MCTBENCH: MULTIMODAL COGNITION TOWARDS TEXT-RICH VISUAL SCENES BENCHMARK

Anonymous authors

Paper under double-blind review

Abstract

The comprehension of text-rich visual scenes has become a focal point for evaluating Multi-modal Large Language Models (MLLMs) due to their widespread applications. Current benchmarks tailored to the scenario emphasize perceptual capabilities, while overlooking the assessment of cognitive abilities. To address this limitation, we introduce a Multimodal benchmark towards Text-rich visual scenes, to evaluate the Cognitive capabilities of MLLMs through visual reasoning and content-creation tasks (MCTBench). To mitigate potential evaluation bias from the varying distributions of datasets, MCTBench incorporates several perception tasks (e.g., scene text recognition) to ensure a consistent comparison of both the cognitive and perceptual capabilities of MLLMs. To improve the efficiency and fairness of content-creation evaluation, we conduct an automatic evaluation pipeline. Evaluations of various MLLMs on MCTBench reveal that, despite their impressive perceptual capabilities, their cognition abilities require enhancement. We hope MCTBench will offer the community an efficient resource to explore and enhance cognitive capabilities towards text-rich visual scenes.

025 026 027

004

010 011

012

013

014

015

016

017

018

019

021

1 INTRODUCTION

Multimodal Large Language Models (MLLMs) OpenAI (2023); Team et al. (2023); Liu et al. (2023b); Li et al. (2024b) have exhibited promising performance across various cross-modal tasks, and revealed potential for widespread real-world applications. In practical applications, many images contain crucial textual elements that are essential for addressing specific challenges, such as key information extraction from receipts. Consequently, the ability to comprehend text-rich visual scenes can significantly enhance the practicality of MLLMs and drive innovative applications across multiple domains.

Recent benchmarks Liu et al. (2024c); Li et al. (2024a); iang Yue et al. (2023) have increasingly focused on evaluating MLLMs towards text-rich visual scenes. Nonetheless, the benchmarks are centred around evaluating perceptual capabilities yet overlook the assessment of cognitive abilities, which are a significant strength of MLLMs (as illustrated in Figure 1).

In this paper, we propose a Multimodal benchmark to evaluate the Cognitive capabilities of MLLMs 040 in Text-rich visual scenes (MCTBench). To assess the cognitive abilities of MLLMs thoroughly, 041 we design two types of tasks in the MCTBench: reasoning tasks for comprehension of the input 042 scenes, and open-ended content-creation tasks for generating output responses. Besides, MCTBench 043 integrates various perception tasks to study the differences with cognition tasks, while avoiding evalu-044 ation biases from varying dataset distributions. Fundamentally, MCTBench curates approximately 5.2k text-rich images from a wide range of public datasets, along with 8.5k rigorously annotated 046 question-answer pairs categorized into three tasks: perception, reasoning and content-creation. The 047 perception and reasoning tasks are formatted as multiple-choice questions for convenient evaluation, 048 following common practices in Fu et al. (2024); Liu et al. (2024b); Li et al. (2023b). Due to the subjectivity and high cost of human evaluation in open-ended content creation, we establish an automated evaluation pipeline by leveraging sophisticated MLLMs (e.g., GPT-4V) as the evaluator, to compare 051 the predictions of models against the provided references. Our experimental results demonstrate that MLLMs exhibit notably lower performance of cognition capabilities compared to perception 052 in text-rich visual scenes, particularly for text-enhanced models. Furthermore, performances in cognition tasks (reasoning and content-creation) are improved with larger parameter scales. Our main



MCTBench: <u>Cognition</u> (Reasoning, Content-Creation) + <u>Perception</u>

Figure 1: The Comparison between previous Benchmarks Singh et al. (2019); Liu et al. (2024c); Li et al. (2024a), and our proposed MCTBench. **Q** and **GT** stand for question and ground truth.

contributions are summarized as follows:

- 1. We propose a brand-new and large-scale benchmark for evaluating the cognitive capability of MLLMs towards text-rich visual scenes.
- 2. The evaluation on MCTBench highlights that MLLMs necessitate enhancements in their cognitive capabilities in text-rich visual scenes.
- 3. We develop an automated evaluation pipeline for the content-creation task, offering researchers an efficient tool for further investigation of cognitive capabilities.

2 RELATED WORK

075 076

077

079 080

081

084

085

090

092

2.1 MULTIMODAL LARGE LANGUAGE MODELS

093 The significant advancements in Large Language Models (LLMs) OpenAI (2023); Touvron et al. 094 (2023); Chiang et al. (2023) have paved the way for recent research Team et al. (2023); Bai et al. 095 (2023); Liu et al. (2023b); Chen et al. (2023a); Dai et al. (2023) into developing Multimodal Large 096 Language Models (MLLMs) that integrate visual capabilities. Early works in this field Alayrac et al. (2022); Li et al. (2023c); Liu et al. (2023b); Chen et al. (2023a) have introduced various vision-098 language projectors such as Q-formerDai et al. (2023), Multi-Layer Perceptron (MLP) Liu et al. 099 (2023b), and PerceiverAlayrac et al. (2022), which act as intermediaries between LLMs and visual encoders. Furthermore, these efforts have also established robust training paradigms for MLLMs. 100 Building upon these foundational paradigms, recent initiatives Chen et al. (2023b); Lu et al. (2024a); 101 Liu et al. (2023a); McKinzie et al. (2024) have focused on scaling the quality of training data, to 102 enhance general visual capabilities effectively. 103

A primary challenge in the recent development of MLLMs is attaining fine-grained comprehension,
exemplified by tasks such as Visual Question Answering (VQA) on text-rich images. To address this
issue, increasing the resolution and integrating fine-grained visual features have been proven effective
across various studies Feng et al. (2023); Liu et al. (2024d;a); Hu et al. (2023); Ye et al. (2023b;a); Li
et al. (2024c); Wei et al. (2023). Additionally, works such as Feng et al. (2023); Hu et al. (2023);

Benchmark	Text-Rich Oriented	#Image	e #QAs I	Percetion	Reasoning	Gontent Creatior	Answer Type
MME Fu et al. (2024)	×	1137	2.2K	~	~		Yes/No
MMBench Liu et al. (2024b)	×	-	3K	~	~		MC
OCRBench Liu et al. (2024c)	v	450	1K	~			Open
SEED-bench-2-plus Li et al. (2024a)	v	-	2.3K	~			ŃС
Contextual Wadhawan et al. (2024)	✓	506	506	~	✓		Open
MMMU iang Yue et al. (2023)	✓	-	11.5K	~	~		MC/Open
MCTBench	~	5.2K	8.5K	~	~	~	MC/Open
							-

Table 1: The comparison between MCTBench and previous benchmarks. **Open** and **MC** respectively present open-ended and multiple choice format for answer type. **QAs** stands for question-answer pairs. Text-Rich Oriented indicates whether the benchmark focuses on text-rich visual scenes.



Figure 2: The pipeline of constructing MCTBench.

Zhang et al. (2024); Li et al. (2024c); Tang et al. (2024) have incorporated high-quality, text-rich visual tuning data to refine these models further.

2.2 MLLM BENCHMARKS

As multimodal large language models (MLLMs) continue to exhibit cross-task generality, single-task
evaluations (e.g., Goyal et al. (2017); Singh et al. (2019); Chen et al. (2015); Lin et al. (2024)) are
inadequate for a comprehensive performance assessment. Recent works Liu et al. (2024b); Li et al.
(2023b); Fu et al. (2024); Yu et al. (2023) present general MLLM benchmarks comprising multiple
tasks. Furthermore, to explore the performance of MLLMs on more complex tasks, MathVista Lu
et al. (2024b) evaluates their mathematical abilities, and MMMU iang Yue et al. (2023) integrates
multiple-discipline questions to benchmark MLLMs in expert domains.

Conversely, text-rich visual scenes are attracting growing attention due to their potential applications. Early works Mathew et al. (2021); Mishra et al. (2019a); Singh et al. (2019) focused on single tasks, while OCRBench Liu et al. (2024c) integrates multiple single-task datasets into five representative OCR(Optical Character Recognition)-based tasks. In contrast, our work evaluates MLLMs on complex tasks beyond OCR-based ones in text-rich visual scenes. A similar work is presented in Wadhawan et al. (2024), which demonstrates the model's performance on reasoning tasks but only on a limited set of test datasets. Our study provides a broader evaluation of cognition in text-rich visual scenes, pushing the boundaries of what MLLMs can achieve in more diverse scenarios such as content-creation. Table 1 demonstrates the detailed comparison between ours and previous benchmarks.

3 MCTBENCH

In this section, we outline the process of constructing the MCTBench. Section 3.1 provides an
 overview of MCTBench and compares it with previous benchmarks. Section 3.2 describes the
 procedure of collecting text-rich image sources from publicly accessible datasets. Finally, Section 3.2
 explains the annotation process applied to the collected images.

162 3.1 OVERVIEW

164 The MCTBench is designed to evaluate the cognitive capabilities of Multimodal Large Language 165 Models (MLLMs) towards text-rich visual scenes. To construct the comprehensive and diverse benchmark, we collected 5,194 images from a variety of public datasets, encompassing a wide 166 array of text-rich scenes such as natural environments, books, scientific contexts, advertisements, 167 e-commerce, and video shots. We meticulously annotate these images with a total of 8.5k question-168 answer pairs categorized into three tasks: perception, reasoning, and content creation. Specifically, MCTBench consists of 2,734 perception multiple-choice samples, 2,602 reasoning multiple-choice 170 samples, and 3,130 content-creation samples. Figure 2 illustrates the overall construction pipeline of 171 MCTBench. 172

173 174

3.2 TEXT-RICH IMAGES COLLECTION

Image Source The images of MCTBench are collected from 10 different publicly available datasets, aiming to incorporate comprehensive visual scenes to evaluate the cognition of Multimodal Large Language Models (MLLMs). We begin with sampling common general natural scenes (e.g., street views, competitions, road signs) from the COCO Lin et al. (2015), Flickr30k Young et al. (2014), GQA Hudson & Manning (2019), SeedBench Li et al. (2023b) and Visual Genome Krishna et al. (2016) datasets.

Furthermore, we select conventional text-rich multimodal datasets: OCR-VQA Mishra et al. (2019b), and VizWiz Gurari et al. (2018) taken by blind photographers. To further diversify MCTBench, we incorporate three domain-specific scenes with broad potential applications: advertising (AutoUnderAds Hussain et al. (2017)), e-commerce (FoodLogoDet-1500 Hou et al. (2021)), and science (ScienceQA Lu et al. (2022)). We randomly extract one frame from each video in the AutoUnderAds dataset. All data sources are specifically selected from the testsets. We adhere to the original licenses stated by all datasets.

188 **OCR-based image filtering**. We select the text-rich images from the sourced images, which are 189 guided by the following guidelines. To maintain the clarity and substantive textual content, we only 190 retain images with valid OCR-recognized characters (with recognition probabilities higher than 0.2) 191 of at least 10 characters. To ensure text contributes to overall visual semantics, we select images 192 where text regions occupy more than 10% of the image area, after validating the impact of valid text 193 lines on semantic expression. These meticulous selection criteria resulted in a curated collection of 194 high-quality and crystal-clear text-rich images, designed to challenge and inspire advancements in perceptual and cognitive understanding within textual domains. 195

196

197 **Annotation** Considering the bias and efficiency of the manual annotation, we employ a GPT-aided 198 approach to generate at least 10 pseudo-questions for each image, and ask annotators to remove low-quality ones. All answers are human-annotated with two rounds: (1) Each image with at least 10 199 GPT-aided pseudo-questions, is randomly assigned to three annotators. Each annotator independently 200 annotates the questions and provides answers. (2) Quality checkers will review the annotation in 201 the first round. If any question or image quality does not meet our annotation guidelines, the set is 202 re-annotated by the corresponding annotators, who also revise their answers. Annotators in the second 203 round are required to have at least 2 years of experience in text-rich multimodal scene annotation. 204 To reach an agreement, we use majority voting to determine the final answer. If majority voting 205 does not reach an agreement (i.e., all three answers are inconsistent), we check if the discrepancy 206 originated in the second round. If so, the question is re-annotated; if inconsistencies persist, it is 207 discarded. The question is discarded if the discrepancy does not arise in the second round. In addition 208 to content-creation tasks, due to the inherent diversity of responses, we do not provide unified answers. Instead, we offer standard references generated by powerful MLLMs (e.g., GPT-4V) and meticulously 209 210 reviewed by humans.

211

Quality control During the image quality assessment, annotators remove low-quality images (e.g.,
 blurry, unclear text, inappropriate, solely tables/documents, and only watermark). For the QA quality
 assessment, annotators eliminate low-quality questions (e.g., ambiguous, overly generalized, too
 simplistic) and check the correctness of annotated answers (e.g., logical errors). On the other hand,
 We filter out multiple-choice options with more than 30% word count disparity and remove the



Figure 3: Visualization of the question for three different tasks using word clouds. In the word cloud, the size of a word indicates how frequently it appears. Best viewed in color.

questions that GPT-4 refuses to answer due to ethical concerns. To detail the specific issues addressed by each type, we assigned a label to each question and visualized the terms in Figure 3.

3.3 DATA CONSTRUCTION

4 EXPERIMENTS

In this section, we conduct comprehensive experiments to evaluate current MLLMs. Firstly, we outline experimental settings for the evaluated models and metrics used on MCTBench in Section 4.1 and 4.2. The results obtained from these experiments allow us to perform analysis of the selected models across three categories of tasks in Section 4.3. Furthermore, we conduct a case study to investigate performance variations among diverse models in Section 4.5.

4.1 MODELS

To evaluate the performance of various MLLMs on MCTBench, we select a diverse range of models categorized into two primary types: **general MLLMs** and **text-enhanced MLLMs** (specifically optimized for text recognition in images). We establish two naive baselines (random choice and frequency choice) as reference points to ensure the robustness and validity of the dataset in Table 2. Random choice involves selecting an option at random as the prediction result, while frequency choice involves selecting the prediction result based on the option with the highest proportion in the ground truth.

256

228

229

230 231 232

233

234 235 236

237 238

239 240

241

242

243

244

245 246 247

248

257 **General MLLMs** The experiment initially evaluates popular closed-source models, specifically 258 Gemini-Pro Team et al. (2023) and GPT-4V(ision) OpenAI (2023). For open-source models, we select several notable general-purpose MLLMs, including Sharegpt4V Chen et al. (2023b), Honeybee Cha 259 et al. (2024), LLaVA Liu et al. (2024a; 2023b), Otter Li et al. (2023a), Yi-VL AI et al. (2024), Qwen-260 VL-chat Bai et al. (2023) and Deepseek-VL Lu et al. (2024a) as competitive baselines. Additionally, 261 we incorporate CogVLM Wang et al. (2023) and SPHINX-v2 Lin et al. (2023), which are enhanced 262 for fine-grained understanding. To assess differentiated performance, we also integrate a larger 263 open-source model, Mini-Gemini Li et al. (2024b). 264

265

Text-enhanced MLLMs Recent researchers propose remarkable works to tackle the understanding of text-rich images via enhancing the textual capabilities. Consequently, we select models mPLUG-DocOwl Ye et al. (2023b), Monkey Liu et al. (2024d); Li et al. (2024c), InternLM-XComposer2-VL Dong et al. (2024), CogAgent Hong et al. (2023) and LLaVA-NeXT Liu et al. (2024a) which have demonstrated strong OCR capabilities in previous evaluations.

270 4.2 METRICS 271

276

MCTBench is constructed with three categories of tasks: perception, reasoning, and content-creation.
Due to the standard answer provided in perception and reasoning tasks, each QA pair is designed
in a multiple-choice format. In contrast, the content-creation task is considered as an open-ended
generation problem due to the diversity of responses.

Perception and reasoning Perception and reasoning tasks entail the acquisition of information from input data, comprehension of images and text, and derivation of conclusions. Consequently, employing multiple-choice question-answering can effectively validate the corresponding capabilities of MLLMs, following Liu et al. (2024b); Li et al. (2023b). In practice, we prioritize original prompts employed by MLLMs and, if not specified, use those prompts which yield optimal results. For lengthy responses, we use regular expressions and supplementary rules to extract the option answers. We use mean accuracy to evaluate MLLMs' perception and reasoning capability.

284 **Content-creation** To ensure consistency and efficiency in evaluation, we implement automatic 285 evaluation for content-creation tasks. However, due to the diversity of answers in content creation, 286 standard responses are not feasible, leading us to establish competitive references. These references 287 are generated through manually crafted responses from text-only GPT-4, based on inputs from OCR 288 recognition and detailed descriptions. The evaluation is grounded on four principal aspects: relevance, faithfulness, creativity, and instruction following. Subsequently, we employ machines (e.g., GPT-4V) 289 to compare other MLLMs against our references with the mentioned principal, and categorize their 290 performance as Good, Same, or Bad (i.e., the GSB metric). We measure each model's performance 291 by calculating the percentage of 'Good' and 'Same' ratings relative to all questions, indicating how 292 many outperform or match the constructed reference (i.e., (G+S)/(G+S+B)). 293

To assess correlations and validate the reliability of machine evaluation, we also conduct a manual evaluation on subsets of the data using the mentioned metrics (GSB) and compare them with machine evaluation. We integrate three powerful MLLMs (GPT-4V OpenAI (2023), Gemini-Pro Team et al. (2023) and LLaVA-NeXT Liu et al. (2024a)) as evaluators. We measure the evaluation correlation between humans and machines using accuracy and Pearson correlation Benesty et al. (2009) on GSB scores. Table 3 illustrates the results between three machine evaluators and human evaluations, indicating that GPT-4V achieves a top correlation score.

4.3 RESULTS

303 **Perception** Firstly, experiments are conducted to verify the perceptual performances of each 304 model as baselines. As shown in Table 2, most models achieved satisfying scores on the perception 305 task. Closed-source models (e.g., GPT-4V OpenAI (2023)) demonstrated excellent accuracy, while 306 some open-source models (e.g., Mini-Gemini Li et al. (2024b)) demonstrated superior perception 307 capabilities, surpassing the their performance. Models with higher resolutions and more parameters (e.g., LLaVA-NeXT Liu et al. (2024a) and Mini-Gemini Li et al. (2024b)), typically performed better. 308 Among similarly-sized models, text-enhanced MLLMs generally outperformed others by effectively 309 extracting text from images and generating precise responses. 310

311

301

302

Reasoning The reasoning task is more challenging than perception. It requires not only effective extraction and fusion of visual and textual features, but also involves comprehensive inference to generate accurate responses. As shown in Table 2, most models exhibited a significant drop in scores due to the increased difficulty of reasoning tasks compared to perception tasks.

Notably, GPT-4V demonstrates exceptional performance among MLLMs, surpassing most models by
 a significant margin. Besides, there still exists a positive correlation between a model's performance
 and the number of its parameters. This phenomenon arises from larger models' enhanced ability, to
 comprehend text and integrate image information for reasoning more effectively.

Nevertheless, text-enhanced MLLMs have not substantially outperformed general models on reasoning tasks. Given that text-enhanced models have achieved superior results in perception tasks, we posit their effectiveness in recognizing text within images. However, achieving higher scores in reasoning tasks necessitates the ability to analyse and summarise effectively. TextMonkey Liu et al. (2024d) shows the least performance gap and achieves results comparable to the perception

355 356

357

366

				Cogn	ition	Ave	rage Sc	ores
Model		Params	Perception	Reasoning	Content- Creation*	MC	Cog	All
		N	aive Baselin	e				
Random choic	e	-	25.00	25.00	-	-	-	-
Frequency cho	ice	-	25.16	25.52	-	-	-	-
		Ge	neral MLLN	4s				
GPT-4V Open	AI (2023)	-	83.58	74.21	87.35	78.90	83.12	81.71
Gemini-Pro Te	eam et al. (2023)	-	78.79	70.18	56.78	74.49	65.63	68.58
Yi-VL AI et al	. (2024)	6B	77.25	72.33	41.45	74.79	58.12	63.68
Deepseek-VL	Lu et al. (2024a)	7B	76.74	68.79	57.25	72.77	65.01	67.59
Honeybee Cha	et al. (2024)	7B	72.60	67.22	73.64	69.91	71.78	71.15
Otter Li et al. ((2023a)	7B	58.12	54.42	31.70	56.27	43.99	48.08
Qwen-VL-cha	t Bai et al. (2023)	7B	77.98	70.68	67.53	74.33	70.93	72.06
Sharegpt4V Cl	nen et al. (2023b)	13B	74.54	69.49	66.19	72.02	69.10	70.07
LLaVA-1.5 Liu	u et al. (2023b)	13B	78.09	72.56	66.47	75.33	70.90	72.37
SPHINX-v2 L	in et al. (2023)	13B	78.02	71.94	62.30	74.98	68.64	70.75
CogVLM Wan	g et al. (2023)	17B	71.40	69.52	65.61	70.46	68.04	68.84
Mini-Gemini I	Li et al. (2024b)	34B	83.83	<u>73.33</u>	<u>86.76</u>	78.58	<u>82.67</u>	<u>81.31</u>
		Text-e	nhanced MI	LMs				
IXC 2 Dong et	al. (2024)	7B	78.05	72.10	74.45	75.08	<u>74.76</u>	74.87
Monkey Li et a	al. (2024c)	7B	<u>79.22</u>	72.64	59.56	75.93	67.75	70.47
TextMonkey L	iu et al. (2024d)	7B	71.80	69.45	22.81	70.63	46.72	54.69
mPLUG-DocC	Owl Ye et al. (2023a)	10B	75.05	70.06	60.87	72.56	66.71	68.66
CogAgent Hor	ng et al. (2023)	34B	58.56	56.46	56.86	57.51	57.19	57.29
LLaVA-NeXT	Liu et al. (2024a)	34B	83.87	71.64	85.30	77.76	81.53	80.27

Table 2: Evaluation results for MLLMs on MCTBench. MC means the average score of the two tasks (perception and reasoning) in multiple-choice format. Cog means the average scores of the two cognitive tasks (reasoning and content-creation). All means the overall average scores of all tasks.
*The content-creation task is scored using the percentage of 'Good' and 'Same' ratings by the GSB metric except for accuracy used in Perception and reasoning tasks. Numbers in Bold and underline represent the top-2 results in each task.

task. However, most text-enhanced models are not explicitly trained in this aspect, and consequently do not outperform general MLLMs.

	GPT-4V	Gemini-Pro	Best open-source (LLaVA-NeXT)
Accuracy	79.38	70.71	65.22
Pearson Correlation	0.558	0.380	0.304

Table 3: Correlation analysis between automatic machine and human evaluation on content-creation using the accuracy and Pearson correlation coefficient.

367 **Content-creation** The open-ended creation task differs from the aforementioned tasks, leading 368 to the observation of a broader array of perspectives. As mentioned in Section 4.2, the evaluation 369 is based on four principal aspects: relevance, faithfulness, creativity, and instruction following. 370 Therefore, as the content-creation task emphasizes the generation of suitable text for images, we also 371 notice significantly enhancement on performances by employing larger language models. However, 372 the performance of text-enhanced MLLMs exhibits considerable variation due to differences in 373 training objectives. Some models are specifically trained to extract structured text information or 374 comprehend lengthy text inputs. Consequently, they may underperform relative to general models 375 when tasked with creative endeavours. We also found that the closed-source model Gemini-Pro Team et al. (2023) is unable to achieve good results. A notable discrepancy in its scores of template 376 following demonstrates that Gemini struggles to create corresponding formats for task specifications. 377 For instance, in the task of generating a slogan, Gemini may tend to produce lengthy paragraphs

7

378	Input	Perception	Reasoning	Creation
379	Image + OCR texts	75.60	71.68	65.45
380	Image (text regions removed)	62.22	62.22	46.89
381	Image (text regions only)	74.69	68.01	58.78
382	Image (baseline)	78.09	72.56	66.47

Table 4: The effect of visual and textual information.

instead of concise phrases. Some text-enhanced MLLMs, such as TextMonkey, also have similar phenomena.

390 **Summary** MLLMs have shown certain multi-modal capabilities in perception tasks, with mainstream models achieving commendable performances. However, in reasoning tasks, existing models still have room for improvement. Performance on content-creation tasks indicates that text-enhanced 392 MLLMs trained for different types of tasks may lose some creative capabilities. We suggest that 393 reasoning and creation tasks serve as distinct dimensions for evaluation, offering insights into the 394 model's comprehension of input and its proficiency in generating output responses. While existing models excel in basic perception tasks, achieving comparable competence in reasoning and creation 396 tasks remains challenging.

397 398 399

400

384 385 386

387

388 389

391

4.4 ABLATION STUDY

401 The effect of visual and textual information. In Table 4, we conduct several experiments to reveal the effect of visual and textual information in images of our benchmark. We use the LLaVA-1.5 402 as the baseline. Firstly, we conduct an experiment that adds OCR texts as input. As shown in the 403 table, even if the OCR texts are contained as input, the model cannot gain explicit improvement. 404 This demonstrates that the MCTBench does not rely solely on texts to get answers. Furthermore, 405 we conduct an experiment that removes all texts from images to check the importance of texts in 406 our benchmark. Specifically, we detected and blurred all the texts in images. We noticed that the 407 performance dropped significantly. This also demonstrates that it is difficult for a model to correctly 408 answer the question without text in the image. Additionally, if we only keep the text region, and 409 delete other background parts in the image. The performance is between the above two experiments, 410 indicating that MCTBench relies both on textual and visual information, to get the final result. To 411 conclude, explicitly adding OCR texts, or removing text/image parts does not help, or even lead to 412 worse performance on MCTBench. MLLMs are required to jointly recognize related textual and visual patterns, to answer the questions correctly. 413

414

415 4.5 CASE STUDY 416

417

For the perception and reasoning tasks, we select representative cases shown in Figure 4. Specifically, 418 we select GPT-4V as a strong reference along with several open-source representative MLLMs, 419 split into three groups: large model size and resolution MLLMs (Mini-Gemini Li et al. (2024b) and 420 LLaVA-NeXT Liu et al. (2024a)), text-enhanced MLLMs (Monkey Li et al. (2024c) and mPLUG-421 DocOwl Ye et al. (2023a)), and general MLLMs (LLaVA-1.5 Liu et al. (2023a) and ShareGPT4V 422 Chen et al. (2023b)). For the selected perception question, both high-resolution MLLMs and text-423 enhanced MLLMs perform well, while general MLLMs fail in fine-grained understanding. On the contrary, text-enhanced MLLMs excelled in perception but performed poorly in reasoning tasks. 424 Larger models like Mini-Gemini Li et al. (2024b), LLaVA-NeXT Liu et al. (2024a), and GPT-4V 425 OpenAI (2023) handle reasoning better by effectively integrating textual and visual elements. 426

427 For the content-creation task, using GPT-4V OpenAI (2023) as a robust reference, we select general and text-enhanced models that performed well on MCTBench, and conduct case studies in three 428 scenarios. As shown in Figure 5, GPT-4V OpenAI (2023) significantly surpassed other models in 429 content creation quality. Mini-Gemini Li et al. (2024b) also showed consistent performance across 430 cognitive tasks, while text-enhanced models like Monkey Li et al. (2024c) were limited to text 431 recognition and basic descriptions, resulting in less attractive content.

Image	11- TOM PRICE 1311- 11- ROEBOURNE 297 - MANUARAN 399 - JULISTICAM			Be Good to Your Goblins /	DRIVE LIKE: Year Grandpa is, in the Crasswalk Remning of Statis Sills
	Identificat	tion	Anno R R R R R R R R R R R R R R R R R R	Reasoning	
Scene	Sign Recognition	Brand identification	Key information Understadnding	g Target audience predication	Emotion analysis
	What is the <u>estimated</u> <u>distance</u> to Roebourne?	What is the <u>name of the</u> <u>product</u> represented in the image?	Does this juice product cater to the needs of <u>vegetarians</u> ?	Who is most likely <u>the target</u> <u>audience</u> for this advertisement?	What emotion is the billboard trying to invoke in drivers?
Question	A. 287 kilometers	A. Dare	A. Yes	A. Children	A. Fear
	B. 189 kilometers	B. Male smoke	B. No	B. Parents	B. Joy
	C. 382 kilometers	C. 3 hours of heat		C. Teenagers	C. Guilt
	D. 131 kilometers	D. Craving relief		D. Elderly	D. Caution
Groud Truth	A. 287 kilometers	A. Dare	A. Yes	B. Parents	D. Caution
GPT-4V	C. 382 kilometers 🗱	C. 3 hours of heat ×	A. Yes	B. Parents	D. Caution
MGM	A. 287 kilometers	A. Dare	A. Yes	B. Parents	D. Caution
LLavA-Nex I	A. 287 kilometers	A. Dare	A. Ies	B. Parents	A Foor #
mPLUG-doc	A 287 kilometers	A Dare	A. Ies A Ves	A Children	A Fear
	A. 207 KIIOIIIEters	A. Daile	A. IES	A. children	A. Feat
LLaVA-15	D 131 kilometers 🗶	C 3 hours of heat 🗶	B No 🗶	A Children 🗶	C Guilt 🗶

Figure 4: The cases of predication from different MLLMs divided into four groups: GPT-4V OpenAI (2023), Mini-Gemini Li et al. (2024b)(MGM) and LLaVA-NeXT Liu et al. (2024a) for larger model size, Monkey Li et al. (2024c) and mPLUG-DocOwl Ye et al. (2023a) for text-enhanced MLLMs, LLaVA-1.5 Liu et al. (2023a) and ShareGPT4V Chen et al. (2023b) for the general MLLMs

5 LIMITATIONS

Our dataset primarily focuses on English, which may limit the generalization of our findings to multilingual scenes. Although we believe that the cognitive capacities of MLLMs should theoretically extend to other languages, we have not empirically substantiated this assertion in the present study. Additionally, we have only selected a subset of representative models for evaluation due to space constraints. This selection may not cover the full spectrum of currently available MLLMs. Our future work aims to provide evaluation results for a more extensive range of models.

6 CONCLUSION

In this work, we introduce MCTBench, a comprehensive benchmark designed to evaluate the cognitive capabilities of MLLMs in text-rich visual scenes. The MCTBench comprises 5.2k images and 8.5k question-answer pairs, covering a range of tasks including reasoning, content creation for cognitive assessment, and conventional perception. Evaluations of MLLMs on MCTBench reveal that current MLLMs still need further advancements in cognitive capabilities, despite their superior perception performance. We hope that MCTBench will motivate researchers to further improve the cognitive capabilities of MLLMs in text-rich visual scenes, thereby enhancing the practical utility of AI in real-world applications.

478 REFERENCES

01. AI, :, 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, 2024.

- 485 Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katie Millican, Malcolm Reynolds, Roman Ring, Eliza Rutherford, Serkan



Figure 5: The cases of predication on content-creation tasks from three representative MLLMs: GPT-4V OpenAI (2023), Mini-Gemini Li et al. (2024b) and Monkey Li et al. (2024c). We mark high-quality sentences in red, words hit the text in the image with underlining, and rate the quality of the generation with stars.

550

569

578

579

580

581

585

586

587

588

591

540 Cabi, Tengda Han, Zhitao Gong, Sina Samangooei, Marianne Monteiro, Jacob Menick, Sebastian 541 Borgeaud, Andrew Brock, Aida Nematzadeh, Sahand Sharifzadeh, Mikolaj Binkowski, Ricardo 542 Barreira, Oriol Vinyals, Andrew Zisserman, and Karen Simonyan. Flamingo: a visual language 543 model for few-shot learning, 2022. 544 Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang 545 Zhou, and Jingren Zhou. Qwen-vl: A frontier large vision-language model with versatile abilities. 546 arXiv preprint arXiv:2308.12966, 2023. 547 548 Jacob Benesty, Jingdong Chen, Yiteng Huang, and Israel Cohen. Pearson correlation coefficient. In 549 Noise reduction in speech processing, pp. 37–40. Springer, 2009.

- Junbum Cha, Wooyoung Kang, Jonghwan Mun, and Byungseok Roh. Honeybee: Locality-enhanced
 projector for multimodal llm. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2024.
- Jun Chen, Deyao Zhu, Xiaoqian Shen, Xiang Li, Zechun Liu, Pengchuan Zhang, Raghuraman Krishnamoorthi, Vikas Chandra, Yunyang Xiong, and Mohamed Elhoseiny. Minigpt-v2: large language model as a unified interface for vision-language multi-task learning. *arXiv preprint arXiv:2310.09478*, 2023a.
- Lin Chen, Jisong Li, Xiaoyi Dong, Pan Zhang, Conghui He, Jiaqi Wang, Feng Zhao, and Dahua
 Lin. Sharegpt4v: Improving large multi-modal models with better captions. *arXiv preprint arXiv:2311.12793*, 2023b.
- 562
 563
 564
 Xinlei Chen, Hao Fang, Tsung-Yi Lin, Ramakrishna Vedantam, Saurabh Gupta, Piotr Dollar, and C. Lawrence Zitnick. Microsoft coco captions: Data collection and evaluation server, 2015.
- Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality, March 2023. URL https://lmsys.org/blog/2023-03-30-vicuna/.
- Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale Fung, and Steven Hoi. Instructblip: Towards general-purpose vision-language models with instruction tuning, 2023.
- 573 Xiaoyi Dong, Pan Zhang, Yuhang Zang, Yuhang Cao, Bin Wang, Linke Ouyang, Xilin Wei, Songyang
 574 Zhang, Haodong Duan, Maosong Cao, Wenwei Zhang, Yining Li, Hang Yan, Yang Gao, Xinyue
 575 Zhang, Wei Li, Jingwen Li, Kai Chen, Conghui He, Xingcheng Zhang, Yu Qiao, Dahua Lin, and Ji576 aqi Wang. Internlm-xcomposer2: Mastering free-form text-image composition and comprehension
 577 in vision-language large model. *arXiv preprint arXiv:2401.16420*, 2024.
 - Hao Feng, Zijian Wang, Jingqun Tang, Jinghui Lu, Wengang Zhou, Houqiang Li, and Can Huang. Unidoc: A universal large multimodal model for simultaneous text detection, recognition, spotting and understanding, 2023.
- Chaoyou Fu, Peixian Chen, Yunhang Shen, Yulei Qin, Mengdan Zhang, Xu Lin, Jinrui Yang, Xiawu
 Zheng, Ke Li, Xing Sun, Yunsheng Wu, and Rongrong Ji. Mme: A comprehensive evaluation
 benchmark for multimodal large language models, 2024.
 - Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. Making the V in VQA matter: Elevating the role of image understanding in Visual Question Answering. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- Danna Gurari, Qing Li, Abigale J. Stangl, Anhong Guo, Chi Lin, Kristen Grauman, Jiebo Luo, and
 Jeffrey P. Bigham. Vizwiz grand challenge: Answering visual questions from blind people, 2018.
- Wenyi Hong, Weihan Wang, Qingsong Lv, Jiazheng Xu, Wenmeng Yu, Junhui Ji, Yan Wang, Zihan
 Wang, Yuxiao Dong, Ming Ding, et al. Cogagent: A visual language model for gui agents. *arXiv* preprint arXiv:2312.08914, 2023.

594 595 596 597 598	Qiang Hou, Weiqing Min, Jing Wang, Sujuan Hou, Yuanjie Zheng, and Shuqiang Jiang. Foodlogodet- 1500: A dataset for large-scale food logo detection via multi-scale feature decoupling network. In <i>Proceedings of the 29th ACM International Conference on Multimedia</i> , MM '21. ACM, October 2021. doi: 10.1145/3474085.3475289. URL http://dx.doi.org/10.1145/3474085. 3475289.
599 600 601 602	Wenbo Hu, Y. Xu, Y. Li, W. Li, Z. Chen, and Z. Tu. Bliva: A simple multimodal llm for better handling of text-rich visual questions. In AAAI Conference on Artificial Intelligence, 2023. URL https://api.semanticscholar.org/CorpusID:261049015.
603 604	Drew A. Hudson and Christopher D. Manning. Gqa: A new dataset for real-world visual reasoning and compositional question answering, 2019.
605 606 607 608 609	Zaeem Hussain, Mingda Zhang, Xiaozhong Zhang, Keren Ye, Christopher Thomas, Zuha Agha, Nathan Ong, and Adriana Kovashka. Automatic understanding of image and video advertisements. 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp. 1100–1110, 2017. URL https://api.semanticscholar.org/CorpusID:11172071.
610 611 612 613 614	iang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, Cong Wei, Botao Yu, Ruibin Yuan, Renliang Sun, Ming Yin, Boyuan Zheng, Zhenzhu Yang, Yibo Liu, Wenhao Huang, Huan Sun, Yu Su, and Wenhu Chen. Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi, 2023.
615 616 617 618	Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A. Shamma, Michael S. Bernstein, and Fei-Fei Li. Visual genome: Connecting language and vision using crowdsourced dense image annotations, 2016.
619 620 621	Bo Li, Yuanhan Zhang, Liangyu Chen, Jinghao Wang, Jingkang Yang, and Ziwei Liu. Otter: A multi-modal model with in-context instruction tuning. <i>ArXiv</i> , abs/2305.03726, 2023a. URL https://api.semanticscholar.org/CorpusID:258547300.
622 623	Bohao Li, Rui Wang, Guangzhi Wang, Yuying Ge, Yixiao Ge, and Ying Shan. Seed-bench: Bench- marking multimodal llms with generative comprehension, 2023b.
625 626	Bohao Li, Yuying Ge, Yi Chen, Yixiao Ge, Ruimao Zhang, and Ying Shan. Seed-bench-2-plus: Benchmarking multimodal large language models with text-rich visual comprehension, 2024a.
627 628 629	Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models, 2023c.
630 631 632	Yanwei Li, Yuechen Zhang, Chengyao Wang, Zhisheng Zhong, Yixin Chen, Ruihang Chu, Shaoteng Liu, and Jiaya Jia. Mini-gemini: Mining the potential of multi-modality vision language models, 2024b.
633 634 635 636 637	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 <i>proceedings of the IEEE/CVF conference on computer vision and pattern recognition</i> , 2024c.
638 639 640	Kun-Yu Lin, Henghui Ding, Jiaming Zhou, Yi-Xing Peng, Zhilin Zhao, Chen Change Loy, and Wei-Shi Zheng. Rethinking clip-based video learners in cross-domain open-vocabulary action recognition. <i>arXiv preprint arXiv:2403.01560</i> , 2024.
641 642 643 644	Tsung-Yi Lin, Michael Maire, Serge Belongie, Lubomir Bourdev, Ross Girshick, James Hays, Pietro Perona, Deva Ramanan, C. Lawrence Zitnick, and Piotr Dollár. Microsoft coco: Common objects in context, 2015.
645 646 647	Ziyi Lin, Chris Liu, Renrui Zhang, Peng Gao, Longtian Qiu, Han Xiao, Han Qiu, Chen Lin, Wenqi Shao, Keqin Chen, Jiaming Han, Siyuan Huang, Yichi Zhang, Xuming He, Hongsheng Li, and Yu Qiao. Sphinx: The joint mixing of weights, tasks, and visual embeddings for multi-modal large language models, 2023.

648 649 650	Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning, 2023a.
651	Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning, 2023b.
652	Hastian Liu, Chunguan Li, Vuhang Li, Ro Li, Vuanhan Zhang, Shang Shang, Shang Jan Jan Jan Jan
653 654	Lava-next: Improved reasoning, ocr, and world knowledge, January 2024a. URL https:
655	//llava-vl.github.io/blog/2024-01-30-llava-next/.
656	Yuan Liu Haodong Duan, Yuanhan Zhang, Bo Li, Songyang Zhang, Wangho Zhao, Yike Yuan, Jiagi
657	Wang Conghui He Ziwei Liu Kai Chen and Dahua Lin Mmbench: Is your multi-modal model
658	an all-around player?, 2024b.
659	Yuliang Liu, Zhang Li, Biao Yang, Chunyuan Li, Xucheng Yin, Cheng lin Liu, Lianwen Jin, and
660 661	Xiang Bai. On the hidden mystery of ocr in large multimodal models, 2024c.
662	Yuliang Liu, Biao Yang, Ojang Liu, Zhang Li, Zhiyin Ma, Shuo Zhang, and Xiang Bai, Textmonkey:
663 664	An ocr-free large multimodal model for understanding document. <i>arXiv preprint arXiv:2403.04473</i> , 2024d
665	202 - 7 u .
666	Haoyu Lu, Wen Liu, Bo Zhang, Bingxuan Wang, Kai Dong, Bo Liu, Jingxiang Sun, Tongzheng Ren,
667	Zhuoshu Li, Hao Yang, Yaofeng Sun, Chengqi Deng, Hanwei Xu, Zhenda Xie, and Chong Ruan.
668	Deepseek-vl: Towards real-world vision-language understanding, 2024a.
669	Pan Lu, Swaroop Mishra, Tony Xia, Liang Oiu, Kai-Wei Chang, Song-Chun Zhu, Oyvind Tafjord,
670	Peter Clark, and Ashwin Kalyan. Learn to explain: Multimodal reasoning via thought chains for
671	science question answering. In The 36th Conference on Neural Information Processing Systems
672	(NeurIPS), 2022.
673	
674	Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng,
675	of foundation models in visual contexts 2024b
676	or roundation models in visual contexts, 20240.
677	Minesh Mathew, Dimosthenis Karatzas, and CV Jawahar. Docvqa: A dataset for vqa on document
678	images. In Proceedings of the IEEE/CVF winter conference on applications of computer vision,
679	pp. 2200–2209, 2021.
680	
681	Brandon McKinzle, Zhe Gan, Jean-Philippe Fauconnier, Sam Dodge, Bowen Zhang, Philipp Duiler,
682	from multimodal llm pre training arYiv preprint arYiv: 2403 00611, 2024
683	nom mutumodal nin pre-training. <i>urxiv preprint urxiv</i> .2405.09011, 2024.
684	Anand Mishra, Shashank Shekhar, Ajeet Kumar Singh, and Anirban Chakraborty. Ocr-vqa: Visual
685	question answering by reading text in images. In 2019 international conference on document
686	analysis and recognition (ICDAR), pp. 947–952. IEEE, 2019a.
687	Arond Michael Shachard, Shakhar Aiget Kumar Singh and Anishar Chalsenbarty, Oger ugge Visual
688	auestion answering by reading text in images. In ICDAP, 2010b
689	question answering by reading text in images. In <i>ICDAR</i> , 20190.
690	OpenAI. Gpt-4 technical report. 2023. URL https://api.semanticscholar.org/
691	CorpusID:257532815.
692	
693	Amanpreet Singh, Vivek Natarajan, Meet Shah, Yu Jiang, Xinlei Chen, Dhruv Batra, Devi Parikh,
694	and Marcus Rohrbach. Towards vqa models that can read. In <i>Proceedings of the IEEE/CVF</i>
695	Conjerence on Computer vision and Pattern Recognition (CVPR), June 2019.
696	Jinggun Tang, Chunhui Lin, Zhen Zhao, Shu Wei, Binghong Wu, Oi Liu, Hao Feng, Yang Li, Sigi
697	Wang, Lei Liao, et al. Textsquare: Scaling up text-centric visual instruction tuning. arXiv preprint
698	arXiv:2404.12803, 2024.
699	
700	Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu
701	Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. Gemini: a family of highly capable multimodal models. <i>arXiv preprint arXiv:2312.11805</i> , 2023.

702 703 704 705 706	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation language models. <i>ArXiv</i> , abs/2302.13971, 2023. URL https://api.semanticscholar.org/ CorpusID:257219404.
707 708 709	Rohan Wadhawan, Hritik Bansal, Kai-Wei Chang, and Nanyun Peng. Contextual: Evaluating context-sensitive text-rich visual reasoning in large multimodal models, 2024.
710 711 712	Weihan Wang, Qingsong Lv, Wenmeng Yu, Wenyi Hong, Ji Qi, Yan Wang, Junhui Ji, Zhuoyi Yang, Lei Zhao, Xixuan Song, et al. Cogvlm: Visual expert for pretrained language models. <i>arXiv</i> preprint arXiv:2311.03079, 2023.
713 714 715 716	Haoran Wei, Lingyu Kong, Jinyue Chen, Liang Zhao, Zheng Ge, Jinrong Yang, Jianjian Sun, Chunrui Han, and Xiangyu Zhang. Vary: Scaling up the vision vocabulary for large vision-language models, 2023.
717 718 719	Jiabo Ye, Anwen Hu, Haiyang Xu, Qinghao Ye, Ming Yan, Yuhao Dan, Chenlin Zhao, Guohai Xu, Chenliang Li, Junfeng Tian, Qian Qi, Ji Zhang, and Fei Huang. mplug-docowl: Modularized multimodal large language model for document understanding, 2023a.
720 721 722	Qinghao Ye, Haiyang Xu, Jiabo Ye, Ming Yan, Anwen Hu, Haowei Liu, Qi Qian, Ji Zhang, Fei Huang, and Jingren Zhou. mplug-owl2: Revolutionizing multi-modal large language model with modality collaboration, 2023b.
723 724 725 726	Peter Young, Alice Lai, Micah Hodosh, and Julia Hockenmaier. From image descriptions to visual denotations: New similarity metrics for semantic inference over event descriptions. <i>TACL</i> , 2:67–78, 2014.
727 728	Weihao Yu, Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Zicheng Liu, Xinchao Wang, and Lijuan Wang. Mm-vet: Evaluating large multimodal models for integrated capabilities, 2023.
729 730 731	Yanzhe Zhang, Ruiyi Zhang, Jiuxiang Gu, Yufan Zhou, Nedim Lipka, Diyi Yang, and Tong Sun. Llavar: Enhanced visual instruction tuning for text-rich image understanding, 2024.
732	
733	
734	
735	
736	
737	
738	
739	
740	
741	
743	
744	
745	
746	
747	
748	
749	
750	
751	
752	
752 753	
752 753 754	