# Gemini in Reasoning: Unveiling Commonsense in Multimodal Large Language Models

Anonymous ACL submission

### Abstract

The burgeoning interest in Multimodal Large Language Models (MLLMs), such as OpenAI's 003 GPT-4V(ision), has significantly impacted both academic and industrial realms. These models enhance Large Language Models (LLMs) with advanced visual understanding capabili-007 ties, facilitating their application in a variety of multimodal tasks. Recently, Google introduced 009 Gemini, a cutting-edge MLLM designed specifically for multimodal integration. Despite its 011 advancements, preliminary benchmarks indicate that Gemini lags behind GPT models in 013 commonsense reasoning tasks. However, this assessment, based on a limited dataset (i.e., HellaSWAG), does not fully capture Gemini's authentic commonsense reasoning potential. To address this gap, our study undertakes a thor-017 ough evaluation of Gemini's performance in complex reasoning tasks that necessitate the integration of commonsense knowledge across modalities. We carry out a comprehensive anal-022 ysis of 12 commonsense reasoning datasets, ranging from general to domain-specific tasks. This includes 11 datasets focused solely on language, as well as one that incorporates multi-026 modal elements. Our experiments across four LLMs and two MLLMs demonstrate Gemini's 027 competitive commonsense reasoning capabilities. We also highlight common challenges faced by current LLMs and MLLMs in commonsense reasoning, emphasizing the need for further advancements.

## 1 Introduction

Commonsense reasoning, integral to human cognition, plays a crucial role in navigating the intricacies of everyday life. Consider a scenario where someone decides what to wear based on the weather. This decision extends beyond the mere selection of attire; it involves understanding weather patterns, the suitability of clothing for different temperatures, and the social context of the occasion. It's about synthesizing diverse pieces of knowledge: a forecast predicting rain, the practical necessity for a raincoat, and the societal expectation of dressing appropriately for an event. This reasoning goes beyond simply processing information; it entails integrating varied pieces of knowledge that humans often take for granted. A major challenge in Natural Language Processing (NLP) research is the ambiguity and under-specification of human language. Individuals rely heavily on their commonsense knowledge and reasoning abilities to decipher these ambiguities and infer missing information. Commonsense reasoning has consistently posed unique challenges in NLP research (Li et al., 2021; Bian et al., 2023), encompassing spatial, physical, social, temporal, and psychological aspects, along with an understanding of social norms, beliefs, values, and the nuances of predicting and interpreting human behavior (Liu and Singh, 2004). Models often lack this innate commonsense, hindering their ability to contextualize data coherently, in stark contrast to the human capacity for effortlessly understanding everyday situations (Shwartz and Choi, 2020; Bhargava and Ng, 2022).

043

044

045

046

047

050

051

052

053

057

058

059

060

061

062

063

064

065

067

068

069

071

072

073

074

075

076

077

078

079

081

Recent advances in Large Language Models (LLMs) have sparked unprecedented enthusiasm in the NLP community and beyond, significantly enhancing a wide array of applications (Min et al., 2021; Zhao et al., 2023; Wang et al., 2023; Kasneci et al., 2023; He et al., 2023). Building on these achievements, Multimodal Large Language Models (MLLMs) have emerged as a pivotal focus in the next wave of AI (Wu et al., 2023b), speculated to advance towards Artificial General Intelligence (AGI), which aims to develop AI systems smarter than humans and beneficial for all of humanity (Rayhan et al., 2023). The rise of MLLMs, particularly OpenAI's GPT-4V(ision) (Yang et al., 2023) and Google's Gemini (Team et al., 2023), marks significant progress in this area. Among these developments, Gemini emerges as a formidable challenger to the state-

of-the-art MLLM, GPT-4V, specially engineered for multimodal integration. Its release has ignited constructive discussions about the current level of 086 AGI achievement. In widely used academic benchmarks, Gemini has attained new state-of-the-art status in the majority of tasks. However, preliminary evaluations of Gemini, especially when com-090 pared to models like the GPT series, have indicated potential shortcomings in its commonsense reasoning capabilities, a fundamental aspect of human cognition. Yet, it is important to consider that basing the assessment of Gemini's commonsense reasoning abilities solely on the HellaSWAG dataset (Zellers et al., 2019b) may not comprehensively reflect Gemini's full scope in this domain.

To address the gap in the comprehensive evaluation of Gemini's real-world performance in commonsense reasoning tasks, our study conducts ex-101 tensive experiments across 12 commonsense rea-102 103 soning datasets, covering a broad spectrum of domains such as general, physical, social, and tem-104 poral reasoning. The definitions of all tasks can 105 be found in Appendix A. We experiment with four 106 popular LLMs for the language dataset evaluation, 107 including Llama2-70b (Touvron et al., 2023), Gem-108 ini Pro (Team et al., 2023), GPT-3.5 Turbo, and 109 GPT-4 Turbo (OpenAI, 2023). For the multimodal 110 dataset, we assess both Gemini Pro Vision and 111 GPT-4V. Our key findings are summarized as fol-112 lows: (1) Overall, Gemini Pro's performance is 113 comparable to that of GPT-3.5 Turbo, demonstrat-114 ing marginally better average results across 11 lan-115 guage datasets (1.4% higher accuracy), though it 116 lags behind GPT-4 Turbo by an average of 8.2% in 117 accuracy. Moreover, Gemini Pro Vision exhibits 118 lower performance than GPT-4V on the multimodal 119 dataset, except for temporal-related questions. (2) 120 Approximately 65.8% of Gemini Pro's reasoning 121 processes are evaluated as logically sound and con-122 textually relevant, indicating its potential for effec-123 tive application in various domains. (3) Gemini Pro 124 encounters significant challenges in temporal and 125 social commonsense reasoning, indicating key ar-126 eas for further development. (4) Our manual error 127 analysis reveals that Gemini Pro often misunder-128 stands provided contextual information, accounting 129 for 30.2% of its total errors. Furthermore, Gemini 130 Pro Vision struggles particularly with identifying 131 emotional stimuli in images, especially those in-132 volving human entities, which constitutes 32.6% of 133 its total errors. 134

In summary, our contributions are threefold:

 We provide the first thorough evaluation of Gemini Pro's efficacy in commonsense reasoning tasks, employing 12 diverse datasets that span both language-based and multimodal scenarios.

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

158

159

160

161

162

163

164

166

167

168

169

170

171

172

173

174

175

176

177

178

179

180

181

182

- (2) Our study reveals that Gemini Pro exhibits performance comparable to GPT-3.5 Turbo in language-only commonsense reasoning tasks, demonstrating logical and contextual reasoning processes. However, it lags behind GPT-4 Turbo in accuracy and encounters challenges in temporal and social reasoning, as well as in emotion recognition in images.
- (3) Our findings lay a valuable foundation for future research in the field of commonsense reasoning within LLMs and MLLMs, highlighting the necessity to enhance specialized domains in these models and the nuanced recognition of mental states and emotions in multimodal contexts.

## 2 Experimental Setup

## 2.1 Datasets

We experiment with 12 datasets related to different types of commonsense reasoning, which include 11 language-based datasets and one multimodal dataset. The language-based datasets encompass three main categories of commonsense reasoning problems. General and Contextual Reasoning: (1) CommonsenseQA (Talmor et al., 2019), focusing on general commonsense knowledge; (2) Cosmos QA (Huang et al., 2019), emphasizing contextual understanding narratives, (3)  $\alpha$ NLI (Bhagavatula et al., 2019), introducing abductive reasoning, which involves inferring the most plausible explanation; and (4) HellaSWAG, centering around reasoning with contextual event sequences. Specialized and Knowledge Reasoning: (1) TRAM (Wang and Zhao, 2023b), testing reasoning about time; (2) NumerSense (Lin et al., 2020), focusing on numerical understanding; (3) PIQA (Bisk et al., 2020), assessing physical interaction knowledge; (4) QASC (Khot et al., 2020), dealing with science-related reasoning; and (5) RiddleSense (Lin et al., 2021), challenging creative thinking through riddles. Social and Ethical Reasoning: (1) Social IQa (Sap et al., 2019), testing the understanding of social interactions; and (2) ETHICS (Hendrycks et al., 2020), evaluating moral and ethical reasoning. For the multimodal dataset (vision and language), we select VCR (Zellers et al., 2019a), a large-scale dataset for cognition-level visual understanding. For datasets like TRAM and ETHICS, which include multiple tasks, we extract the commonsense reasoning part for experiments. We employ accuracy as the performance metric for all datasets. More details about each dataset, as well as example questions, are in Appendix B.

## 2.2 Models

183

184

185

188

189

190

191

192

193

194

196

197

198

204

209

210

211

212

213

214

215

216

217

218

219

222

225

231

We consider four popular LLMs for languagebased dataset evaluation, including the opensource model Llama-2-70b-chat (Touvron et al., 2023) as well as the closed-source models Gemini Pro (Team et al., 2023), GPT-3.5 Turbo, and GPT-4 Turbo (OpenAI, 2023). Each of these models is accessed using its corresponding API key. Specifically, we query Gemini through Google Vertex AI, the GPT models through the OpenAI API, and Llama2 through DeepInfra. For the multimodal dataset, we consider GPT-4V (gpt-4-vision-preview in API) and Gemini Pro Vision (gemini-pro-vision in API) in our experiments. Given the constraints of API costs and rate limitations, we randomly select 200 examples from the validation set for each language-based dataset following (Wang and Zhao, 2023b) and 50 examples from the validation set for the VCR dataset following (Liu and Chen, 2023). For all evaluations, we employ greedy decoding (i.e., temperature = 0) during model response generation. Notably, there are instances where the models decline to respond to certain queries, particularly those involving potentially illegal or unethical content. Sometimes, models provide answers that are outside the scope of the options. In these cases, we categorize these unanswered questions as incorrect.

## 2.3 Prompts

In the evaluation of language-based datasets, we employ two prompting settings: (1) zero-shot standard prompting (SP) (Kojima et al., 2022), which aims to gauge the models' inherent commonsense capabilities in linguistic contexts, and (2) few-shot chain-of-thought (CoT) prompting (Wei et al., 2022), implemented to observe potential enhancements in the models' performance. For the multimodal dataset, we utilize zero-shot standard prompting to assess the authentic end-to-end visual commonsense reasoning abilities of MLLMs.

## **3** Results

## 3.1 Overall Performance Comparison

233

235

236

237

238

239

240

241

242

243

244

245

246

247

248

249

250

251

252

253

254

255

256

257

259

260

261

262

263

264

265

267

268

269

270

271

272

273

274

275

276

277

278

279

280

281

Table 1 demonstrates the accuracy comparison of four LLMs under zero-shot SP and few-shot CoT settings on 11 language-based commonsense reasoning datasets. There are several key takeaways. First, from the model perspective, GPT-4 Turbo outperforms the other models across the majority of datasets on average. Under the zero-shot learning paradigm, it surpasses Gemini Pro, the second-best performing model, by 7.3%, and shows an even greater lead of 9.0% under the few-shot learning paradigm. Gemini Pro exhibits marginally higher average accuracy than GPT-3.5 Turbo, with an increase of 1.3% under zero-shot SP and 1.5% in the few-shot CoT scenario. It also demonstrates substantially better performance than Llama-2-70b. Regarding prompting methods, the CoT approach consistently enhances performance across all datasets, with pronounced gains observed in datasets such as CommonsenseQA, TRAM, and Social IQa. Lastly, from a dataset standpoint, it is apparent that while these models exhibit commendable performance across a broad spectrum of commonsense domains, they encounter challenges in specific areas, particularly those involving temporal (TRAM) and social (Social IQa) dimensions of commonsense reasoning.

For the multimodal VCR dataset, we report the performance of GPT-4V and Gemini Pro Vision in Table 2. The VCR consists of three subtasks: (1)  $O \rightarrow A$ , which involves generating an answer to a question based on the visual context; (2)  $QA \rightarrow R$ , which requires the model to produce a rationale for a given answer; and (3)  $Q \rightarrow AR$ , which challenges the model to both answer the question and justify the response with appropriate rationales. In all subtasks, GPT-4V demonstrates superior performance compared to Gemini Pro Vision, indicating a more robust capacity for integrating visual and textual information to provide coherent responses. In  $Q \rightarrow$ AR, the relatively lower performance of both models, compared to the other two subtasks, suggests that there is considerable room for improvement in understanding the interplay between visual cues and commonsense reasoning.

## 3.2 Effects of Commonsense Domain

Referring to Section 2.1, we have categorized 11 language-based datasets into three groups and presented the performance for each setting within each



Figure 1: Average model performance across three major commonsense reasoning categories over 11 languagebased datasets, including General and Contextual Reasoning (CommonsenseQA, Cosmos QA,  $\alpha$ NLI, HellaSWAG), Specialized and Knowledge Reasoning (TRAM, NumerSense, PIQA, QASC, RiddleSense), and Social and Ethical Reasoning (Social IQa, ETHICS). GPT-4 Turbo consistently exhibits superior performance in all commonsense reasoning categories. Gemini Pro marginally surpasses GPT-3.5 Turbo in the first two categories, except for Social and Ethical Reasoning.

group in Figure 1. Our findings indicate that GPT-4 284 Turbo consistently leads in performance across all categories. The Llama-2-70b model demonstrates marginally lower accuracy in comparison to the other models. Gemini Pro and GPT-3.5 Turbo display comparable performances; however, Gemini Pro slightly outperforms GPT-3.5 Turbo in two of 289 the three categories. Notably, its performance dip in the Social and Ethical Reasoning group may stem from its tendency to refuse to answer questions that could potentially involve unethical con-293 tent, which we have counted as incorrect in our evaluation. Based on our experiments, among the 200 samples, Gemini Pro refuses to answer 3.0% of the problems (6 in total) in the Social IQa dataset and 6.5% of the problems (13 in total) in the ETHICS dataset. Overall, all models exhibit robust capabilities in handling Social and Ethical Reasoning datasets, suggesting a relatively advanced grasp of moral and social norms. However, there is a notable disparity in their performance on General and Contextual Reasoning tasks, indicating a potential gap in their understanding of broader commonsense principles and their application in varied contexts. The Specialized and Knowledge Reasoning category, particularly in the realms of temporal and

riddle-based challenges, highlights specific deficiencies in the models' abilities to process complex temporal sequences and to engage in the abstract and creative thought required to decipher riddles.

309

310

311

312

313

314

315

316

317

318

319

320

321

322

323

324

325

328

330

331

332

333

334

Regarding the multimodal dataset, Figure 2 details the comparative performance between GPT-4V and Gemini Pro Vision across different question types, in alignment with the guidelines of the VCR dataset (Zellers et al., 2019a). We concentrate on the "Q  $\rightarrow$  A" subtask as it most directly assesses the models' visual commonsense capabilities. Considering the data sample for each type, Gemini Pro Vision's performance either matches or is slightly lower than GPT-4V's, except in temporaltype questions, where it surpasses GPT-4V. This suggests its enhanced capability not only in recognizing but also in contextualizing time-related elements within visual scenarios.

#### 3.3 **Reasoning Justification within MLLMs**

To assess the reasoning capabilities of MLLMs, particularly their ability to provide not only correct answers but also sound and contextually grounded reasoning in matters of commonsense, we adopted a systematic sampling approach. For each of the 11 language-based datasets evaluated with four LLMs, we randomly selected 30 questions that

Table 1: Performance comparison of four LLMs across 11 language-based commonsense reasoning datasets. For the k-shot CoT setting, k is set to 5 for the majority of datasets, except HellaSWAG (k=10) and PIQA (k=1). The best results for the k-shot setting are boldfaced, and for the 0-shot setting, underlined. GPT-4 Turbo outperforms other models across the majority of datasets under both settings by a large margin. Gemini Pro and GPT-3.5 Turbo exhibit comparably matched performance overall, with Gemini Pro demonstrating marginally superior commonsense reasoning capabilities compared to GPT-3.5 Turbo on average.

Dataset	Method							
	Llama-2-70b (0-shot, SP)	Llama-2-70b (k-shot, CoT)	Gemini Pro (0-shot, SP)	Gemini Pro (k-shot, CoT)	GPT-3.5 Turbo (0-shot, SP)	GPT-3.5 Turbo (k-shot, CoT)	GPT-4 Turbo (0-shot, SP)	GPT-4 Turbo (k-shot, CoT)
CommonsenseQA	72.0	76.5	76.5	79.0	73.0	76.0	<u>78.0</u>	80.0
Cosmos QA	77.0	81.0	81.5	84.5	75.0	78.5	<u>86.5</u>	88.0
$\alpha$ NLI	77.5	80.5	79.5	81.5	75.5	78.0	<u>87.0</u>	88.0
HellaSWAG	73.0	77.0	76.0	78.5	78.0	80.0	<u>94.0</u>	95.0
TRAM	66.0	70.0	73.5	76.0	68.5	72.0	<u>79.5</u>	82.0
NumerSense	74.0	75.5	80.0	82.0	81.5	82.5	<u>85.0</u>	86.0
PIQA	74.0	78.5	89.0	90.5	87.0	89.5	<u>94.5</u>	95.5
QASC	78.0	82.0	80.0	82.5	83.0	85.0	<u>91.5</u>	92.5
RiddleSense	62.5	66.0	75.0	82.5	71.5	75.0	<u>94.0</u>	95.0
Social IQa	71.0	77.5	73.0	78.5	73.0	78.0	<u>82.0</u>	84.5
ETHICS	88.0	89.5	87.0	87.5	94.0	95.0	<u>97.0</u>	98.0
Average	73.9	77.6	79.2	82.1	78.2	80.9	<u>88.1</u>	89.5



Figure 2: Performance comparison between GPT-4V and Gemini Pro Vision on the VCR dataset, categorized by question type, with a focus on the "Q  $\rightarrow$  A" subtask. Within our sample of 50 questions, the distribution across each type is as follows: activity (12), explanation (16), hypothetical (3), mental (4), role (5), scene (4), and temporal (6). GPT-4V matches or surpasses Gemini Pro Vision in performance across these question types, with the exception of the temporal category.

were correctly answered and 30 questions that were incorrectly answered by each LLM following (Bian et al., 2023). In cases where a dataset presented fewer than 30 incorrect answers, we included all available incorrect responses to ensure comprehensive analysis. After selecting these questions, we prompted each model to explain "*What is the rationale behind the answer to the question?*" The reasoning processes provided by the models were then manually reviewed and classified as either True or False, based on their logical soundness and

335

336

337

339

341

343

Table 2: Performance comparison between GPT-4V and Gemini Pro Vision on the VCR dataset. "Q  $\rightarrow$  A" evaluates question-answering accuracy, "QA  $\rightarrow$  R" assesses answer justification, and "Q  $\rightarrow$  AR" measures the performance of both correctly answering questions and selecting rationales. GPT-4V outperforms Gemini Pro Vision across all subtasks.

Method	$\mathbf{Q} \to \mathbf{A}$	$\mathbf{Q}\mathbf{A}\to\mathbf{R}$	$\boldsymbol{Q} \to \boldsymbol{A}\boldsymbol{R}$
GPT-4V	80.0	72.0	56.0
Gemini Pro Vision	74.0	70.0	48.0

relevance to the question. Figure 3 illustrates a comprehensive view of the average reasoning correctness across the 11 datasets, in terms of the sampled correct and incorrect questions. In fact, not every model had 30 incorrect questions for each dataset. In such scenarios, we scaled the available data up to 30 questions to ensure standardized computation. Figure 3 shows that GPT-4 Turbo's leading performance in both correct and incorrect answers highlights its advanced reasoning mechanisms and its ability to maintain coherent logic, even when the final answers are not accurate. Additionally, Gemini Pro has emerged as a notably proficient model, generally demonstrating commendable reasoning abilities and offering a well-rounded approach to commonsense reasoning. GPT-3.5, while trailing slightly behind Gemini Pro, still demonstrates competitive reasoning abilities. Appendix C presents two real examples from Gemini Pro and GPT-3.5,

364

346

365

369

372

374

385

391

illustrating the cases of a correct answer with a correct rationale and an incorrect answer with an incorrect rationale, respectively.



Figure 3: Average reasoning correctness across 11 language datasets. The comparison among four LLMs is based on a random sample of 30 correct and 30 incorrect questions per dataset. In cases where a dataset contained fewer than 30 incorrect questions, the data were scaled up to maintain consistency in the sample size.

Moving to the multimodal perspective, our analysis of GPT-4V and Gemini Pro Vision on the VCR dataset reveals notable patterns in reasoning correctness. With GPT-4V at 24% and Gemini Pro Vision at 26%, approximately one-quarter of the cases showed both models correctly identifying the answers but failing to provide appropriate rationale. This discrepancy suggests that while the models can often determine the correct outcomes, their ability to understand or explain the underlying reasoning behind these answers is not consistently aligned. Furthermore, in the instances of incorrect answers, GPT-4V and Gemini Pro Vision showed correct rationales 16% and 22% of the time, respectively. This indicates that, despite arriving at incorrect conclusions, the models demonstrate a capacity for effective reasoning or logical processing. However, this does not consistently translate into accurate outcomes, implying that while some aspects of the required knowledge are captured, other crucial elements are likely missed.

## 3.4 Case Study: Gemini Pro in Commonsense Reasoning

Given our focus on evaluating the commonsense reasoning capabilities of the Gemini Pro model, we conduct a qualitative analysis to assess its performance across representative examples in four major categories (three language-based and one multimodal), as described in Section 2.1. To ensure an authentic end-to-end capability evaluation, we present examples under the zero-shot learning setting, employing standard prompting techniques. Due to space constraints, we present two examples here; additional examples are in Appendix D. 395

396

397

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

General (CommonsenseQA). In the general commonsense evaluation (General and Contextual Reasoning category) using the CommonsenseQA dataset, consider the example question: "People are what when you're a stranger? (A) train (B) strange (C) human (D) stupid (E) dangerous." Gemini Pro correctly chose (B) "strange," and its reasoning process is notable. It recognized that while all options relate to the concept of a "stranger", only "strange" accurately encapsulates the neutral and open-ended nature of the question. The model effectively ruled out other options: (A) "train", for being too specific and unrelated; (C) "human", as accurate but not capturing the question's essence; (D) "stupid", for being judgmental and offensive; and (E) "dangerous", due to its negative connotation. This selection of "strange" indicates an understanding of the unfamiliar nature associated with strangers, highlighting Gemini Pro's capability in interpreting and applying general commonsense knowledge appropriately.

Visual (VCR). In the visual commonsense evaluation using the VCR dataset, we analyzed Gemini Pro Vision's response to a scenario involving physical safety and potential danger, as shown in Figure 4. Presented with an image of individuals on the edge of a cliff, the model was questioned: "What would happen if person 4 pushed person 3 at this moment?" In this context, Gemini Pro Vision's response mirrored the logical inference that if the second person from the left (person 4) pushed the third person from the left (person 3), the result would be person 3 falling off the cliff, leading to a fatal outcome. This example from the VCR dataset underscores Gemini Pro Vision's ability to analyze visual scenes and make predictions about the potential consequences of actions within those scenes, a crucial aspect of visual commonsense reasoning. It demonstrates the model's grasp of spatial relations and physical consequences, providing evidence of its capacity to process and reason about complex visual information akin to human cognition.

Overall, the cases presented underscore the advanced reasoning capabilities of Gemini Pro and



Figure 4: Example image from the VCR dataset.

Gemini Pro Vision, while also identifying challenges in achieving human-like inference.

### 3.5 Error Analysis

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

To gain a deeper understanding of the mistakes made by models, we manually analyzed instances where a model made incorrect choices or provided inappropriate answers. We conducted a thorough examination of common error types encountered in commonsense reasoning tasks, with the results averaged across four LLMs. Our focus was on assessing these models in two distinct settings: zero-shot SP and few-shot CoT. Table 3 shows the proportions of five common error types averaged over four LLMs in each setting.

Context misinterpretation emerged as the most frequent error, occurring more often in the zeroshot SP setting (28.6%) compared to the few-shot CoT (23.4%). This trend suggests that the additional context in few-shot CoT helps models better understand scenarios, thereby reducing errors related to contextual misunderstanding. Logical errors were the second most common, accounting for 23.9% in zero-shot SP and slightly less in fewshot CoT (20.1%), indicating that extra examples in the latter setting aid in more consistent logical reasoning. Ambiguity errors, at 16.2% in zeroshot SP, were reduced to 11.6% in few-shot CoT, demonstrating the effectiveness of added context in resolving language ambiguities. In contrast, Overgeneralization errors showed an increase in fewshot CoT (15.6%) from zero-shot SP (11.8%), possibly due to models' overextending patterns learned from the additional examples. Notably, knowledge errors, where models misapplied correct and necessary commonsense knowledge, saw a significant increase in few-shot CoT (29.3%) compared to zero-shot SP (19.5%). This finding suggests that while extra context can be beneficial, it can also lead to inaccuracies, particularly in complex or nuanced scenarios.

Table 3: Proportion of common error types in commonsense reasoning in LLM evaluation. Misinterpret. represents misinterpretation.

Error Type	Zero-shot SP	Few-shot CoT
Context Misinterpret.	28.6%	23.4%
Logical Errors	23.9%	20.1%
Text Ambiguity	16.2%	11.6%
Overgeneralization	11.8%	15.6%
Knowledge Errors	19.5%	29.3%

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

503

504

505

506

507

508

509

510

511

512

513

514

515

In our analysis of the VCR dataset, we focused on instances where either GPT-4V or Gemini Pro Vision chose incorrect answers in the  $Q \rightarrow A$  subtask. The four common error types for each model are summarized in Table 4. Emotion recognition errors were the most common, with GPT-4V encountering these errors in 30.1% of cases and Gemini Pro Vision slightly more at 31.3%. This high incidence suggests that both models find interpreting emotional cues in visual content particularly challenging, underscoring the complexity of deciphering human emotions from visual stimuli. Spatial perception errors were also significant, constituting 22.5% of errors for GPT-4V and 25.2% for Gemini Pro Vision. These figures indicate the models' difficulties in accurately understanding spatial relationships and the arrangement of elements in images. Logical errors were another major error type, more pronounced in GPT-4V (27.7%) than in Gemini Pro Vision (24.9%), pointing to challenges in logical reasoning within visual contexts. Context misinterpretation, although less frequent, was still a notable issue, with GPT-4V at 19.7% and Gemini Pro Vision at 18.6%. These errors demonstrate the models' struggles with grasping the overarching context or narrative depicted in the visual content.

Overall, error analysis sheds light on the specific challenges LLMs and MLLMs face in commonsense reasoning, providing valuable insights for future improvements for future model refinement.

Table 4: Proportion of commmon error types in visual commonsense reasoning in MLLM evaluation (GPT-4V and Gemini Pro Vision). Misinterpret.: and E. represent misinterpretation and errors, respectively.

Error Type	GPT-4V	Gemini Pro Vision
Context Misinterpret.	19.7%	18.6%
Spatial Perception E.	22.5%	25.2%
Emotion Recognition E.	30.1%	31.3%
Logical Errors	27.7%	24.9%

#### **Related Work** 4

516

517

518

519

520

521

526

527

531

532

533

535

536

539

540

551

554

557

561

Commonsense Reasoning in NLP. Commonsense reasoning has gained renewed attention in recent years, especially in the context of advancements in LLMs that have significantly influenced numerous applications in NLP. However, there is a growing concern about their ability to understand and reason about commonsense knowledge (Storks et al., 2019; Tamborrino et al., 2020; Bhargava and Ng, 2022). This concern is echoed in various studies that focus on evaluating the capabilities of LLMs in this area (Bian et al., 2023; Weng et al., 2023; Shen and Kejriwal, 2023). Concurrently, researchers have been exploring diverse strategies to enhance the commonsense reasoning of NLP systems, from leveraging knowledge graphs to commonsense knowledge transfer (Huang et al., 2023; Ye et al., 2023; Zhou et al., 2023). Prior to delving into methodological refinements, a comprehensive evaluation is essential to understand the authentic commonsense reasoning capabilities of LLMs. In our study, we endeavor to advance this line of inquiry by examining how LLMs, particularly focusing on the Gemini model, navigate and implement commonsense reasoning in various NLP contexts.

Training Paradigms in LLMs. In NLP research, 541 pretraining language models on large-scale varied textual datasets has become essential. BERT-based 543 models like BERT (Kenton and Toutanova, 2019) 544 and RoBERTa (Liu et al., 2019) exemplify this, 545 being applied to tasks ranging from disease prediction (Zhao et al., 2021) to text classification (Wang et al., 2022b) and time series analysis (Wang et al., 2022c). The debut of GPT-3 shifted this focus 549 towards more flexible learning methods like zero-550 shot and few-shot learning, showcasing models' adaptability to new tasks with minimal data (Brown et al., 2020). This shift has spurred the development of novel prompting techniques to enhance LLMs' reasoning and understanding capabilities, 555 including CoT prompting (Wei et al., 2022), self-556 consistency with CoT (Wang et al., 2022a), tree-ofthought prompting (Yao et al., 2023), and metacognitive prompting (Wang and Zhao, 2023a). In this work, we evaluate four LLMs for language tasks and two MLLMs for multimodal tasks under zeroshot and few-shot settings to provide an in-depth 562 understanding of their strengths and limitations in diverse commonsense reasoning tasks.

Evaluations on MLLMs. Since the release of

the state-of-the-art MLLM, GPT-4V, several evaluations have been conducted across diverse tasks, including medical imaging (Wu et al., 2023a), visual question answering (Li et al., 2023; Yang et al., 2023), and video understanding (Lin et al., 2023), focusing on either on case-by-case qualitative analyses or on quantitative assessments across diverse tasks. The recent release of Google's Gemini has garnered considerable attention, and early experiments have been conducted to evaluate its capabilities in both language understanding (Akter et al., 2023) and the multimodal domain (Liu and Chen, 2023; Fu et al., 2023). However, a significant gap remains in fully comprehending the commonsense reasoning capabilities of Gemini, a known potential shortcoming since its introduction. Our work comprehensively analyzes Gemini's capabilities in this area, comparing it with other MLLMs to highlight its potential and areas for improvement.

566

567

568

569

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586

587

588

589

590

591

592

593

594

595

596

597

598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

#### 5 Discussion

In this study, we conducted a comprehensive evaluation of state-of-the-art LLMs and MLLMs, focusing particularly on Gemini Pro and Gemini Pro Vision, across 12 diverse commonsense reasoning datasets. Our findings indicate that while these models mark a significant advancement in various domains, demonstrating impressive performance in commonsense reasoning tasks, they still exhibit limitations, particularly in tasks requiring deep contextual understanding or abstract reasoning, such as those involving temporal dynamics, riddles, or intricate social scenarios. Although significant progress has been made, achieving AGI still represents a substantial goal on the horizon. Our work sets the stage for future research in this field, highlighting both the achievements and areas needing improvement in commonsense reasoning.

Looking ahead, addressing these challenges is crucial to enhance the overall effectiveness of LLMs and MLLMs in commonsense reasoning. Future research should aim to refine the models' capabilities in interpreting and reasoning within complex contexts and abstract scenarios. Additionally, there is an emerging need for more holistic evaluation metrics and methodologies capable of accurately assessing the nuances of commonsense reasoning in AI systems. These metrics should evaluate not only the correctness of responses but also their logical coherence and context relevance.

## 6 Limitations

615

637

638

641

648

653

654

656

657

661

663

664

While this study offers valuable insights into the 616 role of LLMs and MLLMs in commonsense reason-617 ing, there are some limitations. Firstly, our evalua-618 tion is heavily dependent on the selected questions and datasets used for analysis. Despite their diver-621 sity, these datasets may not cover all facets of this domain. As a result, the performance and capabilities of Gemini Pro and other models can vary in real-world scenarios or with alternative datasets. Additionally, our analysis is confined to English 625 language datasets, limiting the generalizability of our findings to multilingual contexts, where cultural nuances and linguistic differences are crucial in commonsense reasoning. Finally, our study represents a specific moment in the rapidly evolving 630 landscape of AI, focusing on API-based systems that are subject to change. The introduction of newer models or updates to existing ones might lead to different performance outcomes, highlighting the need for ongoing evaluation and analysis. 635

## References

- Syeda Nahida Akter, Zichun Yu, Aashiq Muhamed, Tianyue Ou, Alex Bäuerle, Ángel Alexander Cabrera, Krish Dholakia, Chenyan Xiong, and Graham Neubig. 2023. An in-depth look at gemini's language abilities. *arXiv preprint arXiv:2312.11444*.
  - Chandra Bhagavatula, Ronan Le Bras, Chaitanya Malaviya, Keisuke Sakaguchi, Ari Holtzman, Hannah Rashkin, Doug Downey, Wen-tau Yih, and Yejin Choi. 2019. Abductive commonsense reasoning. In International Conference on Learning Representations.
  - Prajjwal Bhargava and Vincent Ng. 2022. Commonsense knowledge reasoning and generation with pretrained language models: A survey. In *Proceedings* of the AAAI Conference on Artificial Intelligence, volume 36, pages 12317–12325.
  - Ning Bian, Xianpei Han, Le Sun, Hongyu Lin, Yaojie Lu, and Ben He. 2023. Chatgpt is a knowledgeable but inexperienced solver: An investigation of commonsense problem in large language models. *arXiv preprint arXiv:2303.16421*.
- Yonatan Bisk, Rowan Zellers, Jianfeng Gao, Yejin Choi, et al. 2020. Piqa: Reasoning about physical commonsense in natural language. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pages 7432–7439.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda

Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901. 666

667

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

700

703

704

705

706

707

708

709

710

711

712

713

714

715

718

- Chaoyou Fu, Renrui Zhang, Haojia Lin, Zihan Wang, Timin Gao, Yongdong Luo, Yubo Huang, Zhengye Zhang, Longtian Qiu, Gaoxiang Ye, et al. 2023. A challenger to gpt-4v? early explorations of gemini in visual expertise. *arXiv preprint arXiv:2312.12436*.
- Zhankui He, Zhouhang Xie, Rahul Jha, Harald Steck, Dawen Liang, Yesu Feng, Bodhisattwa Prasad Majumder, Nathan Kallus, and Julian McAuley. 2023. Large language models as zero-shot conversational recommenders. In *Proceedings of the 32nd ACM international conference on information and knowledge management*, pages 720–730.
- Dan Hendrycks, Collin Burns, Steven Basart, Andrew Critch, Jerry Li, Dawn Song, and Jacob Steinhardt. 2020. Aligning ai with shared human values. In *International Conference on Learning Representations*.
- Lifu Huang, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. 2019. Cosmos qa: Machine reading comprehension with contextual commonsense reasoning. In Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), pages 2391–2401.
- Yongfeng Huang, Yanyang Li, Yichong Xu, Lin Zhang, Ruyi Gan, Jiaxing Zhang, and Liwei Wang. 2023. Mvp-tuning: Multi-view knowledge retrieval with prompt tuning for commonsense reasoning. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 13417–13432.
- Enkelejda Kasneci, Kathrin Seßler, Stefan Küchemann, Maria Bannert, Daryna Dementieva, Frank Fischer, Urs Gasser, Georg Groh, Stephan Günnemann, Eyke Hüllermeier, et al. 2023. Chatgpt for good? on opportunities and challenges of large language models for education. *Learning and individual differences*, 103:102274.
- Jacob Devlin Ming-Wei Chang Kenton and Lee Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of naacL-HLT*, volume 1, page 2.
- Tushar Khot, Peter Clark, Michal Guerquin, Peter Jansen, and Ashish Sabharwal. 2020. Qasc: A dataset for question answering via sentence composition. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 8082–8090.
- Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. Large language models are zero-shot reasoners. *Advances in neural information processing systems*, 35:22199– 22213.

814

815

816

817

818

819

820

821

822

823

824

825

826

827

828

774

775

Xiang Lorraine Li, Adhiguna Kuncoro, Cyprien de Masson d'Autume, Phil Blunsom, and Aida Nematzadeh.
2021. Do language models learn commonsense knowledge? arXiv preprint arXiv:2111.00607.

720

721

724

725

726

727

729 730

731

733

735

736

737

741

742 743

745

746

747

748

751

752

754

755

757

758

765

769

770

- Yunxin Li, Longyue Wang, Baotian Hu, Xinyu Chen, Wanqi Zhong, Chenyang Lyu, and Min Zhang. 2023.
  A comprehensive evaluation of gpt-4v on knowledgeintensive visual question answering. *arXiv preprint arXiv:2311.07536*.
- Bill Yuchen Lin, Seyeon Lee, Rahul Khanna, and Xiang Ren. 2020. Birds have four legs?! numersense: Probing numerical commonsense knowledge of pretrained language models. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 6862–6868.
- Bill Yuchen Lin, Ziyi Wu, Yichi Yang, Dong-Ho Lee, and Xiang Ren. 2021. Riddlesense: Reasoning about riddle questions featuring linguistic creativity and commonsense knowledge. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP* 2021, pages 1504–1515.
- Kevin Lin, Faisal Ahmed, Linjie Li, Chung-Ching Lin, Ehsan Azarnasab, Zhengyuan Yang, Jianfeng Wang, Lin Liang, Zicheng Liu, Yumao Lu, et al. 2023. Mmvid: Advancing video understanding with gpt-4v (ision). *arXiv preprint arXiv:2310.19773*.
- Hugo Liu and Push Singh. 2004. Conceptnet—a practical commonsense reasoning tool-kit. *BT technology journal*, 22(4):211–226.
- Mengchen Liu and Chongyan Chen. 2023. An evaluation of gpt-4v and gemini in online vqa. *arXiv preprint arXiv:2312.10637*.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019.
  Roberta: A robustly optimized bert pretraining approach. arXiv preprint arXiv:1907.11692.
- Bonan Min, Hayley Ross, Elior Sulem, Amir Pouran Ben Veyseh, Thien Huu Nguyen, Oscar Sainz, Eneko Agirre, Ilana Heintz, and Dan Roth. 2021. Recent advances in natural language processing via large pre-trained language models: A survey. ACM Computing Surveys.
- OpenAI. 2023. Gpt-4 technical report.
- Abu Rayhan, Rajan Rayhan, and Swajan Rayhan. 2023. Artificial general intelligence: Roadmap to achieving human-level capabilities.
- Maarten Sap, Hannah Rashkin, Derek Chen, Ronan Le Bras, and Yejin Choi. 2019. Social iqa: Commonsense reasoning about social interactions. In Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), pages 4463–4473.

- Ke Shen and Mayank Kejriwal. 2023. An experimental study measuring the generalization of finetuned language representation models across commonsense reasoning benchmarks. *Expert Systems*, page e13243.
- Vered Shwartz and Yejin Choi. 2020. Do neural language models overcome reporting bias? In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 6863–6870.
- Shane Storks, Qiaozi Gao, and Joyce Y Chai. 2019. Commonsense reasoning for natural language understanding: A survey of benchmarks, resources, and approaches. *arXiv preprint arXiv:1904.01172*, pages 1–60.
- Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. 2019. Commonsenseqa: A question answering challenge targeting commonsense knowledge. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), pages 4149–4158.
- Alexandre Tamborrino, Nicola Pellicanò, Baptiste Pannier, Pascal Voitot, and Louise Naudin. 2020. Pretraining is (almost) all you need: An application to commonsense reasoning. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 3878–3887.
- Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. 2023. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.
- Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2022a. Self-consistency improves chain of thought reasoning in language models. *arXiv preprint arXiv:2203.11171*.
- Yuqing Wang and Yun Zhao. 2023a. Metacognitive prompting improves understanding in large language models. *arXiv preprint arXiv:2308.05342*.
- Yuqing Wang and Yun Zhao. 2023b. Tram: Benchmarking temporal reasoning for large language models. *arXiv preprint arXiv:2310.00835*.
- Yuqing Wang, Yun Zhao, Rachael Callcut, and Linda Petzold. 2022b. Integrating physiological time series and clinical notes with transformer for early prediction of sepsis. *arXiv preprint arXiv:2203.14469*.

884

885

886

887

Yuqing Wang, Yun Zhao, and Linda Petzold. 2022c. Enhancing transformer efficiency for multivariate time series classification. *arXiv preprint arXiv:2203.14472*.

829

830

832

834

835

837

838

839

842

843

844

849

851

853

854

858

859

863

868

870

871

872

873

874

875

876

877

878

879

881

- Yuqing Wang, Yun Zhao, and Linda Petzold. 2023. Are large language models ready for healthcare? a comparative study on clinical language understanding. *arXiv preprint arXiv:2304.05368*.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in Neural Information Processing Systems*, 35:24824–24837.
- Yixuan Weng, Minjun Zhu, Fei Xia, Bin Li, Shizhu He, Shengping Liu, Bin Sun, Kang Liu, and Jun Zhao. 2023. Large language models are better reasoners with self-verification. In *Findings of the Association* for Computational Linguistics: EMNLP 2023, pages 2550–2575.
- Chaoyi Wu, Jiayu Lei, Qiaoyu Zheng, Weike Zhao, Weixiong Lin, Xiaoman Zhang, Xiao Zhou, Ziheng Zhao, Ya Zhang, Yanfeng Wang, et al. 2023a. Can gpt-4v (ision) serve medical applications? case studies on gpt-4v for multimodal medical diagnosis. *arXiv preprint arXiv:2310.09909*.
- Jiayang Wu, Wensheng Gan, Zefeng Chen, Shicheng Wan, and Philip S Yu. 2023b. Multimodal large language models: A survey. *arXiv preprint arXiv:2311.13165*.
- Zhengyuan Yang, Linjie Li, Kevin Lin, Jianfeng Wang, Chung-Ching Lin, Zicheng Liu, and Lijuan Wang. 2023. The dawn of lmms: Preliminary explorations with gpt-4v (ision). *arXiv preprint arXiv:2309.17421*, 9(1).
- Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L Griffiths, Yuan Cao, and Karthik Narasimhan. 2023. Tree of thoughts: Deliberate problem solving with large language models. *arXiv preprint arXiv:2305.10601*.
- Shuquan Ye, Yujia Xie, Dongdong Chen, Yichong Xu, Lu Yuan, Chenguang Zhu, and Jing Liao. 2023. Improving commonsense in vision-language models via knowledge graph riddles. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2634–2645.
- Rowan Zellers, Yonatan Bisk, Ali Farhadi, and Yejin Choi. 2019a. From recognition to cognition: Visual commonsense reasoning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 6720–6731.
- Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. 2019b. Hellaswag: Can a machine really finish your sentence? In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4791–4800.

- Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, et al. 2023. A survey of large language models. *arXiv preprint arXiv:2303.18223*.
- Yun Zhao, Yuqing Wang, Junfeng Liu, Haotian Xia, Zhenni Xu, Qinghang Hong, Zhiyang Zhou, and Linda Petzold. 2021. Empirical quantitative analysis of covid-19 forecasting models. In 2021 International Conference on Data Mining Workshops (ICDMW), pages 517–526. IEEE.
- Wangchunshu Zhou, Ronan Le Bras, and Yejin Choi. 2023. Commonsense knowledge transfer for pre-trained language models. *arXiv preprint arXiv:2306.02388*.

## A Commonsense Overview

Commonsense reasoning, a fundamental aspect of human intelligence, facilitates an intuitive understanding and interpretation of the world through basic and often implicit knowledge and beliefs. For instance, it involves understanding that a person carrying an umbrella on a cloudy day likely anticipates rain, or inferring that a closed door in a library signifies a need for quiet. In MLLMs, commonsense reasoning plays a vital role, enabling these models to interact with and interpret human language and visual cues in a manner that mirrors human understanding. In our study, we explore a variety of commonsense reasoning tasks. Definitions for each domain are provided as follows.

General Commonsense. This domain entails an understanding of basic, everyday knowledge about the world, such as recognizing that birds typically fly and fish live in water.

**Contextual Commonsense.** This domain involves interpreting information within specific contexts, such as understanding that a person wearing a coat and shivering is likely cold.

**Abductive Commonsense.** This domain is about formulating the most plausible explanations for observations, such as inferring that wet streets are likely due to recent rain.

**Event Commonsense.** This domain focuses on understanding sequences of events and the causal relationships between them, such as predicting that eating spoiled food can lead to feeling sick.

**Temporal Commonsense.** This domain involves understanding time-related concepts, such as the fact that breakfast is typically eaten in the morning. **Numerical Commonsense.** This domain is about understanding numbers in everyday contexts, such as knowing that a cube has six faces.

- **Physical Commonsense.** This domain concerns 936 understanding the physical world, such as knowing 937 that a glass will break if dropped on a hard floor. 938
- Science Commonsense. This domain involves the application of scientific principles in daily life, such as understanding that water boils at a higher tem-941 perature at sea level than in the mountains. 942
- Riddle Commonsense. This domain challenges creative thinking through riddles, such as deciphering a riddle where the answer is "a shadow", requiring lateral thinking to associate intangible concepts with physical entities. 947
- Social Commonsense. This domain involves understanding social interactions, such as recognizing 949 that a person is likely upset if he/she is crying.
- Moral Commonsense. This domain deals with 951 952 evaluating actions based on moral and ethical standards, such as understanding that stealing is generally considered wrong. 954

Visual Commonsense. This domain involves interpreting and understanding visual information in the context of the physical and social world, such as deducing that a person in a photo is likely running a race if they are wearing a number bib and surrounded by other runners.

### R **Datasets**

955

957

959

960

961

962

964

965

966

967

970 971

972

981

We provide an overview of each dataset as well as example questions in Table 5.

### С **Reasoning Justification Examples**

Figure 5 illustrates examples of correct and incorrect answers with their corresponding rationales for Gemini Pro and GPT-3.5.

### More Examples of Gemini Pro in D **Commonsense Reasoning**

We provide two more case examples for Gemini Pro to get a deeper understanding of its commonsense reasoning capabilities.

Temporal (TRAM). In the temporal commonsense 973 evaluation (Specialized and Knowledge Reasoning 974 category) using the TRAM dataset, consider the 975 example question: "He also promises to 'come to' 976 him. How long does it take for him to 'come to' 977 him? (A) 100 years (B) in a minute's time (C) a few hours." Lacking sufficient context, especially 979 regarding the identities involved and the meaning 980 of 'come to', Gemini Pro was unable to provide a definitive answer. Gemini Pro's response highlights a significant aspect of its temporal reasoning capabilities. It illustrates the model's reliance on specific contextual information to make accurate temporal judgments. While this cautious approach is prudent to avoid incorrect assumptions, it also signifies a limitation in addressing ambiguous or incomplete information – a frequent challenge in real-world communications. This example underlines the difficulties LLMs encounter in temporal reasoning tasks, especially when faced with scenarios that offer limited or unclear context.

984

985

986

987

988

989

990

991

992

993

Social (Social IQa). In assessing Gemini Pro's 994 performance in social commonsense reasoning us-995 ing the Social IQa dataset (Social and Ethical Rea-996 soning category), an interesting scenario was pre-997 sented: "The people bullied Sasha all her life. But 998 Sasha got revenge on the people. What will the peo-999 ple want to do next? (A) Do whatever Sasha says 1000 (B) Get even (C) Flee from Sasha." The correct 1001 answer is (C), but Gemini Pro's response is insight-1002 ful. It chose (B) "Get even" as the most likely 1003 option, reasoning that the desire for revenge is a 1004 strong motivator and Sasha's actions likely ignited 1005 a similar desire in them. Gemini Pro considered 1006 (A) as a less likely option, depending on whether 1007 Sasha's revenge instilled deep fear and assumed complete submission. The least likely option, ac-1009 cording to Gemini Pro, was (C), which it associ-1010 ated with physical harm or an ongoing threat. This 1011 response demonstrates Gemini Pro's nuanced un-1012 derstanding of social dynamics and emotional mo-1013 tivations. However, it also highlights a limitation 1014 in accurately predicting human reactions in com-1015 plex social scenarios, where emotional responses 1016 might not always follow a logical pattern. This 1017 instance reveals the challenges LLMs face in ac-1018 curately interpreting and responding to intricate 1019 social situations, an area that remains challenging 1020 for AI systems. 1021 Table 5: Overview of commonsense datasets used in our experiments. "K-Way MC" signifies a multiple-choice response format with K options. Bold text in the "Example Questions" column represents the correct answers.

Dataset	Domain	Answer Type	Example Questions
			General and Contextual Reasoning
CommonsenseQA	general	5-Way MC	Where is a doormat likely to be in front of? (A). facade; ( <b>B</b> ). front door; (C). doorway; (D). entrance porch; (E). hallway.
Cosmos QA	contextual	4-Way MC	Given the context "It wasn't time for my book to be released I have received about five rejection letters." What may be the reason for your book getting rejected? (A). None of the above choices; (B). I never; (C). I felt; (D). It wasn't finished.
αNLI	abductive	2-Way MC	Given the beginning of the story: Four Outlaws camped in Blood Gulch, and the end of the story: He arrested them, what is the more plausible hypothesis: (A). They found where the sheriff was; ( <b>B</b> ). <b>The sheriff found where they were.</b>
HellaSWAG	event	4-Way MC	Given the context "A boy in an orange shirt is playing a video game. the scene" and the activity label "Washing face", which of the following endings is the most appropriate continuation of the scenario? (A). changes to safety features; (B). changes to the game itself; (C). switches to show; (D). cuts to the boys
			Specialized and Knowledge Reasoning
TRAM	temporal	3-Way MC	Then the green ball told the orange ball that blue ball was stupid. How long was the green ball talking to the orange ball? (A). 5 weeks; (B). 24 hours; (C). 15 seconds.
NumerSense	numerical	Number	Complete the sentence by filling in $<$ mask> with the most appropriate number. A classical guitar has $<$ mask> strings. $\rightarrow$ six
PIQA	physical	2-Way MC	To reach the physical goal: trees, choose the more sensible solution: (A). provide homes for people; (B). provide homes for animals.
QASC	science	8-Way MC	Crabs scavanage and uses dead snail shells for what? (A). RNA; (B). homes; (C). making holes; (D). damage; (E). a hosta; (F). Protein; (G). matter; (H). building a nest.
RiddleSense	riddle	5-Way MC	Something very helpful if you want to go gently down a stream. (A). raft; (B). roll down hill; (C). rowboat; (D). water; (E). roll over.
			Social and Ethical Reasoning
Social IQa	social	3-Way MC	Kai observed and improved Jan's work. What will Kai want to do next? (A). be effective; ( <b>B</b> ). get praise; (C). get incredible.
ETHICS	moral	2-way MC	After I got the flu virus I called my friends to tell them that I had the virus. Is this acceptable? (A). acceptable; (B). unacceptable.
			Visual Understanding and Bassoning

Visual Understanding and Reasoning



VCR

4-Way MC

visual

1. What is wrong with Person 2? (A). He is not happy with what is being said to him over the telephone; (B). He is feeling depressed; (C). He is high on pot; (D). Someone has pushed him and he's falling. 2. Given the question: What is wrong with Person 2?, and the answer to the question: He is feeling depressed, what is the rationale behind this answer? (A). Person 1 is talking to him probably trying to cheer him up; (B). He looks sad and is drinking; (C). He is walking with his head down; (D). He is slumped down on bed and his eyes are closed.

Question: What can be used to find your way if you're lost in the dark? (A) lamphreys (B) Ultraviolet shielding (C) Light shining into it (D) the speed of hearing (E) The north star (F) How to check the temperature (G) Observing tidal waves (H) osculum 🔸 Gemini Pro Response: (E) The north star Rationale Ouestion: What is the rationale behind the answer to the question? **Gemini Pro Response:** (E) The north star is the most suitable answer based on its unique position in the Northern Hemisphere and its ability to serve as a reliable reference point for determining north, which can be helpful in navigating back on track. Here's why the other options are not suitable: • (A) Lampreys: These are eel-like fish and have no navigational value. • (B) Ultraviolet shielding: While it can protect you from UV rays, it won't help you find your way. • (H) Osculum: This refers to a kiss and wouldn't be relevant to finding your way in the dark.

(a) Example of a correct response and rationale explanation from Gemini Pro.

Question: What is a compass used to navigate? (A) density (B) lamphreys (C) marine (D) outer space (E) their speed (F) Antarctica (G) non-magnetic planets (H) the Atlantic. **GPT-3.5 Turbo Response:** (C) marine Rationale Question: What is the rationale behind the answer to the question? GPT-3.5 Turbo Response: The rationale behind the answer (C) marine for the question about what a compass is used to navigate is based on the primary and traditional use of compasses. Compasses work by aligning with the Earth's magnetic field, specifically pointing towards the Earth's magnetic north pole. This property makes them an invaluable tool for navigation at sea. When you're out in the open water, it's easy to lose sight of landmarks, and without the help of a compass, it can be challenging to maintain a sense of direction. By using a compass, mariners can determine their heading and ensure they are on the right course. Note: The bolded sections in black highlight the generalization of compass use and the lack of specificity in relation to the question's requirement for a particular navigational context, leading to an incorrect interpretation.

(b) Example of an incorrect response and rationale explanation from GPT-3.5 Pro.

Figure 5: Model reasoning correctness justification examples. The sample questions are from the QASC dataset, with the correct answers highlighted in bold red. In example (a), Gemini Pro exhibits methodical reasoning by exclusion, carefully considering all options to reach the most logical conclusion. Conversely, example (b) illustrates GPT-3.5 Turbo's tendency towards unconventional logic, which can result in imaginative yet atypical answers. These instances emphasize the diverse strategies different models apply to commonsense reasoning tasks, revealing their distinct capabilities and limitations in such contexts.