significant advancements, the reliance on extensive test-time scaling and slow-thinking [227, 230, 231] to support long-MCoT reasoning poses substantial challenges. The exponential growth in computational resources and

Model	Params (B)	MMMU (Val)	MathVista (mini)	Math-Vision	EMMA (mini)			
Human	_	88.6	60.3	68.82	77.75			
Random Choice	-	22.1	17.9	7.17	22.75			
OpenAI								
o1 [216]	-	78.2	73.9	-	45.75			
GPT-4.5 [262]	-	74.4	-	-	-			
GPT-40 [212]	-	69.1	63.8	30.39	36.00			
GPT-40 mini [212]	-	59.4	56.7	-	-			
GPT-4V [263]	-	56.8	49.9	23.98	-			
		Google & De	epMind					
Gemini 2.0 Pro [264]	-	72.7	-	-	-			
Gemini 2.0 Flash [264]	-	71.7	-	41.3	48.00			
Gemini 1.5 Pro [265]	-	65.8	63.9	19.24	-			
		Anthrop	pic					
Claude 3.7 Sonnet [194]	-	75	-	-	56.50			
Claude 3.5 Sonnet [194]	-	70.4	67.7	37.99	37.00			
Claude 3 Opus [194]	-	59.4	50.5	27.13	-			
Claude 3 Sonnet [194]	-	53.1	47.9	-	-			
		xAI						
Grok-3 [266]	-	78.0	-	-	-			
Grok-2 [267]	-	66.1	69.0	-	-			
Grok-2 mini [267]	-	63.2	68.1	-	-			
		Moonsh	not					
Kimi-k1.5 [268]	-	70	74.9 (test)	38.6	33.75			
		Alibab						
QVQ-72B-Preview [269]	72	70.3	71.4	35.9	32.00			
Qwen2.5-VL-72B [270]	72	70.2	74.8	38.1	-			
Qwen2-VL-72B [11]	72	64.5	70.5	25.9	37.25			
Qwen2.5-VL-7B [270]	7	58.6	68.2	25.1	=			
Qwen2-VL-7B [11]	7	-	-	16.3	-			
<u> </u>		OpenGV.	Lab					
InternVL2.5 [271]	78	70.1	-	-	35.25			
InternVL2 [272]	76	58.2	65.5	-	-			
		LLaM						
Llama-3.2-90B [273]	90	60.3	57.3	-	-			
Llama-3.2-11B [273]	11	-	48.6	-	-			
		LLaV	1					
LLaVA-OneVision [274]	72	56.8	67.5	-	27.25			
LlaVA-NEXT-72B [275]	72	49.9	46.6	-	-			
LLaVA-NEXT-34B [275]	34	48.1	46.5	-	-			
LLaVA-NEXT-8B [275]	8	41.7	37.5	-	-			
LLaVA-Reasoner [138]	8	40.0	50.6	-	-			
LLaVA-1.5 [276]	13	36.4	27.6	11.12	-			
		Commui	ıitv					
Mulberry [94]	7	55.0	63.1	-	-			
MAmmoTH-VL [255]	8	50.8	67.6	24.4	-			
MM-Eureka [246]	8	-	67.1	22.2	-			
MM-Eureka-Zero [246]	38	-	64.2	26.6	-			
Curr-ReFT [244]	7	-	64.5	-	-			
Curr-ReFT [244]	3	-	58.6	-	-			
LMM-R1 [239]	3	-	63.2	26.35	-			
LlamaV-o1 [97]	11	-	54.4	-	-			
R1-Onevision [101]	7	-	-	26.16	-			
Virgo [95]	7	46.7	-	24.0	-			
Insight-V [91]	8	42.0	49.8	-	-			
R1-VL [249]	7	63.5	24.7	-	-			
	•							

Table 4: Performance comparison of MLLMs from various institutions across four benchmarks: MMMU (Val), MathVista (Mini), Math-Vision, and EMMA (Mini).

training data required to sustain deep reasoning processes remains a critical bottleneck, necessitating innovations in algorithm efficiency like RL [137] and hardware acceleration.

Lack of Reasoning in General Scenarios. The existing long-MCoT framework primarily focuses on verifiable data in mathematics and science, yet it lacks robust reasoning capabilities for general scenarios. Tasks in mathematics and science, characterized by their strict logical structure and unique solutions, have been extensively explored for test-time scaling using ORM [99, 246] and

PRM [241, 242]. However, in general scenarios, answers are rarely fixed, often encompassing multiple plausible interpretations and inferences. This variability renders reasoning frameworks rooted in mathematics and science less effective. The deficiency in reasoning manifests as an inability to adequately evaluate the inference process when models address tasks in general contexts involving complex situations, ambiguity, and multifactorial influences. Future research should explore the open-ended reward models to achieve robust long-chain reasoning in multimodal general scenarios.

Error Propagation in Extended Reasoning Chains. A major concern in long-chain reasoning in current MCoT systems is the error snowballing effect [230], where small inaccuracies in early steps can amplify through subsequent stages and lead to catastrophic outcomes. Traditional confidence calibration techniques fail to address the unique challenges of multimodal error propagation where different modalities may contradict each other while maintaining high self-consistency scores. Developing quantitative metrics to diagnose, quantify, and mitigate these cumulative errors is an unresolved issue that requires rigorous investigation.

Symbolic-Neural Integration Gap. While neural models excel at pattern recognition, their inability to perform rigorous symbolic operations undermines complex reasoning. Hybrid neurosymbolic architectures [277–279] can help achieve improvement on mathematical proofs, but might still suffer from knowledge grounding issues. The fundamental challenge lies in developing seamless interfaces between distributed representations and discrete symbolic systems, particularly for crossmodal symbolic manipulation (*e.g.*, converting geometric diagrams to formal proofs).

Dynamic Environment Adaptation and Adaptive Chain Length. Most MCoT systems assume static input conditions, severely limiting real-world applicability. The "frozen reasoning" paradox emerges when processing streaming multimodal inputs—most current architectures cannot revise earlier conclusions upon receiving new evidence without restarting the entire chain. Also, the development of resource-efficient reasoning frameworks that can dynamically adjust chain length or the number of reasoning steps [91] in response to computational constraints is critical. Adaptive strategies based on real-time evaluation and feedback mechanisms will be key to balancing reasoning accuracy with resource efficiency.

Hallucination Prevention. A critical challenge in existing MCoT reasoning can be mitigating hallucinations [89], where models produce plausible yet factually incorrect or inconsistent outputs. This issue undermines the reliability of the reasoning process and is especially pronounced in multimodal settings, where integrating diverse data sources can lead to misaligned contexts and spurious information. Future work should focus on robust cross-modal alignment methods and verification mechanisms at each reasoning step. Techniques such as uncertainty quantification, adversarial training, and leveraging external knowledge bases have shown promise in reducing hallucinations. Additionally, insights from human error detection may inspire novel strategies to enhance both the accuracy and interpretability of multimodal reasoning systems.

Data Selection, Annotation, and Augmentation Strategies. Recent studies have demonstrated that carefully curated datasets can activate long-MCoT reasoning capabilities in models [225, 94, 280]. However, automating the selection and annotation of data suitable for extended reasoning remains an open challenge. Leveraging semi-supervised or self-supervised learning strategies, along-side RL approaches, could reduce the reliance on extensive manual annotations.

Modality Imbalance and High-Dimensional Modal Integration. Current research indicates that progress across different modalities is uneven, with some modalities, such as text and images, advancing more rapidly than others. Future work should address the integration of higher-dimensional modalities (*e.g.*, 3D data, sensor information) to achieve a more balanced and comprehensive multimodal reasoning framework that leverages the unique characteristics of each modality.

Interdisciplinary Integration with Cognitive Science. The inherent complexity of MCoT reasoning calls for an interdisciplinary approach that integrates insights from cognitive science, psychology, and neuroscience. Drawing from human decision-making and cognitive theories [224, 151, 217] can inspire novel reasoning architectures that more closely mimic human thought processes, thereby enhancing both performance and interpretability.

Embodied Reasoning Limitations. Most MCoT systems operate in abstract symbol spaces disconnected from physical embodiment. Closing this simulation-to-reality gap necessitates tight integration of proprioceptive feedback, haptic understanding, and dynamic world modeling—challenges that current architecture designs barely address.

Explainable Reasoning and Theoretical Support. As models become increasingly complex, the interpretability of decision-making processes is essential for building trust and enabling practical deployment. Although CoT or MCoT reasoning provides intermediate steps, its underlying theoretical mechanisms remain largely opaque, leaving the model as a "black box" that simply produces predictions and outputs. Developing methodologies that offer transparent, traceable reasoning paths will not only enhance model explainability but also facilitate debugging and further optimization of MCoT reasoning systems. Strengthening theoretical support is therefore crucial to achieving truly explainable reasoning.

Ethical, Robustness and Safety Reasoning. As MCoT systems become more powerful, ensuring AI safety and robustness against adversarial perturbations is paramount. Integrating more robust reasoning techniques and approaches may enhance system transparency and provide a safety net against potential failures. Also, developing multimodal constitutional AI frameworks that can parse and apply ethical constraints across modalities becomes crucial as these systems approach real-world deployment. Further research is needed to quantify and mitigate risks associated with adversarial attacks and other safety concerns in multimodal reasoning environments.

8 Conclusion

This survey presents the first systematic review of multimodal chain-of-thought (MCoT) reasoning. We start by providing definitions and elucidated foundational concepts, laying the groundwork for methodologies. A comprehensive taxonomy categorizes diverse approaches and reasoning paradigms under various perspectives to tackle challenges unique to modalities such as image, video, speech, audio, 3D, and structured data. We consolidate comprehensive MCoT-related datasets and benchmarks to provide an accurate overview of the current resource landscape. Furthermore, our review examines relevant applications, highlighting successes in important domains including robotics, healthcare, autonomous driving, social science and multimodal generation. Lastly, we outline promising future research directions that aim to overcome these limitations and drive progress toward multimodal AGI. We openly share all related resources and information on MCoT to facilitate follow-up research in this rapidly evolving domain.

References

- [1] An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jialong Tang, Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Ma, Jianxin Yang, Jin Xu, Jingren Zhou, Jinze Bai, Jinzheng He, Junyang Lin, Kai Dang, Keming Lu, Keqin Chen, Kexin Yang, Mei Li, Mingfeng Xue, Na Ni, Pei Zhang, Peng Wang, Ru Peng, Rui Men, Ruize Gao, Runji Lin, Shijie Wang, Shuai Bai, Sinan Tan, Tianhang Zhu, Tianhao Li, Tianyu Liu, Wenbin Ge, Xiaodong Deng, Xiaohuan Zhou, Xingzhang Ren, Xinyu Zhang, Xipin Wei, Xuancheng Ren, Xuejing Liu, Yang Fan, Yang Yao, Yichang Zhang, Yu Wan, Yunfei Chu, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, Zhifang Guo, and Zhihao Fan. Qwen2 technical report. *CoRR*, abs/2407.10671, 2024.
- [2] Aohan Zeng, Bin Xu, Bowen Wang, Chenhui Zhang, Da Yin, Diego Rojas, Guanyu Feng, Hanlin Zhao, Hanyu Lai, Hao Yu, Hongning Wang, Jiadai Sun, Jiajie Zhang, Jiale Cheng, Jiayi Gui, Jie Tang, Jing Zhang, Juanzi Li, Lei Zhao, Lindong Wu, Lucen Zhong, Mingdao Liu, Minlie Huang, Peng Zhang, Qinkai Zheng, Rui Lu, Shuaiqi Duan, Shudan Zhang, Shulin Cao, Shuxun Yang, Weng Lam Tam, Wenyi Zhao, Xiao Liu, Xiao Xia, Xiaohan Zhang, Xiaotao Gu, Xin Lv, Xinghan Liu, Xinyi Liu, Xinyue Yang, Xixuan Song, Xunkai Zhang, Yifan An, Yifan Xu, Yilin Niu, Yuantao Yang, Yueyan Li, Yushi Bai, Yuxiao Dong, Zehan Qi, Zhaoyu Wang, Zhen Yang, Zhengxiao Du, Zhenyu Hou, and Zihan Wang. Chatglm: A family of large language models from GLM-130B to GLM-4 all tools. *CoRR*, abs/2406.12793, 2024.

- [3] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton-Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurélien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. Llama 2: Open foundation and fine-tuned chat models. *CoRR*, abs/2307.09288, 2023.
- [4] Alex Young, Bei Chen, Chao Li, Chengen Huang, Ge Zhang, Guanwei Zhang, Heng Li, Jiangcheng Zhu, Jianqun Chen, Jing Chang, Kaidong Yu, Peng Liu, Qiang Liu, Shawn Yue, Senbin Yang, Shiming Yang, Tao Yu, Wen Xie, Wenhao Huang, Xiaohui Hu, Xiaoyi Ren, Xinyao Niu, Pengcheng Nie, Yuchi Xu, Yudong Liu, Yue Wang, Yuxuan Cai, Zhenyu Gu, Zhiyuan Liu, and Zonghong Dai. Yi: Open foundation models by 01.ai. *CoRR*, abs/2403.04652, 2024.
- [5] Rohan Anil, Andrew M. Dai, Orhan Firat, Melvin Johnson, Dmitry Lepikhin, Alexandre Passos, Siamak Shakeri, Emanuel Taropa, Paige Bailey, Zhifeng Chen, Eric Chu, Jonathan H. Clark, Laurent El Shafey, Yanping Huang, Kathy Meier-Hellstern, Gaurav Mishra, Erica Moreira, Mark Omernick, Kevin Robinson, Sebastian Ruder, Yi Tay, Kefan Xiao, Yuanzhong Xu, Yujing Zhang, Gustavo Hernández Ábrego, Junwhan Ahn, Jacob Austin, Paul Barham, Jan A. Botha, James Bradbury, Siddhartha Brahma, Kevin Brooks, Michele Catasta, Yong Cheng, Colin Cherry, Christopher A. Choquette-Choo, Aakanksha Chowdhery, Clément Crepy, Shachi Dave, Mostafa Dehghani, Sunipa Dev, Jacob Devlin, Mark Díaz, Nan Du, Ethan Dyer, Vladimir Feinberg, Fangxiaoyu Feng, Vlad Fienber, Markus Freitag, Xavier Garcia, Sebastian Gehrmann, Lucas Gonzalez, and et al. Palm 2 technical report. CoRR, abs/2305.10403, 2023.
- [6] Marah I Abdin, Sam Ade Jacobs, Ammar Ahmad Awan, Jyoti Aneja, Ahmed Awadallah, Hany Awadalla, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Harkirat S. Behl, Alon Benhaim, Misha Bilenko, Johan Bjorck, Sébastien Bubeck, Martin Cai, Caio César Teodoro Mendes, Weizhu Chen, Vishrav Chaudhary, Parul Chopra, Allie Del Giorno, Gustavo de Rosa, Matthew Dixon, Ronen Eldan, Dan Iter, Amit Garg, Abhishek Goswami, Suriya Gunasekar, Emman Haider, Junheng Hao, Russell J. Hewett, Jamie Huynh, Mojan Javaheripi, Xin Jin, Piero Kauffmann, Nikos Karampatziakis, Dongwoo Kim, Mahoud Khademi, Lev Kurilenko, James R. Lee, Yin Tat Lee, Yuanzhi Li, Chen Liang, Weishung Liu, Eric Lin, Zeqi Lin, Piyush Madan, Arindam Mitra, Hardik Modi, Anh Nguyen, Brandon Norick, Barun Patra, Daniel Perez-Becker, Thomas Portet, Reid Pryzant, Heyang Qin, Marko Radmilac, Corby Rosset, Sambudha Roy, Olatunji Ruwase, Olli Saarikivi, Amin Saied, Adil Salim, Michael Santacroce, Shital Shah, Ning Shang, Hiteshi Sharma, Xia Song, Masahiro Tanaka, Xin Wang, Rachel Ward, Guanhua Wang, Philipp Witte, Michael Wyatt, Can Xu, Jiahang Xu, Sonali Yadav, Fan Yang, Ziyi Yang, Donghan Yu, Chengruidong Zhang, Cyril Zhang, Jianwen Zhang, Li Lyna Zhang, Yi Zhang, Yue Zhang, Yunan Zhang, and Xiren Zhou. Phi-3 technical report: A highly capable language model locally on your phone. CoRR, abs/2404.14219, 2024.
- [7] Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, Katie Millican, et al. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.
- [8] Zonghao Guo, Ruyi Xu, Yuan Yao, Junbo Cui, Zanlin Ni, Chunjiang Ge, Tat-Seng Chua, Zhiyuan Liu, and Gao Huang. Llava-uhd: An LMM perceiving any aspect ratio and highresolution images. In ECCV, pages 390–406, 2024.
- [9] Guowei Xu, Peng Jin, Hao Li, Yibing Song, Lichao Sun, and Li Yuan. Llava-cot: Let vision language models reason step-by-step. *CoRR*, abs/2411.10440, 2024.

- [10] Chunyuan Li, Cliff Wong, Sheng Zhang, Naoto Usuyama, Haotian Liu, Jianwei Yang, Tristan Naumann, Hoifung Poon, and Jianfeng Gao. Llava-med: Training a large language-and-vision assistant for biomedicine in one day. In *NeurIPS*, 2023.
- [11] Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Yang Fan, Kai Dang, Mengfei Du, Xuancheng Ren, Rui Men, Dayiheng Liu, Chang Zhou, Jingren Zhou, and Junyang Lin. Qwen2-vl: Enhancing vision-language model's perception of the world at any resolution. *CoRR*, abs/2409.12191, 2024.
- [12] 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. In CVPR, pages 26753–26763, 2024.
- [13] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances in neural information processing systems*, 36:34892–34916, 2023.
- [14] Shengqiong Wu, Hao Fei, Leigang Qu, Wei Ji, and Tat-Seng Chua. NExT-GPT: Any-to-any multimodal llm. In *International Conference on Machine Learning*, pages 53366–53397, 2024.
- [15] Hao Fei, Shengqiong Wu, Hanwang Zhang, Tat-Seng Chua, and Shuicheng Yan. Vitron: A unified pixel-level vision LLM for understanding, generating, segmenting, editing. In *Advances in neural information processing systems*, 2024.
- [16] Yunfei Chu, Jin Xu, Qian Yang, Haojie Wei, Xipin Wei, Zhifang Guo, Yichong Leng, Yuan-jun Lv, Jinzheng He, Junyang Lin, et al. Qwen2-audio technical report. *arXiv preprint arXiv:2407.10759*, 2024.
- [17] Muhammad Maaz, Hanoona Rasheed, Salman Khan, and Fahad Shahbaz Khan. Video-chatgpt: Towards detailed video understanding via large vision and language models. arXiv preprint arXiv:2306.05424, 2023.
- [18] Hang Zhang, Xin Li, and Lidong Bing. Video-llama: An instruction-tuned audio-visual language model for video understanding. *arXiv preprint arXiv:2306.02858*, 2023.
- [19] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837, 2022.
- [20] Zhuosheng Zhang, Aston Zhang, Mu Li, and Alex Smola. Automatic chain of thought prompting in large language models. *arXiv preprint arXiv:2210.03493*, 2022.
- [21] Xuezhi Wang and Denny Zhou. Chain-of-thought reasoning without prompting. In NeurIPS, 2024.
- [22] Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegreffe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, Shashank Gupta, Bodhisattwa Prasad Majumder, Katherine Hermann, Sean Welleck, Amir Yazdanbakhsh, and Peter Clark. Self-refine: Iterative refinement with self-feedback. In *NeurIPS*, 2023.
- [23] Maciej Besta, Nils Blach, Ales Kubicek, Robert Gerstenberger, Michal Podstawski, Lukas Gianinazzi, Joanna Gajda, Tomasz Lehmann, Hubert Niewiadomski, Piotr Nyczyk, et al. Graph of thoughts: Solving elaborate problems with large language models. In *AAAI*, pages 17682–17690, 2024.
- [24] Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik Narasimhan. Tree of thoughts: Deliberate problem solving with large language models. In *NeurIPS*, 2023.
- [25] Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec Helyar, Aleksander Madry, Alex Beutel, Alex Carney, Alex Iftimie, Alex Karpenko, Alex Tachard Passos, Alexander Neitz, Alexander Prokofiev, Alexander Wei, Allison Tam, Ally Bennett, Ananya Kumar, Andre Saraiva, Andrea Vallone, Andrew Duberstein, Andrew

Kondrich, Andrey Mishchenko, Andy Applebaum, Angela Jiang, Ashvin Nair, Barret Zoph, Behrooz Ghorbani, Ben Rossen, Benjamin Sokolowsky, Boaz Barak, Bob McGrew, Borys Minaiev, Botao Hao, Bowen Baker, Brandon Houghton, Brandon McKinzie, Brydon Eastman, Camillo Lugaresi, Cary Bassin, Cary Hudson, Chak Ming Li, Charles de Bourcy, Chelsea Voss, Chen Shen, Chong Zhang, Chris Koch, Chris Orsinger, Christopher Hesse, Claudia Fischer, Clive Chan, Dan Roberts, Daniel Kappler, Daniel Levy, Daniel Selsam, David Dohan, David Farhi, David Mely, David Robinson, Dimitris Tsipras, Doug Li, Dragos Oprica, Eben Freeman, Eddie Zhang, Edmund Wong, Elizabeth Proehl, Enoch Cheung, Eric Mitchell, Eric Wallace, Erik Ritter, Evan Mays, Fan Wang, Felipe Petroski Such, Filippo Raso, Florencia Leoni, Foivos Tsimpourlas, Francis Song, Fred von Lohmann, Freddie Sulit, Geoff Salmon, Giambattista Parascandolo, Gildas Chabot, Grace Zhao, Greg Brockman, Guillaume Leclerc, Hadi Salman, Haiming Bao, Hao Sheng, Hart Andrin, Hessam Bagherinezhad, Hongyu Ren, Hunter Lightman, Hyung Won Chung, Ian Kivlichan, Ian O'Connell, Ian Osband, Ignasi Clavera Gilaberte, and Ilge Akkaya. Openai o1 system card. *CoRR*, abs/2412.16720, 2024.

- [26] DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, Aixin Liu, Bing Xue, Bingxuan Wang, Bochao Wu, Bei Feng, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng Wang, Honghui Ding, Huajian Xin, Huazuo Gao, Hui Qu, Hui Li, Jianzhong Guo, Jiashi Li, Jiawei Wang, Jingchang Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, J. L. Čai, Jiaqi Ni, Jian Liang, Jin Chen, Kai Dong, Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong Zhang, Liang Zhao, Litong Wang, Liyue Zhang, Lei Xu, Leyi Xia, Mingchuan Zhang, Minghua Zhang, Minghui Tang, Meng Li, Miaojun Wang, Mingming Li, Ning Tian, Panpan Huang, Peng Zhang, Qiancheng Wang, Qinyu Chen, Qiushi Du, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang, R. J. Chen, R. L. Jin, Ruyi Chen, Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shengfeng Ye, Shiyu Wang, Shuiping Yu, Shunfeng Zhou, Shuting Pan, and S. S. Li. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. CoRR, abs/2501.12948, 2025.
- [27] Xiongtao Zhou, Jie He, Lanyu Chen, Jingyu Li, Haojing Chen, Víctor Gutiérrez-Basulto, Jeff Z. Pan, and Hanjie Chen. Miceval: Unveiling multimodal chain of thought's quality via image description and reasoning steps. *CoRR*, abs/2410.14668, 2024.
- [28] Qiguang Chen, Libo Qin, Jin Zhang, Zhi Chen, Xiao Xu, and Wanxiang Che. M³cot: A novel benchmark for multi-domain multi-step multi-modal chain-of-thought. In *ACL*, pages 8199–8221. Association for Computational Linguistics, 2024.
- [29] Zhuosheng Zhang, Aston Zhang, Mu Li, Hai Zhao, George Karypis, and Alex Smola. Multimodal chain-of-thought reasoning in language models. arXiv preprint arXiv:2302.00923, 2023.
- [30] Chengzu Li, Wenshan Wu, Huanyu Zhang, Yan Xia, Shaoguang Mao, Li Dong, Ivan Vulić, and Furu Wei. Imagine while reasoning in space: Multimodal visualization-of-thought. arXiv preprint arXiv:2501.07542, 2025.
- [31] Hao Fei, Shengqiong Wu, Wei Ji, Hanwang Zhang, Meishan Zhang, Mong-Li Lee, and Wynne Hsu. Video-of-thought: Step-by-step video reasoning from perception to cognition. In *Forty-first International Conference on Machine Learning*, 2024.
- [32] Guangyao Li, Henghui Du, and Di Hu. Avqa-cot: When cot meets question answering in audio-visual scenarios. In *CVPR Workshops*, 2024.
- [33] Eslam Abdelrahman, Mohamed Ayman, Mahmoud Ahmed, Habib Slim, and Mohamed Elhoseiny. Cot3dref: Chain-of-thoughts data-efficient 3d visual grounding. *arXiv preprint arXiv:2310.06214*, 2023.
- [34] Ziyu Guo, Renrui Zhang, Chengzhuo Tong, Zhizheng Zhao, Peng Gao, Hongsheng Li, and Pheng-Ann Heng. Can we generate images with cot? let's verify and reinforce image generation step by step. *arXiv preprint arXiv:2501.13926*, 2025.

- [35] Xiaoyu Tian, Junru Gu, Bailin Li, Yicheng Liu, Chenxu Hu, Yang Wang, Kun Zhan, Peng Jia, Xianpeng Lang, and Hang Zhao. Drivevlm: The convergence of autonomous driving and large vision-language models. *CoRR*, abs/2402.12289, 2024.
- [36] Licheng Wen, Daocheng Fu, Xin Li, Xinyu Cai, Tao Ma, Pinlong Cai, Min Dou, Botian Shi, Liang He, and Yu Qiao. Dilu: A knowledge-driven approach to autonomous driving with large language models. In *ICLR*, 2024.
- [37] Shihao Wang, Zhiding Yu, Xiaohui Jiang, Shiyi Lan, Min Shi, Nadine Chang, Jan Kautz, Ying Li, and José M. Álvarez. Omnidrive: A holistic llm-agent framework for autonomous driving with 3d perception, reasoning and planning. *CoRR*, abs/2405.01533, 2024.
- [38] Yifan Bai, Dongming Wu, Yingfei Liu, Fan Jia, Weixin Mao, Ziheng Zhang, Yucheng Zhao, Jianbing Shen, Xing Wei, Tiancai Wang, and Xiangyu Zhang. Is a 3d-tokenized LLM the key to reliable autonomous driving? *CoRR*, abs/2405.18361, 2024.
- [39] Yao Mu, Qinglong Zhang, Mengkang Hu, Wenhai Wang, Mingyu Ding, Jun Jin, Bin Wang, Jifeng Dai, Yu Qiao, and Ping Luo. Embodiedgpt: Vision-language pre-training via embodied chain of thought. *Advances in Neural Information Processing Systems*, 36:25081–25094, 2023.
- [40] Ming-Yi Lin, Ou-Wen Lee, and Chih-Ying Lu. Embodied AI with large language models: A survey and new HRI framework. In *ICARM*, pages 978–983, 2024.
- [41] Michał Zawalski, William Chen, Karl Pertsch, Oier Mees, Chelsea Finn, and Sergey Levine. Robotic control via embodied chain-of-thought reasoning. arXiv preprint arXiv:2407.08693, 2024.
- [42] Ishika Singh, Valts Blukis, Arsalan Mousavian, Ankit Goyal, Danfei Xu, Jonathan Tremblay, Dieter Fox, Jesse Thomason, and Animesh Garg. Progprompt: Generating situated robot task plans using large language models. In *ICRA*, pages 11523–11530, 2023.
- [43] Xuan Xiao, Jiahang Liu, Zhipeng Wang, Yanmin Zhou, Yong Qi, Qian Cheng, Bin He, and Shuo Jiang. Robot learning in the era of foundation models: A survey. *CoRR*, abs/2311.14379, 2023.
- [44] Krishan Rana, Jesse Haviland, Sourav Garg, Jad Abou-Chakra, Ian D. Reid, and Niko Sünderhauf. Sayplan: Grounding large language models using 3d scene graphs for scalable robot task planning. In *CoRL*, pages 23–72, 2023.
- [45] Brianna Zitkovich, Tianhe Yu, Sichun Xu, Peng Xu, Ted Xiao, Fei Xia, Jialin Wu, Paul Wohlhart, Stefan Welker, Ayzaan Wahid, Quan Vuong, Vincent Vanhoucke, Huong T. Tran, Radu Soricut, Anikait Singh, Jaspiar Singh, Pierre Sermanet, Pannag R. Sanketi, Grecia Salazar, Michael S. Ryoo, Krista Reymann, Kanishka Rao, Karl Pertsch, Igor Mordatch, Henryk Michalewski, Yao Lu, Sergey Levine, Lisa Lee, Tsang-Wei Edward Lee, Isabel Leal, Yuheng Kuang, Dmitry Kalashnikov, Ryan Julian, Nikhil J. Joshi, Alex Irpan, Brian Ichter, Jasmine Hsu, Alexander Herzog, Karol Hausman, Keerthana Gopalakrishnan, Chuyuan Fu, Pete Florence, Chelsea Finn, Kumar Avinava Dubey, Danny Driess, Tianli Ding, Krzysztof Marcin Choromanski, Xi Chen, Yevgen Chebotar, Justice Carbajal, Noah Brown, Anthony Brohan, Montserrat Gonzalez Arenas, and Kehang Han. RT-2: vision-language-action models transfer web knowledge to robotic control. In *CoRL*, pages 2165–2183, 2023.
- [46] Yuting He, Fuxiang Huang, Xinrui Jiang, Yuxiang Nie, Minghao Wang, Jiguang Wang, and Hao Chen. Foundation model for advancing healthcare: Challenges, opportunities, and future directions. *CoRR*, abs/2404.03264, 2024.
- [47] Sagar Goyal, Eti Rastogi, Sree Prasanna Rajagopal, Dong Yuan, Fen Zhao, Jai Chintagunta, Gautam Naik, and Jeff Ward. Healai: A healthcare LLM for effective medical documentation. In *WSDM*, pages 1167–1168, 2024.
- [48] Zhiyao Ren, Yibing Zhan, Baosheng Yu, Liang Ding, and Dacheng Tao. Healthcare copilot: Eliciting the power of general llms for medical consultation. *CoRR*, abs/2402.13408, 2024.

- [49] Yiqiao Jin, Mohit Chandra, Gaurav Verma, Yibo Hu, Munmun De Choudhury, and Srijan Kumar. Better to ask in english: Cross-lingual evaluation of large language models for healthcare queries. In *ACM WWW*, pages 2627–2638, 2024.
- [50] Ziyu Wang, Hao Li, Di Huang, and Amir M. Rahmani. Healthq: Unveiling questioning capabilities of LLM chains in healthcare conversations. *CoRR*, abs/2409.19487, 2024.
- [51] Zhenfang Chen, Qinhong Zhou, Yikang Shen, Yining Hong, Hao Zhang, and Chuang Gan. See, think, confirm: Interactive prompting between vision and language models for knowledge-based visual reasoning. *arXiv* preprint arXiv:2301.05226, 2023.
- [52] Juncheng Yang, Zuchao Li, Shuai Xie, Wei Yu, Shijun Li, and Bo Du. Soft-prompting with graph-of-thought for multi-modal representation learning. *arXiv preprint arXiv:2404.04538*, 2024.
- [53] Junyi Yao, Yijiang Liu, Zhen Dong, Mingfei Guo, Helan Hu, Kurt Keutzer, Li Du, Daquan Zhou, and Shanghang Zhang. Promptcot: Align prompt distribution via adapted chain-of-thought. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7027–7037, 2024.
- [54] Daniel Rose, Vaishnavi Himakunthala, Andy Ouyang, Ryan He, Alex Mei, Yujie Lu, Michael Saxon, Chinmay Sonar, Diba Mirza, and William Yang Wang. Visual chain of thought: bridging logical gaps with multimodal infillings. *arXiv preprint arXiv:2305.02317*, 2023.
- [55] Lei Wang, Yi Hu, Jiabang He, Xing Xu, Ning Liu, Hui Liu, and Heng Tao Shen. T-sciq: Teaching multimodal chain-of-thought reasoning via large language model signals for science question answering. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 19162–19170, 2024.
- [56] Cheng Tan, Jingxuan Wei, Zhangyang Gao, Linzhuang Sun, Siyuan Li, Ruifeng Guo, Bihui Yu, and Stan Z Li. Boosting the power of small multimodal reasoning models to match larger models with self-consistency training. In *European Conference on Computer Vision*, pages 305–322. Springer, 2024.
- [57] Fanglong Yao, Changyuan Tian, Jintao Liu, Zequn Zhang, Qing Liu, Li Jin, Shuchao Li, Xiaoyu Li, and Xian Sun. Thinking like an expert: Multimodal hypergraph-of-thought (hot) reasoning to boost foundation modals. *arXiv preprint arXiv:2308.06207*, 2023.
- [58] Jiajin Tang, Ge Zheng, Jingyi Yu, and Sibei Yang. Cotdet: Affordance knowledge prompting for task driven object detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3068–3078, 2023.
- [59] Ge Zheng, Bin Yang, Jiajin Tang, Hong-Yu Zhou, and Sibei Yang. Ddcot: Duty-distinct chain-of-thought prompting for multimodal reasoning in language models. *Advances in Neural Information Processing Systems*, 36:5168–5191, 2023.
- [60] Lei Li. Cpseg: Finer-grained image semantic segmentation via chain-of-thought language prompting. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 513–522, 2024.
- [61] Pushkal Katara, Zhou Xian, and Katerina Fragkiadaki. Gen2sim: Scaling up robot learning in simulation with generative models. In 2024 IEEE International Conference on Robotics and Automation (ICRA), pages 6672–6679. IEEE, 2024.
- [62] Fanxu Meng, Haotong Yang, Yiding Wang, and Muhan Zhang. Chain of images for intuitively reasoning. *arXiv preprint arXiv:2311.09241*, 2023.
- [63] Lai Wei, Wenkai Wang, Xiaoyu Shen, Yu Xie, Zhihao Fan, Xiaojin Zhang, Zhongyu Wei, and Wei Chen. Mc-cot: A modular collaborative cot framework for zero-shot medical-vqa with LLM and MLLM integration. CoRR, abs/2410.04521, 2024.
- [64] Chancharik Mitra, Brandon Huang, Trevor Darrell, and Roei Herzig. Compositional chain-of-thought prompting for large multimodal models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14420–14431, 2024.

- [65] Shanshan Zhong, Zhongzhan Huang, Shanghua Gao, Wushao Wen, Liang Lin, Marinka Zitnik, and Pan Zhou. Let's think outside the box: Exploring leap-of-thought in large language models with creative humor generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13246–13257, 2024.
- [66] Liqi He, Zuchao Li, Xiantao Cai, and Ping Wang. Multi-modal latent space learning for chain-of-thought reasoning in language models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 18180–18187, 2024.
- [67] Daoan Zhang, Junming Yang, Hanjia Lyu, Zijian Jin, Yuan Yao, Mingkai Chen, and Jiebo Luo. Cocot: Contrastive chain-of-thought prompting for large multimodal models with multiple image inputs. *arXiv* preprint arXiv:2401.02582, 2024.
- [68] Debjyoti Mondal, Suraj Modi, Subhadarshi Panda, Rituraj Singh, and Godawari Sudhakar Rao. Kam-cot: Knowledge augmented multimodal chain-of-thoughts reasoning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 18798–18806, 2024.
- [69] Xuewen Luo, Fan Ding, Yinsheng Song, Xiaofeng Zhang, and Junnyong Loo. Pkrd-cot: A unified chain-of-thought prompting for multi-modal large language models in autonomous driving. *arXiv preprint arXiv:2412.02025*, 2024.
- [70] Zuyan Liu, Yuhao Dong, Yongming Rao, Jie Zhou, and Jiwen Lu. Chain-of-spot: Interactive reasoning improves large vision-language models. *arXiv preprint arXiv:2403.12966*, 2024.
- [71] Zhenyu Pan, Haozheng Luo, Manling Li, and Han Liu. Chain-of-action: Faithful and multi-modal question answering through large language models. arXiv preprint arXiv:2403.17359, 2024.
- [72] Yixuan Wu, Yizhou Wang, Shixiang Tang, Wenhao Wu, Tong He, Wanli Ouyang, Philip Torr, and Jian Wu. Dettoolchain: A new prompting paradigm to unleash detection ability of mllm. In *European Conference on Computer Vision*, pages 164–182. Springer, 2024.
- [73] Changmeng Zheng, Dayong Liang, Wengyu Zhang, Xiao-Yong Wei, Tat-Seng Chua, and Qing Li. A picture is worth a graph: A blueprint debate paradigm for multimodal reasoning. In *Proceedings of the 32nd ACM International Conference on Multimedia*, pages 419–428, 2024.
- [74] Bozhi Luan, Hao Feng, Hong Chen, Yonghui Wang, Wengang Zhou, and Houqiang Li. Textcot: Zoom in for enhanced multimodal text-rich image understanding. arXiv preprint arXiv:2404.09797, 2024.
- [75] M Abdul Khaliq, P Chang, M Ma, Bernhard Pflugfelder, and F Miletić. Ragar, your falsehood radar: Rag-augmented reasoning for political fact-checking using multimodal large language models. *arXiv preprint arXiv:2404.12065*, 2024.
- [76] Timin Gao, Peixian Chen, Mengdan Zhang, Chaoyou Fu, Yunhang Shen, Yan Zhang, Shengchuan Zhang, Xiawu Zheng, Xing Sun, Liujuan Cao, et al. Cantor: Inspiring multimodal chain-of-thought of mllm. In *Proceedings of the 32nd ACM International Conference on Multimedia*, pages 9096–9105, 2024.
- [77] Yushi Hu, Weijia Shi, Xingyu Fu, Dan Roth, Mari Ostendorf, Luke Zettlemoyer, Noah A Smith, and Ranjay Krishna. Visual sketchpad: Sketching as a visual chain of thought for multimodal language models. *arXiv preprint arXiv:2406.09403*, 2024.
- [78] Qiji Zhou, Ruochen Zhou, Zike Hu, Panzhong Lu, Siyang Gao, and Yue Zhang. Image-of-thought prompting for visual reasoning refinement in multimodal large language models. *arXiv preprint arXiv:2405.13872*, 2024.
- [79] Qun Li, Haixin Sun, Fu Xiao, Yiming Wang, Xinping Gao, and Bir Bhanu. Ps-cot-adapter: adapting plan-and-solve chain-of-thought for scienceqa. *Science China Information Sciences*, 68(1):119101, 2025.

- [80] Yingzi Ma, Yulong Cao, Jiachen Sun, Marco Pavone, and Chaowei Xiao. Dolphins: Multimodal language model for driving. In *European Conference on Computer Vision*, pages 403–420. Springer, 2024.
- [81] Yihe Deng, Pan Lu, Fan Yin, Ziniu Hu, Sheng Shen, Quanquan Gu, James Zou, Kai-Wei Chang, and Wei Wang. Enhancing large vision language models with self-training on image comprehension. *arXiv preprint arXiv:2405.19716*, 2024.
- [82] Guangmin Zheng, Jin Wang, Xiaobing Zhou, and Xuejie Zhang. Enhancing semantics in multimodal chain of thought via soft negative sampling. arXiv preprint arXiv:2405.09848, 2024.
- [83] Haohao Luo, Yang Deng, Ying Shen, See-Kiong Ng, and Tat-Seng Chua. Chain-of-exemplar: enhancing distractor generation for multimodal educational question generation. In ACL, 2024.
- [84] Zixi Jia, Jiqiang Liu, Hexiao Li, Qinghua Liu, and Hongbin Gao. Dcot: Dual chain-of-thought prompting for large multimodal models. In *The 16th Asian Conference on Machine Learning (Conference Track)*, 2024.
- [85] Leigang Qu, Shengqiong Wu, Hao Fei, Liqiang Nie, and Tat-Seng Chua. Layoutllm-t2i: Eliciting layout guidance from llm for text-to-image generation. In *Proceedings of the 31st ACM International Conference on Multimedia*, pages 643–654, 2023.
- [86] Hui Zhang, Dexiang Hong, Tingwei Gao, Yitong Wang, Jie Shao, Xinglong Wu, Zuxuan Wu, and Yu-Gang Jiang. Creatilayout: Siamese multimodal diffusion transformer for creative layout-to-image generation. *arXiv* preprint arXiv:2412.03859, 2024.
- [87] Minheng Ni, Yutao Fan, Lei Zhang, and Wangmeng Zuo. Visual-o1: Understanding ambiguous instructions via multi-modal multi-turn chain-of-thoughts reasoning. *arXiv* preprint *arXiv*:2410.03321, 2024.
- [88] Linger Deng, Yuliang Liu, Bohan Li, Dongliang Luo, Liang Wu, Chengquan Zhang, Pengyuan Lyu, Ziyang Zhang, Gang Zhang, Errui Ding, et al. R-cot: Reverse chain-of-thought problem generation for geometric reasoning in large multimodal models. *arXiv* preprint arXiv:2410.17885, 2024.
- [89] Haojie Zheng, Tianyang Xu, Hanchi Sun, Shu Pu, Ruoxi Chen, and Lichao Sun. Thinking before looking: Improving multimodal llm reasoning via mitigating visual hallucination. *arXiv* preprint arXiv:2411.12591, 2024.
- [90] Chi Xie, Shuang Liang, Jie Li, Zhao Zhang, Feng Zhu, Rui Zhao, and Yichen Wei. Relationlmm: Large multimodal model as open and versatile visual relationship generalist. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2025.
- [91] Yuhao Dong, Zuyan Liu, Hai-Long Sun, Jingkang Yang, Winston Hu, Yongming Rao, and Ziwei Liu. Insight-v: Exploring long-chain visual reasoning with multimodal large language models. *arXiv preprint arXiv:2411.14432*, 2024.
- [92] Mahtab Bigverdi, Zelun Luo, Cheng-Yu Hsieh, Ethan Shen, Dongping Chen, Linda G Shapiro, and Ranjay Krishna. Perception tokens enhance visual reasoning in multimodal language models. *arXiv preprint arXiv:2412.03548*, 2024.
- [93] Guanting Dong, Chenghao Zhang, Mengjie Deng, Yutao Zhu, Zhicheng Dou, and Ji-Rong Wen. Progressive multimodal reasoning via active retrieval. *arXiv preprint arXiv:2412.14835*, 2024.
- [94] Huanjin Yao, Jiaxing Huang, Wenhao Wu, Jingyi Zhang, Yibo Wang, Shunyu Liu, Yingjie Wang, Yuxin Song, Haocheng Feng, Li Shen, et al. Mulberry: Empowering mllm with o1-like reasoning and reflection via collective monte carlo tree search. *arXiv preprint arXiv:2412.18319*, 2024.

- [95] Yifan Du, Zikang Liu, Yifan Li, Wayne Xin Zhao, Yuqi Huo, Bingning Wang, Weipeng Chen, Zheng Liu, Zhongyuan Wang, and Ji-Rong Wen. Virgo: A preliminary exploration on reproducing o1-like mllm. *arXiv* preprint arXiv:2501.01904, 2025.
- [96] Wanpeng Hu, Haodi Liu, Lin Chen, Feng Zhou, Changming Xiao, Qi Yang, and Changshui Zhang. Socratic questioning: Learn to self-guide multimodal reasoning in the wild. *arXiv* preprint arXiv:2501.02964, 2025.
- [97] Omkar Thawakar, Dinura Dissanayake, Ketan More, Ritesh Thawkar, Ahmed Heakl, Noor Ahsan, Yuhao Li, Mohammed Zumri, Jean Lahoud, Rao Muhammad Anwer, et al. Llamavol: Rethinking step-by-step visual reasoning in llms. *arXiv preprint arXiv:2501.06186*, 2025.
- [98] Ruilin Luo, Zhuofan Zheng, Yifan Wang, Yiyao Yu, Xinzhe Ni, Zicheng Lin, Jin Zeng, and Yujiu Yang. Ursa: Understanding and verifying chain-of-thought reasoning in multimodal mathematics. *arXiv preprint arXiv:2501.04686*, 2025.
- [99] EvolvingLMMs Lab. Multimodal open r1, 2025. URL https://github.com/ EvolvingLMMs-Lab/open-r1-multimodal. Accessed: 2025-02-28.
- [100] Jinyang Wu, Mingkuan Feng, Shuai Zhang, Ruihan Jin, Feihu Che, Zengqi Wen, and Jianhua Tao. Boosting multimodal reasoning with mcts-automated structured thinking. *arXiv* preprint *arXiv*:2502.02339, 2025.
- [101] R1-onevision: open-source multimodal large language model with reasoning ability, 2025. URL https://yangyi-vai.notion.site/r1-onevision# 198b1e4047f780c78306fb451be7160d.
- [102] Simon A. Aytes, Jinheon Baek, and Sung Ju Hwang. Sketch-of-thought: Efficient llm reasoning with adaptive cognitive-inspired sketching, 2025.
- [103] Jiapeng Li, Ping Wei, Wenjuan Han, and Lifeng Fan. Intentqa: Context-aware video intent reasoning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 11963–11974, 2023.
- [104] Xiaohan Wang, Yuhui Zhang, Orr Zohar, and Serena Yeung-Levy. Videoagent: Long-form video understanding with large language model as agent. In *European Conference on Com*puter Vision, pages 58–76. Springer, 2024.
- [105] Yiwei Sun, Zhihang Liu, Chuanbin Liu, Bowei Pu, Zhihan Zhang, and Hongtao Xie. Hallucination mitigation prompts long-term video understanding. arXiv preprint arXiv:2406.11333, 2024.
- [106] Yunlong Tang, Gen Zhan, Li Yang, Yiting Liao, and Chenliang Xu. Cardiff: Video salient object ranking chain of thought reasoning for saliency prediction with diffusion. *arXiv* preprint *arXiv*:2408.12009, 2024.
- [107] Rongpei Hong, Jian Lang, Jin Xu, Zhangtao Cheng, Ting Zhong, and Fan Zhou. Following clues, approaching the truth: Explainable micro-video rumor detection via chain-of-thought reasoning. In *THE WEB CONFERENCE* 2025, 2025.
- [108] Qi Zhao, Shijie Wang, Ce Zhang, Changcheng Fu, Minh Quan Do, Nakul Agarwal, Kwonjoon Lee, and Chen Sun. Antgpt: Can large language models help long-term action anticipation from videos? *arXiv preprint arXiv:2307.16368*, 2023.
- [109] Houlun Chen, Xin Wang, Hong Chen, Zihan Song, Jia Jia, and Wenwu Zhu. Grounding-prompter: Prompting llm with multimodal information for temporal sentence grounding in long videos. *arXiv preprint arXiv:2312.17117*, 2023.
- [110] Vaishnavi Himakunthala, Andy Ouyang, Daniel Rose, Ryan He, Alex Mei, Yujie Lu, Chinmay Sonar, Michael Saxon, and William Yang Wang. Let's think frame by frame with vip: A video infilling and prediction dataset for evaluating video chain-of-thought. *arXiv preprint arXiv:2305.13903*, 2023.

- [111] Zhifei Xie, Daniel Tang, Dingwei Tan, Jacques Klein, Tegawend F Bissyand, and Saad Ezzini. Dreamfactory: Pioneering multi-scene long video generation with a multi-agent framework. *arXiv preprint arXiv:2408.11788*, 2024.
- [112] Jian Hu, Zixu Cheng, Chenyang Si, Wei Li, and Shaogang Gong. Cos: Chain-of-shot prompting for long video understanding. *arXiv preprint arXiv:2502.06428*, 2025.
- [113] Leonardo Plini, Luca Scofano, Edoardo De Matteis, Guido Maria D'Amely di Melendugno, Alessandro Flaborea, Andrea Sanchietti, Giovanni Maria Farinella, Fabio Galasso, and Antonino Furnari. Ti-prego: Chain of thought and in-context learning for online mistake detection in procedural egocentric videos. *arXiv preprint arXiv:2411.02570*, 2024.
- [114] Dong Zhang, Xin Zhang, Jun Zhan, Shimin Li, Yaqian Zhou, and Xipeng Qiu. Speechgpt-gen: Scaling chain-of-information speech generation. *arXiv preprint arXiv:2401.13527*, 2024.
- [115] Yexing Du, Ziyang Ma, Yifan Yang, Keqi Deng, Xie Chen, Bo Yang, Yang Xiang, Ming Liu, and Bing Qin. Cot-st: Enhancing llm-based speech translation with multimodal chain-of-thought. *arXiv preprint arXiv:2409.19510*, 2024.
- [116] Jingran Xie, Shun Lei, Yue Yu, Yang Xiang, Hui Wang, Xixin Wu, and Zhiyong Wu. Lever-aging chain of thought towards empathetic spoken dialogue without corresponding question-answering data. arXiv preprint arXiv:2501.10937, 2025.
- [117] Peiwen Sun, Sitong Cheng, Xiangtai Li, Zhen Ye, Huadai Liu, Honggang Zhang, Wei Xue, and Yike Guo. Both ears wide open: Towards language-driven spatial audio generation. *arXiv* preprint arXiv:2410.10676, 2024.
- [118] Ziyang Ma, Zhuo Chen, Yuping Wang, Eng Siong Chng, and Xie Chen. Audio-cot: Exploring chain-of-thought reasoning in large audio language model. *arXiv preprint arXiv:2501.07246*, 2025.
- [119] Zhifei Xie, Mingbao Lin, Zihang Liu, Pengcheng Wu, Shuicheng Yan, and Chunyan Miao. Audio-reasoner: Improving reasoning capability in large audio language models. *arXiv* preprint arXiv:2503.02318, 2025.
- [120] Chuwei Luo, Yufan Shen, Zhaoqing Zhu, Qi Zheng, Zhi Yu, and Cong Yao. Layoutllm: Layout instruction tuning with large language models for document understanding. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 15630–15640, 2024.
- [121] Liangyu Zha, Junlin Zhou, Liyao Li, Rui Wang, Qingyi Huang, Saisai Yang, Jing Yuan, Changbao Su, Xiang Li, Aofeng Su, et al. Tablegpt: Towards unifying tables, nature language and commands into one gpt. *arXiv preprint arXiv:2307.08674*, 2023.
- [122] Zilong Wang, Hao Zhang, Chun-Liang Li, Julian Martin Eisenschlos, Vincent Perot, Zifeng Wang, Lesly Miculicich, Yasuhisa Fujii, Jingbo Shang, Chen-Yu Lee, et al. Chain-of-table: Evolving tables in the reasoning chain for table understanding. *arXiv preprint* arXiv:2401.04398, 2024.
- [123] Xingyu Fu, Minqian Liu, Zhengyuan Yang, John Corring, Yijuan Lu, Jianwei Yang, Dan Roth, Dinei Florencio, and Cha Zhang. Refocus: Visual editing as a chain of thought for structured image understanding. *arXiv* preprint arXiv:2501.05452, 2025.
- [124] Zeqing Yuan, Haoxuan Lan, Qiang Zou, and Junbo Zhao. 3d-premise: Can large language models generate 3d shapes with sharp features and parametric control? *arXiv preprint arXiv:2401.06437*, 2024.
- [125] Yutaro Yamada, Khyathi Chandu, Yuchen Lin, Jack Hessel, Ilker Yildirim, and Yejin Choi. L3go: Language agents with chain-of-3d-thoughts for generating unconventional objects. arXiv preprint arXiv:2402.09052, 2024.

- [126] Yanjun Chen, Yirong Sun, Xinghao Chen, Jian Wang, Xiaoyu Shen, Wenjie Li, and Wei Zhang. Integrating chain-of-thought for multimodal alignment: A study on 3d vision-language learning, 2025.
- [127] Yaoting Wang, Peiwen Sun, Yuanchao Li, Honggang Zhang, and Di Hu. Can textual semantics mitigate sounding object segmentation preference? In *European Conference on Computer Vision*, pages 340–356. Springer, 2024.
- [128] Han Zhang, Zixiang Meng, Meng Luo, Hong Han, Lizi Liao, Erik Cambria, and Hao Fei. Towards multimodal empathetic response generation: A rich text-speech-vision avatar-based benchmark. *arXiv preprint arXiv:2502.04976*, 2025.
- [129] Yan Li, Xiangyuan Lan, Haifeng Chen, Ke Lu, and Dongmei Jiang. Multimodal pear chain-of-thought reasoning for multimodal sentiment analysis. *ACM Transactions on Multimedia Computing, Communications and Applications*, 2024.
- [130] Jiaxing Zhao, Xihan Wei, and Liefeng Bo. R1-omni: Explainable omni-multimodal emotion recognition with reinforcement learning, 2025.
- [131] Kye Gomez. Multimodal-tot. https://github.com/kyegomez/MultiModal-ToT, 2023.
- [132] Jiaxin Ge, Hongyin Luo, Siyuan Qian, Yulu Gan, Jie Fu, and Shanghang Zhang. Chain of thought prompt tuning in vision language models. *arXiv preprint arXiv:2304.07919*, 2023.
- [133] Xiwen Liang, Min Lin, Weiqi Ruan, Yuecheng Liu, Yuzheng Zhuang, and Xiaodan Liang. Memory-driven multimodal chain of thought for embodied long-horizon task planning. *Openreview*, 2025.
- [134] Yao Yao, Zuchao Li, and Hai Zhao. Beyond chain-of-thought, effective graph-of-thought reasoning in language models. *arXiv preprint arXiv:2305.16582*, 2023.
- [135] You Li, Heyu Huang, Chi Chen, Kaiyu Huang, Chao Huang, Zonghao Guo, Zhiyuan Liu, Jinan Xu, Yuhua Li, Ruixuan Li, et al. Migician: Revealing the magic of free-form multi-image grounding in multimodal large language models. *arXiv preprint arXiv:2501.05767*, 2025.
- [136] Fei Ni, Jianye Hao, Shiguang Wu, Longxin Kou, Jiashun Liu, Yan Zheng, Bin Wang, and Yuzheng Zhuang. Generate subgoal images before act: Unlocking the chain-of-thought reasoning in diffusion model for robot manipulation with multimodal prompts. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13991–14000, 2024.
- [137] Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
- [138] Ruohong Zhang, Bowen Zhang, Yanghao Li, Haotian Zhang, Zhiqing Sun, Zhe Gan, Yinfei Yang, Ruoming Pang, and Yiming Yang. Improve vision language model chain-of-thought reasoning. *arXiv preprint arXiv:2410.16198*, 2024.
- [139] Xiaoqi Li, Mingxu Zhang, Yiran Geng, Haoran Geng, Yuxing Long, Yan Shen, Renrui Zhang, Jiaming Liu, and Hao Dong. Manipllm: Embodied multimodal large language model for object-centric robotic manipulation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18061–18070, 2024.
- [140] Qi Sun, Pengfei Hong, Tej Deep Pala, Vernon Toh, U Tan, Deepanway Ghosal, Soujanya Poria, et al. Emma-x: An embodied multimodal action model with grounded chain of thought and look-ahead spatial reasoning. *arXiv preprint arXiv:2412.11974*, 2024.
- [141] Yuecheng Liu, Dafeng Chi, Shiguang Wu, Zhanguang Zhang, Yaochen Hu, Lingfeng Zhang, Yingxue Zhang, Shuang Wu, Tongtong Cao, Guowei Huang, et al. Spatialcot: Advancing spatial reasoning through coordinate alignment and chain-of-thought for embodied task planning. arXiv preprint arXiv:2501.10074, 2025.

- [142] Zhixuan Shen, Haonan Luo, Kexun Chen, Fengmao Lv, and Tianrui Li. Enhancing multirobot semantic navigation through multimodal chain-of-thought score collaboration. *arXiv* preprint arXiv:2412.18292, 2024.
- [143] Zhuosheng Zhang and Aston Zhang. You only look at screens: Multimodal chain-of-action agents. *arXiv preprint arXiv:2309.11436*, 2023.
- [144] Jiaqi Zhang, Chen Gao, Liyuan Zhang, Yong Li, and Hongzhi Yin. Smartagent: Chain-of-user-thought for embodied personalized agent in cyber world. *arXiv preprint* arXiv:2412.07472, 2024.
- [145] Tianqi Wang, Enze Xie, Ruihang Chu, Zhenguo Li, and Ping Luo. Drivecot: Integrating chain-of-thought reasoning with end-to-end driving. arXiv preprint arXiv:2403.16996, 2024.
- [146] Can Cui, Yunsheng Ma, Xu Cao, Wenqian Ye, and Ziran Wang. Receive, reason, and react: Drive as you say, with large language models in autonomous vehicles. *IEEE Intelligent Transportation Systems Magazine*, 2024.
- [147] Yunsheng Ma, Xu Cao, Wenqian Ye, Can Cui, Kai Mei, and Ziran Wang. Learning autonomous driving tasks via human feedbacks with large language models. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 4985–4995, 2024.
- [148] Rui Zhao, Qirui Yuan, Jinyu Li, Haofeng Hu, Yun Li, Chengyuan Zheng, and Fei Gao. Sce2drivex: A generalized mllm framework for scene-to-drive learning. *arXiv preprint* arXiv:2502.14917, 2025.
- [149] Ming Nie, Renyuan Peng, Chunwei Wang, Xinyue Cai, Jianhua Han, Hang Xu, and Li Zhang. Reason2drive: Towards interpretable and chain-based reasoning for autonomous driving. In *European Conference on Computer Vision*, pages 292–308. Springer, 2024.
- [150] Haicheng Liao, Hanlin Kong, Bonan Wang, Chengyue Wang, Wang Ye, Zhengbing He, Chengzhong Xu, and Zhenning Li. Cot-drive: Efficient motion forecasting for autonomous driving with llms and chain-of-thought prompting, 2025.
- [151] Yi Dai. Interpretable video based stress detection with self-refine chain-of-thought reasoning. *arXiv preprint arXiv:2410.09449*, 2024.
- [152] Nan Xi, Jingjing Meng, and Junsong Yuan. Chain-of-look prompting for verb-centric surgical triplet recognition in endoscopic videos. In *Proceedings of the 31st ACM International Conference on Multimedia*, pages 5007–5016, 2023.
- [153] Jiaxiang Liu, Yuan Wang, Jiawei Du, Joey Tianyi Zhou, and Zuozhu Liu. Medcot: Medical chain of thought via hieratrchical expert. *arXiv preprint arXiv:2412.13736*, 2024.
- [154] Jiazhen Pan, Che Liu, Junde Wu, Fenglin Liu, Jiayuan Zhu, Hongwei Bran Li, Chen Chen, Cheng Ouyang, and Daniel Rueckert. Medvlm-r1: Incentivizing medical reasoning capability of vision-language models (vlms) via reinforcement learning. *arXiv* preprint *arXiv*:2502.19634, 2025.
- [155] Meng Luo, Hao Fei, Bobo Li, Shengqiong Wu, Qian Liu, Soujanya Poria, Erik Cambria, Mong-Li Lee, and Wynne Hsu. Panosent: A panoptic sextuple extraction benchmark for multimodal conversational aspect-based sentiment analysis. In *Proceedings of the 32nd ACM International Conference on Multimedia*, pages 7667–7676, 2024.
- [156] Ling Yang, Zhaochen Yu, Chenlin Meng, Minkai Xu, Stefano Ermon, and CUI Bin. Mastering text-to-image diffusion: Recaptioning, planning, and generating with multimodal llms. In *Forty-first International Conference on Machine Learning*, 2024.
- [157] 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, 2022.

- [158] Dustin Schwenk, Apoorv Khandelwal, Christopher Clark, Kenneth Marino, and Roozbeh Mottaghi. A-okvqa: A benchmark for visual question answering using world knowledge. In European conference on computer vision, pages 146–162. Springer, 2022.
- [159] Yan Wang, Yawen Zeng, Jingsheng Zheng, Xiaofen Xing, Jin Xu, and Xiangmin Xu. Videocot: A video chain-of-thought dataset with active annotation tool. arXiv preprint arXiv:2407.05355, 2024.
- [160] Songhao Han, Wei Huang, Hairong Shi, Le Zhuo, Xiu Su, Shifeng Zhang, Xu Zhou, Xiaojuan Qi, Yue Liao, and Si Liu. Videoespresso: A large-scale chain-of-thought dataset for fine-grained video reasoning via core frame selection. arXiv preprint arXiv:2411.14794, 2024.
- [161] Renrui Zhang, Xinyu Wei, Dongzhi Jiang, Ziyu Guo, Shicheng Li, Yichi Zhang, Chengzhuo Tong, Jiaming Liu, Aojun Zhou, Bin Wei, et al. Mavis: Mathematical visual instruction tuning with an automatic data engine. *arXiv preprint arXiv:2407.08739*, 2024.
- [162] Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, et al. Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi. In *Proceedings of the* IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 9556–9567, 2024.
- [163] Bohao Li, Yuying Ge, Yixiao Ge, Guangzhi Wang, Rui Wang, Ruimao Zhang, and Ying Shan. Seed-bench: Benchmarking multimodal large language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13299–13308, 2024.
- [164] 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. arXiv preprint arXiv:2310.02255, 2023.
- [165] Renrui Zhang, Dongzhi Jiang, Yichi Zhang, Haokun Lin, Ziyu Guo, Pengshuo Qiu, Aojun Zhou, Pan Lu, Kai-Wei Chang, Yu Qiao, et al. Mathverse: Does your multi-modal llm truly see the diagrams in visual math problems? In *European Conference on Computer Vision*, pages 169–186. Springer, 2024.
- [166] Ke Wang, Junting Pan, Weikang Shi, Zimu Lu, Houxing Ren, Aojun Zhou, Mingjie Zhan, and Hongsheng Li. Measuring multimodal mathematical reasoning with math-vision dataset. *Advances in Neural Information Processing Systems*, 37:95095–95169, 2025.
- [167] Yunzhuo Hao, Jiawei Gu, Huichen Will Wang, Linjie Li, Zhengyuan Yang, Lijuan Wang, and Yu Cheng. Can mllms reason in multimodality? emma: An enhanced multimodal reasoning benchmark. *arXiv preprint arXiv:2501.05444*, 2025.
- [168] Yaoting Wang, Peiwen Sun, Dongzhan Zhou, Guangyao Li, Honggang Zhang, and Di Hu. Ref-avs: Refer and segment objects in audio-visual scenes. In *European Conference on Computer Vision*, pages 196–213. Springer, 2024.
- [169] Jihan Yang, Shusheng Yang, Anjali W Gupta, Rilyn Han, Li Fei-Fei, and Saining Xie. Thinking in space: How multimodal large language models see, remember, and recall spaces. *arXiv* preprint arXiv:2412.14171, 2024.
- [170] Henghui Ding, Chang Liu, Shuting He, Xudong Jiang, and Chen Change Loy. Mevis: A large-scale benchmark for video segmentation with motion expressions. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 2694–2703, 2023.
- [171] Tianrui Guan, Fuxiao Liu, Xiyang Wu, Ruiqi Xian, Zongxia Li, Xiaoyu Liu, Xijun Wang, Lichang Chen, Furong Huang, Yaser Yacoob, et al. Hallusionbench: an advanced diagnostic suite for entangled language hallucination and visual illusion in large vision-language models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14375–14385, 2024.

- [172] Sanjoy Chowdhury, Sayan Nag, Subhrajyoti Dasgupta, Yaoting Wang, Mohamed Elhoseiny, Ruohan Gao, and Dinesh Manocha. Avtrustbench: Assessing and enhancing reliability and robustness in audio-visual llms. *arXiv preprint arXiv:2501.02135*, 2025.
- [173] Kim Sung-Bin, Oh Hyun-Bin, JungMok Lee, Arda Senocak, Joon Son Chung, and Tae-Hyun Oh. Avhbench: A cross-modal hallucination benchmark for audio-visual large language models. *arXiv preprint arXiv:2410.18325*, 2024.
- [174] Zihui Cheng, Qiguang Chen, Jin Zhang, Hao Fei, Xiaocheng Feng, Wanxiang Che, Min Li, and Libo Qin. Comt: A novel benchmark for chain of multi-modal thought on large vision-language models. *arXiv preprint arXiv:2412.12932*, 2024.
- [175] Yuanhan Zhang, Kaichen Zhang, Bo Li, Fanyi Pu, Christopher Arif Setiadharma, Jingkang Yang, and Ziwei Liu. Worldqa: Multimodal world knowledge in videos through long-chain reasoning. arXiv preprint arXiv:2405.03272, 2024.
- [176] Xiongtao Zhou, Jie He, Lanyu Chen, Jingyu Li, Haojing Chen, Víctor Gutiérrez-Basulto, Jeff Z Pan, and Hanjie Chen. Miceval: Unveiling multimodal chain of thought's quality via image description and reasoning steps. arXiv preprint arXiv:2410.14668, 2024.
- [177] Chaoqun He, Renjie Luo, Yuzhuo Bai, Shengding Hu, Zhen Leng Thai, Junhao Shen, Jinyi Hu, Xu Han, Yujie Huang, Yuxiang Zhang, et al. Olympiadbench: A challenging benchmark for promoting agi with olympiad-level bilingual multimodal scientific problems. *arXiv* preprint arXiv:2402.14008, 2024.
- [178] Dongzhi Jiang, Renrui Zhang, Ziyu Guo, Yanwei Li, Yu Qi, Xinyan Chen, Liuhui Wang, Jianhan Jin, Claire Guo, Shen Yan, et al. Mme-cot: Benchmarking chain-of-thought in large multimodal models for reasoning quality, robustness, and efficiency. *arXiv preprint arXiv:2502.09621*, 2025.
- [179] Yizhi Li, Ge Zhang, Yinghao Ma, Ruibin Yuan, Kang Zhu, Hangyu Guo, Yiming Liang, Jiaheng Liu, Zekun Wang, Jian Yang, et al. Omnibench: Towards the future of universal omni-language models. *arXiv preprint arXiv:2409.15272*, 2024.
- [180] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- [181] 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 arXiv:2303.08774*, 2023.
- [182] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- [183] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- [184] Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, et al. A survey of large language models. arXiv preprint arXiv:2303.18223, 2023.
- [185] Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, et al. Emergent abilities of large language models. *arXiv preprint arXiv:2206.07682*, 2022.
- [186] Libo Qin, Qiguang Chen, Hao Fei, Zhi Chen, Min Li, and Wanxiang Che. What factors affect multi-modal in-context learning? an in-depth exploration. *arXiv preprint arXiv:2410.20482*, 2024.

- [187] Zheng Chu, Jingchang Chen, Qianglong Chen, Weijiang Yu, Tao He, Haotian Wang, Weihua Peng, Ming Liu, Bing Qin, and Ting Liu. A survey of chain of thought reasoning: Advances, frontiers and future. *arXiv preprint arXiv:2309.15402*, 2023.
- [188] Maciej Besta, Florim Memedi, Zhenyu Zhang, Robert Gerstenberger, Nils Blach, Piotr Nyczyk, Marcin Copik, Grzegorz Kwaśniewski, Jürgen Müller, Lukas Gianinazzi, et al. Topologies of reasoning: Demystifying chains, trees, and graphs of thoughts. *arXiv preprint arXiv:2401.14295*, 2024.
- [189] Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large language models are zero-shot reasoners. *Advances in neural information processing systems*, 35:22199–22213, 2022.
- [190] Wenhu Chen, Xueguang Ma, Xinyi Wang, and William W Cohen. Program of thoughts prompting: Disentangling computation from reasoning for numerical reasoning tasks. *arXiv* preprint arXiv:2211.12588, 2022.
- [191] Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik Narasimhan. Tree of thoughts: Deliberate problem solving with large language models. *Advances in neural information processing systems*, 36:11809–11822, 2023.
- [192] Jieyi Long. Large language model guided tree-of-thought. *arXiv preprint arXiv:2305.08291*, 2023.
- [193] Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models. *arXiv* preprint arXiv:2203.11171, 2022.
- [194] Anthropic. The claude 3 model family: Opus, sonnet, haiku. https://www-cdn.anthropic.com/de8ba9b01c9ab7cbabf5c33b80b7bbc618857627/Model_Card_Claude_3.pdf, 2024. Preprint.
- [195] 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 *International conference on machine learning*, pages 19730–19742. PMLR, 2023.
- [196] Anas Awadalla, Irena Gao, Josh Gardner, Jack Hessel, Yusuf Hanafy, Wanrong Zhu, Kalyani Marathe, Yonatan Bitton, Samir Gadre, Shiori Sagawa, et al. Openflamingo: An open-source framework for training large autoregressive vision-language models. arXiv preprint arXiv:2308.01390, 2023.
- [197] Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. Minigpt-4: Enhancing vision-language understanding with advanced large language models. arXiv preprint arXiv:2304.10592, 2023.
- [198] KunChang Li, Yinan He, Yi Wang, Yizhuo Li, Wenhai Wang, Ping Luo, Yali Wang, Limin Wang, and Yu Qiao. Videochat: Chat-centric video understanding. *arXiv preprint* arXiv:2305.06355, 2023.
- [199] Yunfei Chu, Jin Xu, Xiaohuan Zhou, Qian Yang, Shiliang Zhang, Zhijie Yan, Chang Zhou, and Jingren Zhou. Qwen-audio: Advancing universal audio understanding via unified large-scale audio-language models. *arXiv preprint arXiv:2311.07919*, 2023.
- [200] Yu Shu, Siwei Dong, Guangyao Chen, Wenhao Huang, Ruihua Zhang, Daochen Shi, Qiqi Xiang, and Yemin Shi. Llasm: Large language and speech model. arXiv preprint arXiv:2308.15930, 2023.
- [201] 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 arXiv:2306.14824, 2023.
- [202] Jing Yu Koh, Daniel Fried, and Russ R Salakhutdinov. Generating images with multimodal language models. Advances in Neural Information Processing Systems, 36:21487–21506, 2023.

- [203] Quan Sun, Qiying Yu, Yufeng Cui, Fan Zhang, Xiaosong Zhang, Yueze Wang, Hongcheng Gao, Jingjing Liu, Tiejun Huang, and Xinlong Wang. Emu: Generative pretraining in multimodality. *arXiv preprint arXiv:2307.05222*, 2023.
- [204] Kaizhi Zheng, Xuehai He, and Xin Eric Wang. Minigpt-5: Interleaved vision-and-language generation via generative vokens. *arXiv preprint arXiv:2310.02239*, 2023.
- [205] Dong Zhang, Shimin Li, Xin Zhang, Jun Zhan, Pengyu Wang, Yaqian Zhou, and Xipeng Qiu. Speechgpt: Empowering large language models with intrinsic cross-modal conversational abilities. *arXiv preprint arXiv:2305.11000*, 2023.
- [206] Dong Zhang, Xin Zhang, Jun Zhan, Shimin Li, Yaqian Zhou, and Xipeng Qiu. Speechgpt-gen: Scaling chain-of-information speech generation. arXiv preprint arXiv:2401.13527, 2024.
- [207] Paul K Rubenstein, Chulayuth Asawaroengchai, Duc Dung Nguyen, Ankur Bapna, Zalán Borsos, Félix de Chaumont Quitry, Peter Chen, Dalia El Badawy, Wei Han, Eugene Kharitonov, et al. Audiopalm: A large language model that can speak and listen. *arXiv* preprint arXiv:2306.12925, 2023.
- [208] Wenyi Hong, Ming Ding, Wendi Zheng, Xinghan Liu, and Jie Tang. Cogvideo: Large-scale pretraining for text-to-video generation via transformers. *arXiv preprint arXiv:2205.15868*, 2022.
- [209] Dan Kondratyuk, Lijun Yu, Xiuye Gu, José Lezama, Jonathan Huang, Grant Schindler, Rachel Hornung, Vighnesh Birodkar, Jimmy Yan, Ming-Chang Chiu, et al. Videopoet: A large language model for zero-shot video generation. arXiv preprint arXiv:2312.14125, 2023.
- [210] Yang Jin, Zhicheng Sun, Kun Xu, Liwei Chen, Hao Jiang, Quzhe Huang, Chengru Song, Yuliang Liu, Di Zhang, Yang Song, et al. Video-lavit: Unified video-language pre-training with decoupled visual-motional tokenization. *arXiv* preprint arXiv:2402.03161, 2024.
- [211] Roberto Henschel, Levon Khachatryan, Daniil Hayrapetyan, Hayk Poghosyan, Vahram Tadevosyan, Zhangyang Wang, Shant Navasardyan, and Humphrey Shi. Streamingt2v: Consistent, dynamic, and extendable long video generation from text. *arXiv preprint arXiv:2403.14773*, 2024.
- [212] Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. Gpt-4o system card. arXiv preprint arXiv:2410.21276, 2024.
- [213] Jun Zhan, Junqi Dai, Jiasheng Ye, Yunhua Zhou, Dong Zhang, Zhigeng Liu, Xin Zhang, Ruibin Yuan, Ge Zhang, Linyang Li, et al. Anygpt: Unified multimodal llm with discrete sequence modeling. *arXiv preprint arXiv:2402.12226*, 2024.
- [214] Zhifei Xie and Changqiao Wu. Mini-omni: Language models can hear, talk while thinking in streaming. *arXiv preprint arXiv:2408.16725*, 2024.
- [215] Zhifei Xie and Changqiao Wu. Mini-omni2: Towards open-source gpt-40 with vision, speech and duplex capabilities. *arXiv* preprint arXiv:2410.11190, 2024.
- [216] Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec Helyar, Aleksander Madry, Alex Beutel, Alex Carney, et al. Openai o1 system card. arXiv preprint arXiv:2412.16720, 2024.
- [217] Wenshan Wu, Shaoguang Mao, Yadong Zhang, Yan Xia, Li Dong, Lei Cui, and Furu Wei. Mind's eye of llms: Visualization-of-thought elicits spatial reasoning in large language models. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024.
- [218] Jaewook Lee, Yeajin Jang, Hongjin Kim, Woojin Lee, and Harksoo Kim. Analyzing key factors influencing emotion prediction performance of vllms in conversational contexts. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 5801–5816, 2024.

- [219] Yuxuan Lei, Dingkang Yang, Zhaoyu Chen, Jiawei Chen, Peng Zhai, and Lihua Zhang. Large vision-language models as emotion recognizers in context awareness. *arXiv preprint arXiv:2407.11300*, 2024.
- [220] Yue Dai, Soyeon Caren Han, and Wei Liu. Multimodal graph constrastive learning and prompt for chartqa. *arXiv preprint arXiv:2501.04303*, 2025.
- [221] Rombach Robin, Blattmann Andreas, Lorenz Dominik, Esser Patrick, and Ommer Björn. High-resolution image synthesis with latent diffusion models, 2021.
- [222] Charlie Snell, Jaehoon Lee, Kelvin Xu, and Aviral Kumar. Scaling Ilm test-time compute optimally can be more effective than scaling model parameters. arXiv preprint arXiv:2408.03314, 2024.
- [223] Edward Yeo, Yuxuan Tong, Morry Niu, Graham Neubig, and Xiang Yue. Demystifying long chain-of-thought reasoning in llms, 2025.
- [224] Yifan Wu, Pengchuan Zhang, Wenhan Xiong, Barlas Oguz, James C Gee, and Yixin Nie. The role of chain-of-thought in complex vision-language reasoning task. *arXiv preprint arXiv:2311.09193*, 2023.
- [225] Guowei Xu, Peng Jin, Li Hao, Yibing Song, Lichao Sun, and Li Yuan. Llava-o1: Let vision language models reason step-by-step. *arXiv* preprint arXiv:2411.10440, 2024.
- [226] Jonathan St BT Evans. In two minds: dual-process accounts of reasoning. *Trends in cognitive sciences*, 7(10):454–459, 2003.
- [227] Zhong-Zhi Li, Duzhen Zhang, Ming-Liang Zhang, Jiaxin Zhang, Zengyan Liu, Yuxuan Yao, Haotian Xu, Junhao Zheng, Pei-Jie Wang, Xiuyi Chen, et al. From system 1 to system 2: A survey of reasoning large language models. *arXiv preprint arXiv:2502.17419*, 2025.
- [228] Jinhao Jiang, Zhipeng Chen, Yingqian Min, Jie Chen, Xiaoxue Cheng, Jiapeng Wang, Yiru Tang, Haoxiang Sun, Jia Deng, Wayne Xin Zhao, et al. Technical report: Enhancing llm reasoning with reward-guided tree search. *arXiv preprint arXiv:2411.11694*, 2024.
- [229] Yingqian Min, Zhipeng Chen, Jinhao Jiang, Jie Chen, Jia Deng, Yiwen Hu, Yiru Tang, Jiapeng Wang, Xiaoxue Cheng, Huatong Song, et al. Imitate, explore, and self-improve: A reproduction report on slow-thinking reasoning systems. arXiv preprint arXiv:2412.09413, 2024.
- [230] Zeyu Gan, Yun Liao, and Yong Liu. Rethinking external slow-thinking: From snowball errors to probability of correct reasoning. *arXiv* preprint arXiv:2501.15602, 2025.
- [231] Qiguang Chen, Libo Qin, Jinhao Liu, Dengyun Peng, Jiannan Guan, Peng Wang, Mengkang Hu, Yuhang Zhou, Te Gao, and Wangxiang Che. Towards reasoning era: A survey of long chain-of-thought for reasoning large language models. *arXiv preprint arXiv:2503.09567*, 2025.
- [232] Qwen Team. Qwq: Reflect deeply on the boundaries of the unknown, November 2024. URL https://qwenlm.github.io/blog/qwq-32b-preview/.
- [233] Yu Zhao, Huifeng Yin, Bo Zeng, Hao Wang, Tianqi Shi, Chenyang Lyu, Longyue Wang, Weihua Luo, and Kaifu Zhang. Marco-o1: Towards open reasoning models for open-ended solutions. *arXiv preprint arXiv:2411.14405*, 2024.
- [234] Haotian Xu, Xing Wu, Weinong Wang, Zhongzhi Li, Da Zheng, Boyuan Chen, Yi Hu, Shijia Kang, Jiaming Ji, Yingying Zhang, et al. Redstar: Does scaling long-cot data unlock better slow-reasoning systems? *arXiv* preprint arXiv:2501.11284, 2025.
- [235] Hugging Face. open-r1. https://github.com/huggingface/open-r1, 2025. GitHub Repository.
- [236] Jiayi Pan. Tinyzero. https://github.com/Jiayi-Pan/TinyZero, 2025. GitHub Repository.

- [237] Liang Chen, Lei Li, Haozhe Zhao, Yifan Song, and Vinci. R1-v: Reinforcing super generalization ability in vision-language models with less than \$3. https://github.com/Deep-Agent/R1-v, 2025. Accessed: 2025-02-02.
- [238] Haozhan Shen, Zilun Zhang, Qianqian Zhang, Ruochen Xu, and Tiancheng Zhao. Vlmr1: A stable and generalizable r1-style large vision-language model. https://github.com/ om-ai-lab/VLM-R1, 2025. Accessed: 2025-02-15.
- [239] Peng Yingzhe, Zhang Gongrui, Zhang Miaosen, You Zhiyuan, Liu Jie, Zhu Qipeng, Yang Kai, Xu Xingzhong, Geng Xin, and Yang Xu. Lmm-r1: Empowering 3b lmms with strong reasoning abilities through two-stage rule-based rl, 2025.
- [240] Zheng Yaowei, Lu Junting, Wang Shenzhi, Feng Zhangchi, Kuang Dongdong, and Xiong Yuwen. Easyr1: An efficient, scalable, multi-modality rl training framework. https://github.com/hiyouga/EasyR1, 2025.
- [241] Wei Liu, Junlong Li, Xiwen Zhang, Fan Zhou, Yu Cheng, and Junxian He. Diving into self-evolving training for multimodal reasoning. *arXiv* preprint arXiv:2412.17451, 2024.
- [242] Wang Weiyun, Gao Zhangwei, Chen Lianjie, Chen Zhe, Zhu Jinguo, Zhao Xiangyu, Liu Yangzhou, Cao Yue, Ye Shenglong, Zhu Xizhou, Lu Lewei, Duan Haodong, Qiao Yu, Dai Jifeng, and Wang Wenhai. Visualprm: An effective process reward model for multimodal reasoning. 2025.
- [243] Wang Xiaodong and Peng Peixi. Open-r1-video. https://github.com/Wang-Xiaodong1899/ Open-R1-Video, 2025.
- [244] Deng Huilin, Zou Ding, Ma Rui, Luo Hongchen, Cao Yang, and Kang Yu. Boosting the generalization and reasoning of vision language models with curriculum reinforcement learning, 2025. URL https://arxiv.org/abs/2503.07065.
- [245] Liu Yuqi, Peng Bohao, Zhong Zhisheng, Yue Zihao, Lu Fanbin, Yu Bei, and Jia Jiaya. Seg-zero: Reasoning-chain guided segmentation via cognitive reinforcement, 2025. URL https://arxiv.org/abs/2503.06520.
- [246] Fanqing Meng, Lingxiao Du, Zongkai Liu, Zhixiang Zhou, Quanfeng Lu, Daocheng Fu, Botian Shi, Wenhai Wang, Junjun He, Kaipeng Zhang, Ping Luo, Yu Qiao, Qiaosheng Zhang, and Wenqi Shao. Mm-eureka: Exploring visual aha moment with rule-based large-scale reinforcement learning, 2025. URL https://github.com/ModalMinds/MM-EUREKA.
- [247] Hengguang Zhou, Xirui Li, Ruochen Wang, Minhao Cheng, Tianyi Zhou, and Cho-Jui Hsieh. R1-zero's "aha moment" in visual reasoning on a 2b non-sft model, 2025. URL https://arxiv.org/abs/2503.05132.
- [248] Zhangquan Chen, Xufang Luo, and Dongsheng Li. Visrl: Intention-driven visual perception via reinforced reasoning. *arXiv preprint arXiv:2503.07523*, 2025.
- [249] Jingyi Zhang, Jiaxing Huang, Huanjin Yao, Shunyu Liu, Xikun Zhang, Shijian Lu, and Dacheng Tao. R1-vl: Learning to reason with multimodal large language models via stepwise group relative policy optimization. *arXiv preprint arXiv:2503.12937*, 2025.
- [250] Monica AI. Manus, 2025. URL https://manus.im/.
- [251] Liang Xinbin, Xiang Jinyu, Yu Zhaoyang, Zhang Jiayi, and Hong Sirui. Openmanus: An open-source framework for building general ai agents. https://github.com/mannaandpoem/OpenManus, 2025.
- [252] Hanjia Lyu, Ryan Rossi, Xiang Chen, Md Mehrab Tanjim, Stefano Petrangeli, Somdeb Sarkhel, and Jiebo Luo. X-reflect: Cross-reflection prompting for multimodal recommendation. *arXiv preprint arXiv:2408.15172*, 2024.
- [253] Yongsheng Yu and Jiebo Luo. Chain-of-thought prompting for demographic inference with large multimodal models. In 2024 IEEE International Conference on Multimedia and Expo (ICME), pages 1–7. IEEE, 2024.

- [254] Fang Rongyao, Duan Chengqi, Wang Kun, Huang Linjiang, Li Hao, Yan Shilin, Tian Hao, Zeng Xingyu, Zhao Rui, Dai Jifeng, Liu Xihui, and Li Hongsheng. Got: Unleashing reasoning capability of multimodal large language model for visual generation and editing, 2025. URL https://arxiv.org/abs/2503.10639.
- [255] Jarvis Guo, Tuney Zheng, Yuelin Bai, Bo Li, Yubo Wang, King Zhu, Yizhi Li, Graham Neubig, Wenhu Chen, and Xiang Yue. Mammoth-vl: Eliciting multimodal reasoning with instruction tuning at scale. *arXiv preprint arXiv:2412.05237*, 2024.
- [256] Linzhuang Sun, Hao Liang, Jingxuan Wei, Bihui Yu, Tianpeng Li, Fan Yang, Zenan Zhou, and Wentao Zhang. Mm-verify: Enhancing multimodal reasoning with chain-of-thought verification. *arXiv preprint arXiv:2502.13383*, 2025.
- [257] Tianbao Xie, Danyang Zhang, Jixuan Chen, Xiaochuan Li, Siheng Zhao, Ruisheng Cao, Jing Hua Toh, Zhoujun Cheng, Dongchan Shin, Fangyu Lei, et al. Osworld: Benchmarking multimodal agents for open-ended tasks in real computer environments. *Advances in Neural Information Processing Systems*, 37:52040–52094, 2024.
- [258] Samuel Schmidgall, Rojin Ziaei, Carl Harris, Eduardo Reis, Jeffrey Jopling, and Michael Moor. Agentclinic: a multimodal agent benchmark to evaluate ai in simulated clinical environments. *arXiv preprint arXiv:2405.07960*, 2024.
- [259] Kaixiong Gong, Kaituo Feng, Bohao Li, Yibing Wang, Mofan Cheng, Shijia Yang, Jiaming Han, Benyou Wang, Yutong Bai, Zhuoran Yang, et al. Av-odyssey bench: Can your multimodal llms really understand audio-visual information? *arXiv preprint arXiv:2412.02611*, 2024.
- [260] S Sakshi, Utkarsh Tyagi, Sonal Kumar, Ashish Seth, Ramaneswaran Selvakumar, Oriol Nieto, Ramani Duraiswami, Sreyan Ghosh, and Dinesh Manocha. Mmau: A massive multi-task audio understanding and reasoning benchmark. In *The Thirteenth International Conference on Learning Representations*, 2024.
- [261] Xiangru Tang, Daniel Shao, Jiwoong Sohn, Jiapeng Chen, Jiayi Zhang, Jinyu Xiang, Fang Wu, Yilun Zhao, Chenglin Wu, Wenqi Shi, et al. Medagentsbench: Benchmarking thinking models and agent frameworks for complex medical reasoning. *arXiv* preprint arXiv:2503.07459, 2025.
- [262] OpenAI. Openai gpt-4.5 system card, 2025. URL https://cdn.openai.com/gpt-4-5-system-card-2272025.pdf.
- [263] OpenAI. Openai gpt-4v system card, 2024. URL https://cdn.openai.com/papers/GPTV_System_Card.pdf.
- [264] Google. Introducing gemini 2.0: our new ai model for the agentic era, 2024. URL https://blog.google/technology/google-deepmind/google-gemini-ai-update-december-2024/#ceo-message.
- [265] Gemini Team, Petko Georgiev, Ving Ian Lei, Ryan Burnell, Libin Bai, Anmol Gulati, Garrett Tanzer, Damien Vincent, Zhufeng Pan, Shibo Wang, et al. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv preprint arXiv:2403.05530*, 2024.
- [266] xAI. Grok-3 beta release, 2025. URL https://x.ai/news/grok-3.
- [267] xAI. Grok-2 beta release, 2024. URL https://x.ai/news/grok-2.
- [268] Kimi Team, Angang Du, Bofei Gao, Bowei Xing, Changjiu Jiang, Cheng Chen, Cheng Li, Chenjun Xiao, Chenzhuang Du, Chonghua Liao, et al. Kimi k1. 5: Scaling reinforcement learning with llms. *arXiv preprint arXiv:2501.12599*, 2025.
- [269] Qwen Team. Qvq: To see the world with wisdom, December 2024. URL https://qwenlm.github.io/blog/qvq-72b-preview/.
- [270] Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, et al. Qwen2. 5-vl technical report. arXiv preprint arXiv:2502.13923, 2025.

- [271] Zhe Chen, Weiyun Wang, Yue Cao, Yangzhou Liu, Zhangwei Gao, Erfei Cui, Jinguo Zhu, Shenglong Ye, Hao Tian, Zhaoyang Liu, et al. Expanding performance boundaries of open-source multimodal models with model, data, and test-time scaling. *arXiv preprint arXiv:2412.05271*, 2024.
- [272] Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo Chen, Sen Xing, Muyan Zhong, Qinglong Zhang, Xizhou Zhu, Lewei Lu, et al. Internvl: Scaling up vision foundation models and aligning for generic visual-linguistic tasks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 24185–24198, 2024.
- [273] Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- [274] Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Peiyuan Zhang, Yanwei Li, Ziwei Liu, et al. Llava-onevision: Easy visual task transfer. *arXiv* preprint arXiv:2408.03326, 2024.
- [275] Feng Li, Renrui Zhang, Hao Zhang, Yuanhan Zhang, Bo Li, Wei Li, Zejun Ma, and Chunyuan Li. Llava-next-interleave: Tackling multi-image, video, and 3d in large multimodal models. arXiv preprint arXiv:2407.07895, 2024.
- [276] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 26296–26306, 2024.
- [277] Saeed Amizadeh, Hamid Palangi, Alex Polozov, Yichen Huang, and Kazuhito Koishida. Neuro-symbolic visual reasoning: Disentangling. In *International Conference on Machine Learning*, pages 279–290, 2020.
- [278] Jundong Xu, Hao Fei, Liangming Pan, Qian Liu, Mong-Li Lee, and Wynne Hsu. Faithful logical reasoning via symbolic chain-of-thought. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 13326–13365, 2024.
- [279] Lauren Nicole DeLong, Ramon Fernández Mir, and Jacques D Fleuriot. Neurosymbolic ai for reasoning over knowledge graphs: A survey. *IEEE Transactions on Neural Networks and Learning Systems*, 2024.
- [280] Yixin Ye, Zhen Huang, Yang Xiao, Ethan Chern, Shijie Xia, and Pengfei Liu. Limo: Less is more for reasoning. *arXiv preprint arXiv:2502.03387*, 2025.