Benchmarking Chinese Commonsense Reasoning of LLMs: From Chinese-Specifics to Reasoning-Memorization Correlations

Anonymous ACL submission

Abstract

We introduce CHARM, the first benchmark for 002 comprehensively and in-depth evaluating the commonsense reasoning ability of large language models (LLMs) in Chinese, which cov-004 005 ers both globally known and Chinese-specific 006 commonsense. We evaluated 7 English and 12 Chinese-oriented LLMs on CHARM, em-007 ploying 5 representative prompt strategies for improving LLMs' reasoning ability, such as Chain-of-Thought. Our findings indicate that 011 the LLM's language orientation and the task's domain influence the effectiveness of the prompt strategy, which enriches previous research find-013 ings. We built closely-interconnected reason-015 ing and memorization tasks, and found that some LLMs struggle with memorizing Chinese commonsense, affecting their reasoning 017 ability, while others show differences in reason-019 ing despite similar memorization performance. We also evaluated the LLMs' memorizationindependent reasoning abilities and analyzed the typical errors. Our study precisely identified the LLMs' strengths and weaknesses, providing the clear direction for optimization. It can also serve as a reference for studies in other fields. We will release CHARM on GitHub.

1 Introduction

027

Commonsense reasoning is important for the enhancement of the large language models (LLMs) (Bommasani et al., 2021; Achiam et al., 2023) towards artificial general intelligence (AGI) (Davis and Marcus, 2015), therefore requires thorough evaluations. Numerous benchmarks evaluate the commonsense reasoning of LLMs, but most are English-based, limiting non-English evaluations (Davis, 2023). This paper focuses on assessing LLMs' commonsense reasoning in a Chinese context. Currently, some commonsense reasoning benchmarks in Chinese are simply English translations (Conneau et al., 2018; Ponti et al., 2020;

Lin et al., 2022), which overlooks unique Chinese cultural, linguistic, regional, and historical aspects. These factors matter when Chinese users use the LLM, hence should be included in benchmarks. To effectively tackle this, we introduce CHARM, the benchmark designed to thoroughly and in-depth assess the abilities of LLMs in Chinese commonsense reasoning. It covers two domains: globally accepted commonsense (global domain) and Chinese-specific commonsense (Chinese domain). The latter includes 7 aspects: History (H), Traditional Culture and Arts (CA), Daily Life and Customs (LC), Entertainment (E), Public Figures (F), Geography (G), and Chinese Language (L). Therefore CHARM allows a thorough evaluation of LLMs' reasoning in a Chinese context.

041

042

043

044

045

047

049

052

053

055

057

059

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

075

076

077

078

081

Prompt strategies like Chain of Thought (CoT) (Wei et al., 2022) can significantly improve LLMs' reasoning performance (Wang et al., 2022, 2023b). Particularly, as the training corpus of LLMs is primarily in English (Touvron et al., 2023a), studies (Shi et al., 2022; Huang et al., 2023a; Zhang et al., 2023a) have shown that for non-English reasoning tasks, some LLMs perform better when reasoning in English than the native language. We evaluate 7 English and 12 Chinese-oriented LLMs on CHARM, employing 5 representative prompt strategies. The result shows that prompt strategies' effectiveness depends on the LLMs' orientation and the benchmark task's domain, which enriches prior research and guides performance assessment and strategy choice for non-English LLMs.

LLMs' commonsense reasoning relies on memorization. Exploring the correlation between memorization and reasoning offers insights into LLMs, aiding deeper understanding and suggesting ways to enhance these abilities(Bian et al., 2023). Some benchmarks (Yu et al., 2023; Wang et al., 2023a; Fei et al., 2023) aid the research of memorizationreasoning relationships by incorporate tasks for assessing knowledge memorization and applica-



Figure 1: Construction of CHARM. CHARM encompasses both global and Chinese-specific commonsense. CHARM consists closely-interconnected reasoning and memorization tasks.

tion (like reasoning). However, they often use the existing and disparate datasets for different tasks, resulting in a lack of intrinsic connections between 084 these tasks. For instance, the question Q_{rea} tests 085 the LLM's reasoning with the knowledge piece K. However, in memorization tasks, there probably is not any matching questions to determine if the LLM has effectively memorized K. Hence, if the LLM fails on Q_{rea} , it's unclear whether due to poor reasoning or forgetfulness of K. This results in the disjointed evaluation of memorization and reasoning, failing to uncover their intrinsic links. To address this limitation, we selected suitable reasoning tasks from CHARM's Chinese domain, and built related memorization questions for each reasoning question (see Figure 1). This design produces the closely-interconnected reasoning and memorization tasks, therefore allows for not only the concurrent evaluation of the two abilities, but 100 also the assessment of memorization-independent 101 reasoning, providing the clear guidance for the 102 LLMs' enhancement. 103

The contributions of this paper are as follows:

• We present CHARM, the first benchmark for comprehensively evaluating the LLMs' commonsense reasoning ability in Chinese, by encompassing not only the global but also the Chinese-specific commonsense.

107

108

109

110

111

112

 We evaluated the representative prompt strategies on CHARM. Results showed that LLMs' orientation and the task's domain affect prompt strategy performance, which enriches previous research findings.

113

114

115

116

117

118

119

120

121

122

123

125

126

127

128

129

130

131

132

134

135

136

137

138

139

140

141

142

• In CHARM, we built closely-interconnected reasoning and memorization tasks in Chinese commonsense domain, allowing for in-depth understanding the correlation between these abilities and precisely identifying the LLMs' strengths and weaknesses. The design approach could serve as the reference for other fields.

2 Related Work

Commonsense Reasoning Benchmarks There are lots of English commonsense reasoning benchmarks or datasets (Davis, 2023), some of which have been translated to various languages, including Chinese (Conneau et al., 2018; Ponti et al., 2020; Lin et al., 2022). There are also some native Chinese benchmarks related to commonsense reasoning. LogiQA (Liu et al., 2020, 2023), sourced from the Chinese civil servant exam questions, evaluates the LLMs' reading comprehension and logical reasoning across various knowledge types, including basic commonsense knowledges. CLUE (Xu et al., 2020), the comprehensive Chinese language understanding benchmark, includes the Natural Language Inference (NLI) tasks, which require the LLMs' commonsense reasoning abilities. CMMLU (Li et al., 2023) is the comprehensive Chinese benchmark that encompasses exam questions from the 67 subjects. It includes tasks that

234

test elementary commonsense and logical reasoning abilities. CORECODE (Shi et al., 2023) is the dataset designed for commonsense reasoning and conflict detection in Chinese, presented in the dialogue format.

143

144

145

146

147

148

149

150

151

152

153

155

156

157

158

160

161

163

164

165

166

168

169

170

171

172

173

174

176

177

178

181

185

186

191

Prompt Strategy Prompt strategies such as CoT (Wei et al., 2022) can effectively boost the reasoning capabilities of LLMs (Wang et al., 2022, 2023b). Notably, as the LLM training corpus is primarily in English (Touvron et al., 2023a), research reveals that for reasoning tasks in non-English languages, some LLMs exhibit superior performance when reasoning in English as opposed to the native language (Shi et al., 2022; Zhang et al., 2023a). (Kim et al., 2023) proposed a novel cross-language transfer prompt method, which uses both the source and target languages to construct examples.

Benchmarks on Correlations of Memorization and Reasoning There are benchmarks which assess both the knowledge memorization and reasoning capabilities of the LLMs within specific domains. For instance, KoLA (Yu et al., 2023), with its focus on world knowledge, includes tasks related to knowledge memorization and application (reasoning). SeaEval (Wang et al., 2023a), emphasizing cross-language consistency and multicultural reasoning, involves tasks for cultural understanding and complex reasoning. There are also benchmarks aimed at specialized fields, like LawBench (Fei et al., 2023), which include tasks for both memorization and application.

3 CHARM

CHARM is built for comprehensive and in-depth evaluation of LLMs in Chinese commonsense reasoning and revealing the intrinsic correlation between memorization and reasoning. Therefore, CHARM covers two domains, global and Chinese, using carefully selected tasks for comprehensive coverage. In addition, we chose reasoning tasks and constructed the closely-tied memorization tasks. The construction and main features of CHARM are in Figure 1. The detailed composition of CHARM is in Table 1.

3.1 Commonsense Domain

Global commonsense domain consists of universally understood commonsense. It covers objects and aspects of modern life that an individual should be aware of. It includes foundational knowledge that someone with a basic modern education is expected to know. When it involves individuals, they are globally recognized figures.

Chinese commonsense domain encompasses Chinese-specific elements. We categoried them into 7 aspects:

History (H) includes important events and figures in Chinese history, China's dynasties, and other basic facts and shared knowledge about the history of China.

Traditional Culture and Arts (CA) encompasses Chinese traditional cultural arts, literary works, and traditional lifestyles.

Daily Life and Customs (LC) includes modern Chinese daily routines, clothing, food, housing, transportation festivals and so on.

Entertainment (E) includes the movies, television programs, music, and other entertainments in modern Chinese daily life.

Public Figures (F) encompasses the public figures well-known in Chinese society.

Geography (G) includes China's geographical distribution, natural landscapes, and characteristic regional cultures.

Chinese Language (L) includes the fundamentals of the Chinese language, such as Chinese characters, idioms and so on.

For the two domains, especially for the above 7 aspects, we collected corresponding entities, forming the lists as shown in the Figure 1 and 5. Most of the entities are selected from Gaokao Bench¹(Zhang et al., 2023b), Douban², Hupu³. Some entities are collected with the help of searching engines. We only collect the entities that are well-known in China. These entities are then used to create the commonsense reasoning questions, which belong to the corresponding domain and aspect.

3.2 Reasoning Tasks

When designing the reasoning tasks in CHARM, we bear two criteria in mind. First, the tasks should span both commonsense domains, particularly the 7 Chinese aspects. Second, the global and Chinese tasks should have identical types and settings, differing only in their commonsense domains. From the

¹Gaokao Bench is the collection of China's university entrance exam questions, which contributes to all the 7 aspects. ²https://www.douban.com/ is the popular user-centric

cultural review platform in China, which mainly contributes to the *Entertainment* aspect.

³https://www.hupu.com/ is the large sports community popular in China, which mainly contributes to the *Public Figures* aspect.

Task Type	Task	Domain	Chinese aspects	Construction	Question Type	# Question
	Anachronisms Judgment (AI)	Chinese	H, AC, LC, F	[H]	2-option MCQ	150
	i internomono vauginent (i to)	global	-	[[T][H]	2-option MCQ	150
	Time Understanding (TII)	Chinese	H, AC, LC	[H]	4-option MCQ	100
	Time Onderstanding (10)	global	-	[T]	5or6-option MCQ	100
	Sequence Understanding (SaU)	Chinese	H, CA , LC , G , L	[H]	4-option MCQ	100
	Sequence Onderstanding (SqU)		-	[T][H]	4-option MCQ	100
Dessening	Mayia and Musia Decommondation (MMD)	Chinese	Ε	[H]	4-option MCQ	50
Reasoning	Movie and Music Recommendation (MINIK)	global	-	[T]	4-option MCQ	50
	Court He least a l'an (Coll)	Chinese	F	[H]	2-option MCQ	200
	sport Understanding (SpU)	global	-	[H]	2-option MCQ	200
	Network Control of the control of th	Chinese	G, E, L,	[S][H]	3-option MCQ	100
	Natural Language Interence (NLI)	global	-	[S]	3-option MCQ	100
	Basding Comprehension (BC)	Chinese	all 7 aspects	[S]	4-option MCQ	200
	Reading Comprehension (RC)	global	-	[S]	4-option MCQ	200
	Anachronisms Judgment (AJ)	Chinese	H,AC,LC,F	[H]	Free-form QA	135
Manadardian	Time Understanding (TU)	Chinese	H, AC, LC	[H]	Free-form QA	83
Memorization	Iemorization Movie and Music Recommendation (MMR)		Ε	[H]	Free-form QA	399
	Sport Understanding (SpU)	Chinese	F	[H]	Free-form QA	127

Table 1: Overview of CHARM. The question numbers of reasoning and memorization tasks are 1800 and 744.

existing English commonsense reasoning datasets (Davis, 2023; Suzgun et al., 2022), we selected the following 7 tasks:

Anachronisms Judgment (AJ) necessitates the LLM to identify anachronisms in provided sentences. This involves the LLM understanding the era associated with well-known figures, items, and events to facilitate commonsense-based reasoning. Global domain questions are the mix of translations⁴ and handcrafted, while all Chinese domain questions are handcrafted.

Time Understanding (TU) requires the LLM infers a time (including year, date, moment, etc.) based on a given context, which necessitates the fundamental understanding of time-related commonsense and the capacity for mathematical reasoning. All question in the global domain are translations⁵ and all in Chinese domain are handcrafted.

Sequence Understanding (SqU) requires the LLM sort a series of entities according to time or occurrence order, requiring logical reasoning based on commonsense. The global domain questions are the mix of translations⁶ and handcrafted; while all in the Chinese domain are handcrafted.

Movie and Music Recommendation (MMR) necessitates the LLM identifies the most similar matches to a variety of movies or music tracks, requiring the understanding of these popular movies and music and ability to identify their commonalities. All global domain questions are translations⁷, and all in the Chinese domain are handcrafted.

266

267

268

269

270

271

272

273

274

275

276

277

278

279

282

284

285

287

290

291

292

293

294

295

297

300

301

Sport Understanding (SpU) involves a crafted sentence with a famous athlete and a common sport action, and the LLM must assess its credibility, which demands understanding of sports and commonsense judgement. The questions in both domains are handcrafted, refering (Suzgun et al., 2022).

Natural Language Inference (NLI) gives two sentences and asks the LLM to classify their relationship as entailment, contradiction, or neutral, necessitating commonsense-based reasoning and judgement. All global domain questions are selected from CLUE (Xu et al., 2020); the questions in the Chinese domain are partly from CLUE, and partly handcrafted.

Reading Comprehension (RC) gives a passage of text, and the LLM is required to reason based on it. All question in both domains are selected from LogiQA (Liu et al., 2020, 2023).

The chosen tasks adequately cover both the commonsense domains, particularly the 7 aspects of the Chinese commonsense domain. This coverage enables a comprehensive assessment of LLMs' commonsense reasoning ability in Chinese. Moreover, the Chinese-domain questions can be created following the similar types and settings as their global counterpart, facilitating the cleaner comparison of the LLMs' performance across the domains.

All questions in the CHARM reasoning tasks are multiple-choice questions. Detailed information is in Table 1. Question examples of the tasks are in Figure 6 in Appendix B. We use regular expressions to extract the preferred choice from the generation of the LLMs (Huang et al., 2023b; Li et al., 2023) and use *accuracy* as the metric.

⁴https://github.com/google/BIG-bench/tree/ main/bigbench/benchmark_tasks/anachronisms

⁵https://github.com/google/BIG-bench/tree/ main/bigbench/benchmark_tasks/date_understanding ⁶https://github.com/google/BIG-bench/tree/

main/bigbench/benchmark_tasks/logical_sequence

⁷https://github.com/google/BIG-bench/ tree/main/bigbench/benchmark_tasks/movie_ recommendation

Task	AJ	TU	SpU	MMR
Avg. # related memorization questions	1.9	3.2	2.0	8.0

Table 2: Averaged number of related memorization questions per reasoning question for each task.



Figure 2: Distribution of CHARM construction.

3.3 Memorization Task

302

303

304

307

310

312

313

314

315

317

321

323

324

325

326

327

333

339

Shared commonsense knowledge pieces serve as a link between reasoning and memorization questions. From the 7 reasoning tasks, we chose 4 that can be readily associated in this manner, AJ, TU, MMR, SpU, refered as the *Memorization-Reasoning-Interconnected (MRI) tasks*, and built the reltated memorization questions.

Construction We first extract the commonsense knowledge pieces related to the entities in the corresponding reasoning questions. Information about each entity was collected to the degree sufficient to address the associated reasoning question, and then used to formulate the memorization questions. Following the *Knowledge Memorization* task in KoLA (Yu et al., 2023), we choose free-form QA instead of multiple-choice or true/false questions, which can effectively avoid the impact of randomness. All memorization questions are handcrafted. Question examples are in Figure 7 in Appendix C. The averaged number of related memorization questions for each reasoning question are shown in Table 2.

Judgement and Metric For the MMR task, we use a rule-based matching method for evaluation; for the other three tasks, we use Chat-GPT for judgement. We use *accuracy* as the metric.

3.4 Construction and Quality Assurance

The question construction of CHARM involves the following three methods:

Handcraft [H] The questions are created mannualy by the authors, based on the entities in §3.1. For reasoning questions in Chinese commonsense domain, we usually refer to the global counter-part and construct them with the same design idea and style, as shown in Figure 6.

Translate [T] We translate the existing English questions and only keep those consistent with the global commonsense domain.

Models	Open Source?	Model Size	Primary Language
LLaMA-2	Yes	7B, 13B, 70B	English
Vicuna	Yes	7B, 13B	English
GPT-3.5	No	undisclosed	English
GPT-4	No	undisclosed	English
ChatGLM3	Yes	6B	Chinese
Baichuan2	Yes	7B, 13B	Chinese
InternLM2	Yes	7B, 20B	Chinese
Yi	Yes	6B, 34B	Chinese
DeepSeek	Yes	7B, 67B	Chinese
Qwen	Yes	7B, 14B, 72B	Chinese

Table 3:	LLMs	evaluated i	n our ex	periments

Select [S] We select commonsense reasoning questions from existing Chinese datasets.

Detailed information of the construction of all tasks are shown in Table 1. The distribution of the three methods is shown in Figure 2.

We ensure the quality through these steps: First, two annotators independently verify each question, whether handcrafted or translated. Second, annotators are given only questions, not our answers, and use external resources to respond. They can also flag problematic questions. Last, a question passes if both annotators' answers match ours; if not, it's reviewed and modified by multiple authors.

4 Experimental Setup

4.1 Language Models

We evaluated the currently commonly used LLMs, which can be divided into two categories: (1) 7 English LLMs, including GPT series (Achiam et al., 2023), LLaMA-2 (Touvron et al., 2023b), and Vicuna (2) 12 Chinese-oriented LLMs, including ChatGLM3⁸, Baichuan2 (Yang et al., 2023), InternLM2 (Team, 2023), Yi⁹, DeepSeek(Bi et al., 2024) and Qwen (Bai et al., 2023). For open-source models, we choose the chat version instead of the base version. For closed-source models, we use the official API¹⁰. Detailed information is in Table 3.

4.2 **Prompt Strategies**

We selected 5 commonly used prompt strategies, and assessed the performance of the 19 LLMs on CHARM reasoning task:

Direct: The LLM does not perform intermediate reasoning and directly predicts the answer.

ZH-CoT: The LLM conducts intermediate reasoning (Wei et al., 2022) in Chinese before producing the answer.

340

341

342

360 361

363

364

365

366

367

368

369

370

371

372

373

374

356

357

358

⁸https://huggingface.co/THUDM/chatglm3-6b-32k

⁹https://github.com/01-ai/Yi ¹⁰We used the gpt-3.5-turbo-1106 version for GPT-3.5 and the gpt-4-1106-preview version for GPT-4.

LLM			Chine	se Comm	onsense l	Domain					Glob	al Commo	onsense I	Domain		
	AJ	TU	SqU	MMR	SpU	NLI	RC	Avg.	AJ	TU	SqU	MMR	SpU	NLI	RC	Avg.
GPT-3.5-1106	82.00	39.0	65.0	42.0	80.5	61.0	50.5	60.00	90.00	94.0	87.0	46.0	88.5	66.0	49.5	74.43
GPT-4-1106	95.33	60.0	85.0	74.0	86.0	77.0	62.5	77.12	95.33	98.0	97.0	66.0	90.0	72.0	72.0	84.33
LLaMA-2-7B	50.67	36.0	11.0	14.0	49.5	52.0	8.0	31.60	62.67	17.0	14.0	16.0	49.5	22.0	13.0	27.74
LLaMA-2-13B	54.67	33.0	38.0	30.0	58.0	47.0	38.0	42.67	66.67	24.0	39.0	50.0	53.5	57.0	33.5	46.24
LLaMA-2-70B	61.33	37.0	52.0	32.0	55.0	56.0	41.5	47.83	72.67	84.0	73.0	42.0	64.0	61.0	41.5	62.60
Vicuna-7B-v1.5-16k	50.67	29.0	34.0	32.0	51.0	49.0	35.5	40.17	45.33	64.0	37.0	26.0	58.5	52.0	32.5	45.05
Vicuna-13B-v1.5-16k	64.67	25.0	32.0	26.0	51.5	60.0	40.0	42.74	72.67	74.0	41.0	50.0	68.0	61.0	36.0	57.52
ChatGLM3-6B-32k	65.33	40.0	59.0	38.0	77.0	72.0	37.5	55.55	34.00	69.0	71.0	28.0	75.5	63.0	34.0	53.50
Baichuan2-7B	72.67	41.0	48.0	38.0	72.0	53.0	49.5	53.45	55.33	65.0	54.0	26.0	60.5	59.0	29.0	49.83
Baichuan2-13B	83.33	40.0	48.0	46.0	72.5	66.0	51.5	58.19	77.33	74.0	58.0	40.0	71.0	61.0	39.0	60.05
InternLM2-7B	88.00	38.0	58.0	38.0	76.0	81.0	25.0	57.71	74.67	80.0	62.0	20.0	78.0	76.0	23.5	59.17
InternLM2-20B	86.00	55.0	54.0	44.0	74.5	80.0	23.0	59.50	82.67	83.0	61.0	14.0	74.5	72.0	27.0	59.17
Yi-6B	72.67	32.0	47.0	32.0	75.0	50.0	42.0	50.10	79.33	63.0	43.0	14.0	70.5	57.0	33.5	51.48
Yi-34B	94.00	55.0	89.0	76.0	88.5	72.0	51.5	75.14	88.67	92.0	87.0	56.0	89.0	70.0	47.5	75.74
DeepSeek-7B	79.33	34.0	50.0	50.0	79.5	57.0	31.5	54.48	68.00	76.0	47.0	50.0	72.5	59.0	32.5	57.86
DeepSeek-67B	93.33	57.0	83.0	92.0	87.5	77.0	34.5	74.90	90.00	95.0	86.0	22.0	88.0	73.0	39.0	70.43
Qwen-7B	73.33	38.0	55.0	48.0	71.0	57.0	49.5	55.98	74.67	78.0	69.0	50.0	72.5	55.0	36.0	62.17
Qwen-14B	88.00	54.0	77.0	60.0	82.5	66.0	55.0	68.93	84.00	83.0	83.0	44.0	84.5	71.0	40.0	69.93
Qwen-72B	96.67	59.0	91.0	84.0	86.5	84.0	67.5	81.24	94.00	92.0	93.0	64.0	93.0	71.0	63.5	81.50

Table 4: Accuracy on CHARM reasoning tasks. We selected the empirically optimal prompt strategy: XLT for English LLMs and ZH-CoT for Chinese-oriented LLMs. **Bold** and <u>underline</u> represent the first and second place respectively. Detailed results are in Table 8 and Table 9 of Appendix E.

EN-CoT: The reasoning process of CoT is in English for the Chinese questions(Shi et al., 2022).

Translate-EN: We use the DeepL api¹¹ to translate our benchmark into English, and then use English CoT for reasoning (Zhang et al., 2023a).

XLT: The template prompt (Huang et al., 2023a) is used to change the original question into an English request, solve it step by step, and finally format the answer for output .

The examples for each prompt strategy are in Figure 8 in Appendix D. For all prompt strategies, we use the 3-shot setting.

5 Results and Analysis

5.1 Integrated Reasoning Performance

We show the performance of the 19 LLMs on CHARM reasoning tasks in Table 4. We only choose one representative prompt strategy: XLT for English LLMs and ZH-CoT for Chinese-oriented LLMs, which is based on our empirical conclusion in §5.2. The LLMs' performance on the 7 aspects of the Chinese commonsense domain are shown in Table 10 in Appendix F.

Commonsense Domain We found that the LLMs exhibit inconsistent performance in the global and Chinese commonsense domains. The rankings of the English LLMs dropped in the Chinese domain compared to the global domain. For instance, GPT4 ranks first in the global domain, but in the Chinese domain, Qwen-72B outperforms all, pushing GPT4 to the second. In the Chinese domain, the performance of LLaMA-2-70B is even worse than many

	Prompt	Avg. all LLMs	Avg. CN-LLMs	Avg. EN-LLMs
Avg.	Direct	46.36	48.51	42.68
all	ZH-CoT	56.59	62.33	46.75
domains	EN-CoT	54.39	58.18	47.89
	Translate-EN	53.77	55.39	50.98
	XLT	56.73	58.98	52.86
Avg.	Direct	45.59	47.97	41.51
Chinese	ZH-CoT	56.22	62.10	46.15
domain	EN-CoT	51.92	56.34	44.35
	Translate-EN	47.04	47.59	46.10
	XLT	53.63	56.41	48.87
Avg.	Direct	47.13	49.05	43.85
global	ZH-CoT	56.96	62.57	47.35
domain	EN-CoT	56.85	60.01	51.44
	Translate-EN	60.50	63.20	55.87
	XLT	59.82	61.56	56.84

Table 5: Averaged accuracy on CHARM reasoning tasks. "CN-LLMs" means the 12 Chinese-oriented LLMs, "EN-LLMs" means the 7 English LLMs.

Chinese-oriented LLMs in the 6B-7B size range. However, in the global domain, LLaMA-2-70B is better than all Chinese-oriented LLMs up to 20B in size, except for Qwen-14B.

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

5.2 Prompt Strategy Selection

We tested the combinations of the 19 LLMs and the 5 prompt strategies in CHARM reasoning tasks. Detailed results are in Table 8 and Table 9 in Appendix E. To draw some empirical conclusions, we analyze along the following two dimensions:

- Dim1: global or Chinese commonsense domain.
- Dim2: English or Chinese-oriented LLMs.

We average the 19×5 LLM-prompt combinations along the above two dimensions, and the obtained results are in the Table 5. *From the LLM dimension*, it's clear that various LLMs prefer different prompt strategies: XLT consistently excels for English LLMs among the 5 strategies, while

400

401

402

403

404

405

[&]quot;https://www.deepl.com/translator

for Chinese-oriented LLMs, despite some complexity, ZH-CoT generally performs best. *From the commonsense domain dimension*, strategies that use English for reasoning (like XLT, Translate-EN, etc.) are suitable for the global domain; however, ZH-CoT generally performs better in the Chinese domain.

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

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

The conclusion here differs from previous studies (Shi et al., 2022; Huang et al., 2023a), which suggested that employing English for non-English reasoning tasks was more effective than using the native language. These previous studies had limitations, focusing only on English LLMs and neglecting the many Chinese-oriented LLMs developed since 2023. Furthermore, most benchmarks in these studies were merely translations from English, lacking unique cultural and linguistic characteristics in Chinese. The empirical findings with CHARM in this paper have somewhat alleviated those limitations, leading to more current and comprehensive conclusions, and of course still have the limitations, which are detailed in §6.

5.3 Integrated Reasoning vs Memorization

We evaluated the correlation between the integrated reasoning and the memorization on the *MRI* tasks, as mentioned in §3.2. The average performance of the LLMs on the 4 *MRI* tasks is in Figure 3. Detailed performance on each task is in Figure 9 in Appendix G.1.

As shown in Figure 3, the 19 LLMs can be roughly divided into the three types:

• Type I: Low memorization and low integrated reasoning ability. We found that apart from OpenAI's GPT series, all other English LLMs belong to this type.

• Type II: High memorization and medium integrated reasoning ability. GPT3.5 and all Chinese-oriented LLMs below 30B belong to this type. It's worth noting that some LLMs have high memorization performance, but relatively poor integrated reasoning ability.

• Type III: Ultra-high memorization and high integrated reasoning ability. This category includes GPT4 and the three Chinese-oriented LLMs that exceed a size of 30B.

The above findings offer clear guidance for the enhancement of LLMs' reasoning abilities in Chinese commonsense domain. For Type I, the limitation lies in the memorization. For Type II, there should be further improvement in understanding, applying knowledge, and reasoning abilities.



Figure 3: **Averaged** accuracy across the 4 *MRI* tasks in the Chinese commonsense domain.

rank	Integrated Reasoning	Memorization-inde	pendent Reasoning
	8	FRMM	MIB
1	Qwen-72B	Yi-34B (†3)	GPT-4 (<u></u>)
2	DeepSeek-67B	GPT-4 (1)	Qwen-72B (1)
3	GPT-4	DeepSeek-67B (↓1)	Yi-34B (1)
4	Yi-34B	Qwen-72B (13)	DeepSeek-67B (↓2)
5	Qwen-14B	Qwen-14B (-)	Qwen-14B (-)
6	InternLM2-20B	InternLM2-7B (12)	InternLM2-7B (12)
7	GPT-3.5	GPT-3.5 (-)	GPT-3.5 (-)
8	InternLM2-7B	InternLM2-20B (12)	InternLM2-20B (12)
9	DeepSeek-7B	Baichuan2-13B (1)	DeepSeek-7B (-)
10	Baichuan2-13B	DeepSeek-7B (1)	Baichuan2-13B (-)
11	Qwen-7B	Yi-6B (13)	ChatGLM3-6B(²)
12	Baichuan2-7B	ChatGLM3-6B (1)	Baichuan2-7B (-)
13	ChatGLM3-6B	Qwen-7B $(\downarrow 2)$	Yi-6B (1)
14	Yi-6B	Baichuan2-7B (12)	Qwen-7B $(\downarrow 3)$
15	LLaMA-2-70B	LLaMA-2-70B (-)	LLaMA-2-70B (-)
16	LLaMA-2-13B	LLaMA-2-13B (-)	LLaMA-2-13B (-)
17	Vicuna-13B	Vicuna-13B (-)	Vicuna-13B (-)
18	Vicuna-7B	LLaMA-2-7B (1)	Vicuna-7B (-)
19	LLaMA-2-7B	Vicuna-7B (1)	LLaMA-2-7B (-)

Table 6: Leaderboard on the *MRI* tasks. We propose two methods, i.e. FRMM and MIC, to compare the LLMs' **memorization-independent reasoning**, as detailed in Appendix H. The arrows and numbers in brackets in the last two columns indicate changes in ranking order relative to the second column.

In addition, we also evaluated the correlation between memorization and integrated reasoning during the LLM pre-training process, details can be found in Figure 10 in Appendix G.2.

The results clearly indicate that strong memorization is the foundation of integrated reasoning. Weak memorization leads to poor reasoning, as shown by Type I LLMs. Also, factors other than memorization can cause significant differences in reasoning abilities among LLMs with similar memorization.

5.4 Memorization-Independent Reasoning

We propose two methods to compare the LLMs' memorization-independent reasoning on the *MRI* tasks. The detail of the *Filtering Reasoning Questions based on Mono-LLM-Memorization* (FRMM) is in Appendix H.1, and *Memorization-Independent*

Battles among LLMs (MIB) in Appendix H.2. The results are in Table 6. When comparing the leaderboards for integrated and memorizationindependent reasoning, Type III LLMs rank at the forefront and the Type I rank at the end in all leaderboards. There is a slight variation in the ranking order within the three types of LLMs.

491

492

493

494

496

497

498

499

500

501

502

503

509

510

511

513

514

515

517

518

519

521

523

524

525

526

527

529

530

531

532

534

536

538

541

542

For the in-depth analysis, we chose Vicuna-13B, Qwen-7B, Qwen-72B as representatives for Type I, II, and III LLMs, and filtered out the reasoning questions in the *MRI* tasks, only keeping those with correct answers to the related memorization questions, same as the FRMM in Appendix H.1. This ensures the LLM has sufficiently memorized the commonsense knowledge required for the retained reasoning questions, thereby minimizing the impact of memorization on reasoning. There are totally 500 reasoning questions in the 4 *MRI* tasks, and the numbers of the retained are 126, 332 and 411 for Vicuna-13B, Qwen-7B and Qwen-72B respectively, as shown in Table 7.

Error Types If the LLMs provided incorrect answers for the retained reasoning questions, these errors can be referred to as memorization-independent reasoning errors. We conducted the manual review and analysis of their reasoning process, and classified the errors into 4 main categories.

• Understanding Error In this case, the LLM was unable to accurately comprehend the question, including misunderstanding the content, ignoring or even modifying important information in the premise, and failing to grasp the core query of the question.

• **Knowledge Error** The LLM incorporated inaccurate knowledge during the reasoning process. It's important to highlight that the knowledge pieces related to the reasoning question were previously examined in the related memorization questions, which the LLM answered **correctly**. However, the LLM output incorrect information during the reasoning phase.

• Logical Error The LLM made logical reasoning errors, such as mathematical reasoning errors, inability to reach the correct conclusion based on sufficient information, or reaching the correct conclusion but outputing the wrong option.

• **Other Errors** These are other scattered, relatively rare types of errors.

We show examples of each type of errors in Figure 12 in Appendix I. The distribution of these error types are shown in Figure 4.

Discussions Obviously, the majority of errors

Models	LLM type	# Original	# Retained	# Incorrect
Vicuna-13B	Type I	500	126	58
Qwen-7B	Type II	500	332	115
Qwen-72B	Type III	500	411	65

Table 7: Memorization-based filtering of reasoning questions. "Incorrect" means the incorrectly answered questions among the **retained**.



Figure 4: Distribution of the memorization-independent reasoning errors

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

562

563

564

565

566

567

568

569

570

571

572

574

575

are from logical reasoning mistakes and knowledge inaccuracies, which further provides the directions for LLMs' enhancement. As for knowledge errors, prior studies (Bian et al., 2023; Allen-Zhu and Li, 2023) have indicated that the way LLMs remember and master knowledge is a relatively complex topic. Simple memorization doesn't guarantee that LLMs can apply this knowledge accurately and skillfully during the reasoning process.

6 Conclusion

This paper introduces CHARM, the first benchmark designed to comprehensively and thoroughly evaluate LLMs' commonsense reasoning in Chinese. CHARM encompasses two counterpart commonsense domains, global and Chinese-specific, with the carefully selected tasks. We evaluated the representative prompt strategies for improving LLMs' reasoning ability, and the empirical findings significantly enhances and supplements the conclusions of previous studies. CHARM comprises closelyinterconnected reasoning and memorization tasks, helping to reveal the intrinsic correlation between memorization and reasoning of LLMs. We have evaluated the strengths and weaknesses of different LLMs and conducted the detailed analysis of memorization-independent reasoning abilities. We hope that CHARM's approach to studying the correlations between memoriztion and reasoning can serve as a reference for similar research in other fields.

Limitations

This study conducted tests on combinations of the 19 LLMs and the 5 prompt strategies, resulting

in empirical conclusions. However, many existing 576 LLMs and prompt strategies have not yet been tested. 577 Furthermore, the best prompt strategy for the commonsense reasoning task for the LLMs, particularly in Chinese or other non-English languages, is not static and should progress with LLM technology. 581 This is influenced by three elements: (1) The new 582 prompt strategies are continuously proposed, which are likely more effective. (2) The new LLMs may have different prompt strategy preference, or be less 585 sensitive to prompt. (3) For other non-English lan-586 guages with high resources, future LLMs would be 587 continuously evolving and updating, and necessitate 588 ongoing updates in evaluation.

The automation of the construction and evaluation of CHARM needs further improvement, including the following: (1) Most of the questions in CHARM Chinese domain are manually constructed by the author. This limits the number of benchmark questions and the range of knowledge pieces covered. (2) Regarding memorization-independent reasoning, we chose only 3 LLMs as representative and manually categorized the types of errors within CHARM. In future research, we could employ robust LLMs, like GPT4, for automated error classification and statistical analysis.

Ethical Consideration

590

594

595

596

610

611

612

613

614

615

618

619

623

This work involves human annotation. We have provided appropriate compensation for all annotators. The total cost of annotation for the project is about 2.2k RMB. For all annotators, we explicitly inform them about the use of the data and require them to ensure that the questions included in CHARM do not involve any social bias, ethical issues or privacy concerns during the annotation process.

References

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
 - Zeyuan Allen-Zhu and Yuanzhi Li. 2023. Physics of language models: Part 3.2, knowledge manipulation. *arXiv preprint arXiv:2309.14402*.
 - Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, et al. 2023. Qwen technical report. *arXiv preprint arXiv:2309.16609*.

Xiao Bi, Deli Chen, Guanting Chen, Shanhuang Chen, Damai Dai, Chengqi Deng, Honghui Ding, Kai Dong, Qiushi Du, Zhe Fu, et al. 2024. Deepseek llm: Scaling open-source language models with longtermism. *arXiv preprint arXiv:2401.02954*.

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

671

672

673

674

675

676

- 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.
- Rishi Bommasani, Drew A Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, et al. 2021. On the opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258*.
- Alexis Conneau, Guillaume Lample, Ruty Rinott, Adina Williams, Samuel R Bowman, Holger Schwenk, and Veselin Stoyanov. 2018. Xnli: Evaluating crosslingual sentence representations. *arXiv preprint arXiv:1809.05053*.
- Ernest Davis. 2023. Benchmarks for automated commonsense reasoning: A survey. *arXiv preprint arXiv:2302.04752*.
- Ernest Davis and Gary Marcus. 2015. Commonsense reasoning and commonsense knowledge in artificial intelligence. *Communications of the ACM*, 58(9):92–103.
- Zhiwei Fei, Xiaoyu Shen, Dawei Zhu, Fengzhe Zhou, Zhuo Han, Songyang Zhang, Kai Chen, Zongwen Shen, and Jidong Ge. 2023. Lawbench: Benchmarking legal knowledge of large language models. *arXiv preprint arXiv:2309.16289*.
- Haoyang Huang, Tianyi Tang, Dongdong Zhang, Wayne Xin Zhao, Ting Song, Yan Xia, and Furu Wei. 2023a. Not all languages are created equal in llms: Improving multilingual capability by cross-lingual-thought prompting. *arXiv preprint arXiv:2305.07004*.
- Yuzhen Huang, Yuzhuo Bai, Zhihao Zhu, Junlei Zhang, Jinghan Zhang, Tangjun Su, Junteng Liu, Chuancheng Lv, Yikai Zhang, Jiayi Lei, et al. 2023b. Ceval: A multi-level multi-discipline chinese evaluation suite for foundation models. *arXiv preprint arXiv:2305.08322*.
- Sunkyoung Kim, Dayeon Ki, Yireun Kim, and Jinsik Lee. 2023. Boosting cross-lingual transferability in multilingual models via in-context learning. *arXiv preprint arXiv:2305.15233*.
- Haonan Li, Yixuan Zhang, Fajri Koto, Yifei Yang, Hai Zhao, Yeyun Gong, Nan Duan, and Timothy Baldwin. 2023. Cmmlu: Measuring massive multitask language understanding in chinese. *arXiv preprint arXiv:2306.09212*.

678

- 731 732

- Xi Victoria Lin, Todor Mihaylov, Mikel Artetxe, Tianlu Wang, Shuohui Chen, Daniel Simig, Myle Ott, Naman Goyal, Shruti Bhosale, Jingfei Du, et al. 2022. Fewshot learning with multilingual generative language models. In Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, pages 9019-9052.
- Hanmeng Liu, Jian Liu, Leyang Cui, Zhiyang Teng, Nan Duan, Ming Zhou, and Yue Zhang. 2023. Logiqa 2.0-an improved dataset for logical reasoning in natural language understanding. IEEE/ACM Transactions on Audio, Speech, and Language Processing.
- Jian Liu, Leyang Cui, Hanmeng Liu, Dandan Huang, Yile Wang, and Yue Zhang. 2020. Logiqa: A challenge dataset for machine reading comprehension with logical reasoning. arXiv preprint arXiv:2007.08124.
- Edoardo Maria Ponti, Goran Glavaš, Olga Majewska, Oianchu Liu, Ivan Vulić, and Anna Korhonen. 2020. Xcopa: A multilingual dataset for causal commonsense reasoning. arXiv preprint arXiv:2005.00333.
- Dan Shi, Chaobin You, Jiantao Huang, Taihao Li, and Deyi Xiong. 2023. Corecode: A common sense annotated dialogue dataset with benchmark tasks for chinese large language models. arXiv preprint arXiv:2312.12853.
- Freda Shi, Mirac Suzgun, Markus Freitag, Xuezhi Wang, Suraj Srivats, Soroush Vosoughi, Hyung Won Chung, Yi Tay, Sebastian Ruder, Denny Zhou, et al. 2022. Language models are multilingual chain-of-thought reasoners. arXiv preprint arXiv:2210.03057.
- Mirac Suzgun, Nathan Scales, Nathanael Schärli, Sebastian Gehrmann, Yi Tay, Hyung Won Chung, Aakanksha Chowdhery, Quoc V Le, Ed H Chi, Denny Zhou, et al. 2022. Challenging big-bench tasks and whether chain-of-thought can solve them. arXiv preprint arXiv:2210.09261.
- InternLM Team. 2023. InternIm: A multilingual language model with progressively enhanced capabilities.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023a. Llama: Open and efficient foundation language models. arXiv preprint arXiv:2302.13971.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023b. Llama 2: Open foundation and finetuned chat models. arXiv preprint arXiv:2307.09288.
- Bin Wang, Zhengyuan Liu, Xin Huang, Fangkai Jiao, Yang Ding, Ai Ti Aw, and Nancy F Chen. 2023a. Seaeval for multilingual foundation models: From cross-lingual alignment to cultural reasoning. arXiv preprint arXiv:2309.04766.

Lei Wang, Wanyu Xu, Yihuai Lan, Zhiqiang Hu, Yunshi Lan, Roy Ka-Wei Lee, and Ee-Peng Lim. 2023b. Planand-solve prompting: Improving zero-shot chain-ofthought reasoning by large language models. arXiv *preprint arXiv:2305.04091*.

733

734

735

736

737

738

739

740

741

742

743

744

745

746

747

748

750

751

752

753

754

755

756

757

758

759

760

761

762

763

764

765

766

767

769

770

772

773

774

775

778

779

782

- Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2022. Self-consistency improves chain of thought reasoning in language models. arXiv preprint arXiv:2203.11171.
- 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.
- Liang Xu, Hai Hu, Xuanwei Zhang, Lu Li, Chenjie Cao, Yudong Li, Yechen Xu, Kai Sun, Dian Yu, Cong Yu, et al. 2020. Clue: A chinese language understanding evaluation benchmark. arXiv preprint arXiv:2004.05986.
- Aiyuan Yang, Bin Xiao, Bingning Wang, Borong Zhang, Ce Bian, Chao Yin, Chenxu Lv, Da Pan, Dian Wang, Dong Yan, et al. 2023. Baichuan 2: Open large-scale language models. arXiv preprint arXiv:2309.10305.
- Jifan Yu, Xiaozhi Wang, Shangqing Tu, Shulin Cao, Daniel Zhang-Li, Xin Lv, Hao Peng, Zijun Yao, Xiaohan Zhang, Hanming Li, et al. 2023. Kola: Carefully benchmarking world knowledge of large language models. arXiv preprint arXiv:2306.09296.
- Xiang Zhang, Senyu Li, Bradley Hauer, Ning Shi, and Grzegorz Kondrak. 2023a. Don't trust chatgpt when your question is not in english: A study of multilingual abilities and types of llms. In Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, pages 7915–7927.
- Xiaotian Zhang, Chunyang Li, Yi Zong, Zhengyu Ying, Liang He, and Xipeng Qiu. 2023b. Evaluating the performance of large language models on gaokao benchmark. arXiv preprint arXiv:2305.12474.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. 2023. Judging llm-as-a-judge with mt-bench and chatbot arena. arXiv preprint arXiv:2306.05685.

Entity and Question Examples of the 7 Α **Chinese Commonsense Aspects**

Figure 5 shows the number of questions and partial entities of each Chinese commonsense aspect we propose, as well as corresponding question examples.

865

866

867

868

869

870

871

872

873

874

875

827

828

829

830

831

832

833

784

786

787

791

793

794

795

796

803

804

810

811

813

814

815

816

819

821

822

823

825

826

B Question Examples of the Reasoning Task in CHARM

Figure 6 shows the question examples of the 7 reasoning tasks in CHARM, including both Chinese and global domains.

C Question Examples of the Memorization Task in CHARM

Figure 7 shows the questions examples of the memorization tasks in CHARM.

D Examples of Prompt Strategies

Figure 8 shows the examples of the 5 prompt strategies.

E Detailed Evaluation Results of 19 LLMs with 5 Prompt Strageties on Reasoning Tasks

We conducted a detailed evaluation of 19 different LLMs using 5 distinct prompt strategies. Table 8 and Table 9 respectively display the performance of various prompt strategies on 7 reasoning tasks in the CHARM's Chinese commonsense domain and global commonsense domain.

F Performance of LLMs on Chinese Commonsense Knowledge Aspects

Table 10 displays the performance of LLMs in the 7 Chinese commonsense aspects. We only choose one representative prompt strategy: XLT for English LLMs and ZH-CoT for Chinese LLMs, which is based on our empirical conclusion in §5.2.

G Correlation of Memorization and Integrated Reasoning

G.1 Detailed Correlations of Memorization and Integrated Reasoning on the 4 *MRI* Tasks

The detailed performances of the 19 LLMs on the 4 *MRI* tasks are in Figure 9.

G.2 Correlation of Memorization and Integrated Reasoning throughout the LLM pretraining

We tested the intermediate checkpoint models of Baichuan2 and DeepSeek on the memorization and reasoning questions on the 4 *MRI* tasks. The results are shown in Figure 10.

With the increase in the number of tokens during the training process, the model's memorization ability quickly reaches a high level (in fact, there is no particularly obvious difference between the results of the first checkpoint and the final results). This is because the knowledge involved in our task setting is the most basic commonsense, and thus widely and abundantly exists in various Chinese training corpora.

However, the improvement in reasoning performance significantly lags behind memorization. This is because to complete a reasoning task in CHARM is actually a multi-step process, requiring memorization of relevant knowledge, understanding of the question, use of knowledge for reasoning, and answering according to the requirements of the question and the demonstration of few-shot examples, etc. If an error occurs in any step of the above complex process, the reasoning task will fail.

H Leaderboard of Memorization-Independent Reasoning

It is non-trivial to acquire and compare the memorization-independent reasoning abilities of the LLMs. Intuitively, we can filter the reasoning questions by only retaining those whose related memorization questions are all correctedly answered by every LLMs. This approach ensures that each LLM has memorized the commonsense knowledge necessary for the retained reasoning questions. However, when we apply this process to all the 19 LLMs, only 28 reasoning questions remain out of the original 500 in the *MRI* tasks, which is obviously insufficient in number and lacks diversity, thereby introducing a high degree of uncertainty due to randomness.

Therefore, we propose two slightly more complex methods, one called *Filtering Reasoning Questions based on Mono-LLM-Memorization* (FRMM), the other is *Memorization-Independent Battles among LLMs* (MIB).

H.1 Filtering Reasoning Questions based on Mono-LLM-Memorization (FRMM)

This method is relatively simple, but has some flaws to a certain extent. For each LLM, we select reasoning questions based on its performance in memorization tasks: only retaining reasoning questions for which all related memorization questions are answered correctly. It's clear that, after individual filtration, different LLMs will retain different reasoning questions, and even differ in the number of retained reasoning questions, as shown in the "# retained" column in Table 11.

876

877

878

879

887

891

893

896

897

901

902

903

904

905

906

907

908

909

910

911

912

913

914

915

916

917

918

919

Then, we calculate the accuracy of the retained reasoning questions for each LLM, and the results are shown in the "Retained Acc" column in Table 11. The LLMs are then ranked based on the accuracy, producing the leaderboard shown in the penultimate column of Table 6.

As mentioned above, while this method can reflect the memorization-independent reasoning abilities of LLMs to some extent, its drawback lies in that the denominator used in calculating the final ranking accuracy differs for different LLMs.

To overcome this, we have proposed the MIB method.

H.2 Memorization-Independent Battles among LLMs (MIB)

To overcome the shortcomings of the FRMM method, we refer to the pairwise battle method adopted in LLM evaluation (Zheng et al., 2023). By tallying each LLM's performance in a "round-robin" tournament of pairwise match-ups, we rank the performance of the LLMs.

Specifically, we select two LLMs at a time and filter the MRI task's reasoning questions based on the performance of these two LLMs in memorization tasks. We only retain the reasoning questions whose related memorization questions are correctly answered by both LLMs. In this way, the two LLMs are battled under fair conditions. We then calculate the accuracy of these two LLMs on the retained reasoning questions separately, and compute the difference in accuracy as the battle score between the two models.

As shown in Figure 11, the element E_{ij} represents the accuracy of LLM *i* minus the accuracy of LLM *j* during the battle between the two LLMs. For a total of 19 LLMs, we average each model's scores from the 18 battles they participated in as their final score, as shown in Table 12.

Finally, we rank the LLMs based on these scores to produce the leaderboard shown in the last column of Table 6.

I Memorization-Independent Reasoning Errors

LLMs can answer memorization questions correctly,
but they make mistakes when it comes to reasoning
problems composed of these knowledge points. Figure 12 shows the examples of three memorizationindependent reasoning errors of LLMs.

LLM	Prompt	AJ	TU	SqU	MMR	SpU	NLI	RC	Avg.
GPT-3.5-1106	Direct	34.67	33.0	46.0	48.0	60.0	62.0	52.5	48.02
	ZH-CoT	21.33	49.0	69.0	46.0	69.0	66.0	55.0	53.62
	Translate-EN	82.00	44.0	54.0	40.0	76.5	64.0	49.5	58.57
CDT 4 1106	XLT	82.00	39.0	65.0	42.0	80.5	61.0	50.5	60.00
GP1-4-1106	ZH-CoT	92.00 93.33	50.0 67.0	81.0 81.0	68.0	73.0 85.0	86.0 83.0	63.0 34.0	73.05
	EN-CoT	95.33	65.0	79.0	70.0	83.0	83.0	50.5	75.12
	Translate-EN XLT	93.33 95.33	43.0 60.0	85.0	42.0 74.0	75.5 86.0	81.0 77.0	53.5 62.5	65.48 77.12
LLaMA-2-7B	Direct	48.00	21.0	20.0	30.0	49.0	6.0	31.5	29.36
	ZH-CoT EN-CoT	50.67 48.00	34.0 34.0	23.0 20.0	22.0 4.0	51.0 50.5	34.0 40.0	29.5 26.5	34.88 31.86
	Translate-EN	11.33	32.0	20.0	20.0	39.0	40.0	20.0	26.05
LLaMA-2-13B	XLT	50.67	36.0	34.0	14.0	49.5	52.0	8.0	31.60
EEMMA 2 150	ZH-CoT	52.00	38.0	38.0	30.0	53.0	34.0	23.0	38.29
	EN-CoT Translate_EN	52.00 52.00	34.0 34.0	38.0	2.0	49.0 62.0	1.0	35.5	30.21
	XLT	54.67	33.0	38.0	30.0	58.0	47.0	38.0	42.67
LLaMA-2-70B	Direct	48.00	23.0	25.0	26.0	49.5	44.0	33.0	35.50
	EN-CoT	52.67	27.0	31.0	24.0	59.0	56.0	41.0	41.81
	Translate-EN	69.33	26.0	46.0	42.0	66.5	65.0	48.0	51.83
Vicuna-7B-v1.5-16k	Direct	52.00	25.0	30.0	16.0	14.5	19.0	21.5	25.43
	ZH-CoT	52.67	25.0	39.0	26.0	49.5	56.0	33.0	40.17
	EN-Col Translate-EN	52.67 68.00	28.0 25.0	26.0 31.0	18.0 30.0	40.5 60.0	55.0 57.0	40.5 33.0	37.24 43.43
	XLT	50.67	29.0	34.0	32.0	51.0	49.0	35.5	40.17
Vicuna-13B-v1.5-16k	Direct ZH-CoT	48.00 67.33	34.0 33.0	34.0 31.0	30.0 34.0	48.5	52.0 54.0	46.0 36.0	41.79 43.76
	EN-CoT	57.33	30.0	32.0	24.0	50.0	50.0	33.0	39.48
	Translate-EN	70.67	23.0	26.0	32.0	63.0	68.0	40.5	46.17
ChatGLM3-6B-32k	Direct	44.67	35.0	48.0	46.0	58.5	73.0	60.5	52.24
	ZH-CoT	65.33	40.0	59.0	38.0	77.0	72.0	37.5	55.55
	Translate-EN	53.55 62.67	24.0	39.0	40.0	53.5	75.0	32.5 46.0	52.40 48.02
D 1 4 7D	XLT	44.67	44.0	43.0	36.0	68.5	65.0	41.0	48.88
Baichuan2-7B	Direct ZH-CoT	44.00 72.67	31.0 41.0	37.0 48.0	24.0 38.0	59.0 72.0	35.0 53.0	56.0 49.5	40.86 53.45
	EN-CoT	53.33	36.0	44.0	30.0	69.0	53.0	41.0	46.62
	Translate-EN XLT	56.67 54.67	21.0 35.0	26.0 44.0	24.0 28.0	41.5 68.0	52.0 48.0	33.5 44.0	36.38 45.95
Baichuan2-13B	Direct	59.33	23.0	42.0	30.0	67.0	36.0	23.5	40.12
	ZH-CoT FN-CoT	83.33 74.67	40.0 40.0	48.0 50.0	46.0 34.0	72.5 68.0	66.0 64.0	51.5 42.5	58.19 53.31
	Translate-EN	61.33	38.0	40.0	32.0	58.5	49.0	36.0	44.98
InternI M2 7P	XLT	62.00	33.0	38.0	34.0	67.0	61.0	46.0	48.71
Internetwiz-7B	ZH-CoT	88.00	38.0	58.0	38.0	76.0	81.0	25.0	57.71
	EN-CoT Translata EN	76.67	42.0	59.0 45.0	38.0	73.0	78.0	38.5	57.88
	XLT	73.33	38.0	60.0	30.0	66.5	72.0	53.5	56.19
InternLM2-20B	Direct	26.67	42.0	61.0	50.0	39.5	54.0	46.5	45.67
	EN-CoT	72.67	40.0	48.0	44.0	67.0	68.0	25.0 25.0	59.50 51.81
	Translate-EN	76.00	34.0	54.0	36.0	53.5	71.0	53.0	53.93
Yi-6B	Direct	14.00	17.0	20.0	30.0	48.0	19.0	35.5	26.21
	ZH-CoT	72.67	32.0	47.0	32.0	75.0	50.0	42.0	50.10
	EN-Col Translate-EN	59.33 56.00	18.0 25.0	34.0 23.0	30.0 26.0	58.0 24.0	52.0 15.0	48.5 23.0	42.83
	XLT	54.67	36.0	35.0	28.0	68.5	56.0	43.0	45.88
Yi-34B	Direct ZH-CoT	86.00 94.00	28.0 55.0	85.0 89.0	56.0 76.0	70.0 88.5	51.0 72.0	68.0 51.5	63.43 75.14
	EN-CoT	90.00	42.0	78.0	66.0	84.5	67.0	50.0	68.21
	Translate-EN XLT	84.00 93.33	28.0 48.0	55.0 87.0	34.0 72.0	71.0 84.0	65.0 66.0	41.5 61.0	54.07 73.05
DeepSeek-7B	Direct	52.00	27.0	21.0	30.0	48.0	40.0	27.5	35.07
	ZH-CoT EN-CoT	79.33 72.67	34.0 33.0	50.0 33.0	50.0 24.0	79.5 73.0	57.0 47.0	31.5 35.5	54.48 45.45
	Translate-EN	65.33	18.0	28.0	36.0	59.0	72.0	40.5	45.55
DeenSeek 67P	XLT	55.33	32.0	39.0	36.0	51.0	37.0	35.0	40.76
Dupseek-0/D	ZH-CoT	93.33	+8.0 57.0	83.0	92.0	87.5	77.0	34.5	74.90
	EN-CoT	81.33	53.0	73.0	58.0	82.0	73.0	35.0	65.05
	XLT	91.33	43.0 59.0	80.0	58.0 66.0	87.5	76.0	40.0 54.5	58.69 73.48
Qwen-7B	Direct	48.67	28.0	41.0	50.0	60.5	56.0	55.5	48.52
	ZH-CoT EN-CoT	73.33 58.67	58.0 40.0	55.0 48.0	48.0 32.0	/1.0 68.5	57.0 58.0	49.5 43.0	55.98 49.74
	Translate-EN	63.33	23.0	31.0	26.0	60.0	54.0	45.0	43.19
Owen-14B	XLT Direct	63.33	26.0	47.0	40.0	63.0	50.0	50.5	48.55
×	ZH-CoT	88.00	54.0	77.0	60.0	82.5	66.0	55.0	68.93
	EN-CoT Translate_FN	85.33 75.33	48.0 29.0	68.0 45.0	56.0 22.0	77.5 60.5	76.0 63.0	53.0 46.5	66.26 48.76
	XLT	82.67	36.0	66.0	56.0	76.5	65.0	50.0	61.74
Qwen-72B	Direct	89.33	36.0	85.0	78.0	80.0	84.0	77.5	75.69
	EN-CoT	92.67	55.0	88.0	64.0	86.0	84.0 78.0	72.0	76.52
	Translate-EN	90.00	37.0	53.0	32.0	73.0	76.0	57.0	59.71
	XLI	89.33	50.0	86.0	/0.0	83.0	/5.0	58.5	/3.12

Table 8: Accuracy of reasoning tasks in the **Chinese** commonsense domain of CHARM.

LLM	Prompt	AJ	TU	SqU	MMR	SpU	NLI	RC	Avg.
GPT-3.5-1106	Direct	41.33	58.0	59.0	42.0	61.0	64.0	45.0	52.90
	ZH-CoT	25.33	89.0	90.0	28.0	77.5	73.0	48.5	61.62
	Translate-EN	39.33 88.67	85.0 86.0	80.0	42.0	83.5	65.0	58.0	72.74
CDT 4 1100	XLT	90.00	94.0	87.0	46.0	88.5	66.0	49.5	74.43
GP1-4-1106	ZH-CoT	90.67 92.67	100.0	92.0 90.0	70.0 34.0	88.0 88.0	78.0 76.0	74.0 50.5	82.24 75.88
	EN-CoT	95.33	97.0	97.0	52.0	89.5	70.0	61.5	80.33
	Translate-EN XLT	92.67 95.33	93.0 98.0	91.0 97.0	48.0 66.0	90.0	68.0 72.0	62.0 72.0	75.52 84.33
LLaMA-2-7B	Direct	43.33	20.0	20.0	28.0	51.5	18.0	20.5	28.76
	ZH-CoT EN-CoT	51.33 56.67	18.0 20.0	22.0 22.0	26.0 22.0	50.0 51.5	31.0 10.0	22.5 27.0	31.55 29.88
	Translate-EN	4.67	20.0	35.0	2.0	51.5	38.0	28.0	25.60
LLaMA-2-13B	XLT Direct	62.67 53.33	21.0	14.0	24.0	49.5	22.0	20.5	34.26
	ZH-CoT	54.67	15.0	35.0	14.0	51.5	32.0	27.0	32.74
	EN-CoT Translate-EN	52.67 34.67	19.0 19.0	35.0 35.0	16.0 20.0	51.5 57.0	35.0 37.0	38.0 28.0	35.31
	XLT	66.67	24.0	39.0	50.0	53.5	57.0	33.5	46.24
LLaMA-2-70B	Direct	48.67	33.0	33.0	30.0	50.0	58.0	20.5	39.02
	EN-CoT	46.00	82.0	53.0	34.0	51.5	64.0	48.5	54.14
	Translate-EN	84.67	76.0	58.0	52.0	71.0	64.0	57.5	66.17
Vicuna-7B-v1.5-16k	Direct	15.33	22.0	33.0	42.0	49.5	22.0	41.5	24.33
	ZH-CoT	52.00	50.0	50.0	10.0	50.5	53.0	31.0	42.36
	EN-Col Translate-EN	49.33 76.00	45.0 58.0	44.0 47.0	16.0 36.0	51.0 66.5	23.0 57.0	31.5 37.5	37.12 54.00
	XLT	45.33	64.0	37.0	26.0	58.5	52.0	32.5	45.05
Vicuna-13B-v1.5-16k	Direct ZH-CoT	55.33 71.33	58.0 71.0	39.0 49.0	28.0 14.0	48.5 53.0	59.0 62.0	30.0 33.0	45.40 50.48
	EN-CoT	66.00	82.0	42.0	38.0	65.0	55.0	40.5	55.50
	Translate-EN	84.00 72.67	66.0 74.0	60.0 41.0	66.0 50.0	71.0	60.0 61.0	42.0	64.14 57.52
ChatGLM3-6B-32k	Direct	44.00	33.0	57.0	42.0	63.0	80.0	38.0	51.00
	ZH-CoT	34.00	69.0	71.0	28.0	75.5	63.0	34.0	53.50
	Translate-EN	66.67	59.0	70.0	42.0	66.5	70.0	39.5	59.24
Delaharan 2 7D	XLT	52.67	58.0	70.0	46.0	66.0	66.0	42.5	57.31
Baichuanz-/B	ZH-CoT	46.00 55.33	20.0 65.0	47.0 54.0	8.0 26.0	58.0 60.5	55.0 59.0	38.0 29.0	49.83
	EN-CoT	44.00	64.0	49.0	20.0	58.5	56.0	31.5	46.14
	Translate-EN XLT	73.33 48.67	59.0 18.0	48.0 49.0	28.0 34.0	64.0 56.0	54.0 50.0	36.5 23.0	51.83 39.81
Baichuan2-13B	Direct	64.00	17.0	55.0	20.0	58.0	37.0	23.5	39.21
	ZH-CoT EN-CoT	77.33 78.67	74.0 70.0	58.0 55.0	40.0 30.0	71.0 57.0	61.0 66.0	39.0 37.5	60.05 56.31
	Translate-EN	73.33	68.0	51.0	36.0	61.5	61.0	42.0	56.12
InternI M2-7B	XLT Direct	70.67	75.0	49.0	42.0	69.5	61.0	31.0	56.88
InternEM2 7D	ZH-CoT	74.67	80.0	62.0	20.0	78.0	76.0	23.5	59.17
	EN-CoT Translate-EN	72.00 70.67	87.0 81.0	70.0 75.0	44.0 60.0	76.0 78.0	73.0 73.0	38.5 48.5	65.79 69.45
	XLT	66.00	87.0	72.0	52.0	76.5	66.0	43.5	66.14
InternLM2-20B	Direct ZH-CoT	81.33	54.0 83.0	78.0	50.0 14.0	63.5 74.5	46.0	48.0	60.12
	EN-CoT	73.33	83.0	63.0	14.0	75.0	73.0	26.5	58.26
	Translate-EN	82.00	84.0 84.0	89.0	40.0	76.0	68.0 70.0	46.5	69.36
Yi-6B	Direct	47.33	17.0	47.0	14.0	25.0	11.0	23.0	26.33
	ZH-CoT	79.33	63.0	43.0	14.0	70.5	57.0	33.5	51.48
	Translate-EN	72.67	57.0	62.0	32.0	69.0	37.0	35.5	52.17
	XLT	54.67	44.0	60.0	62.0	70.5	59.0	41.5	55.95
Y1-34B	Direct ZH-CoT	82.67 88.67	67.0 92.0	85.0 87.0	58.0 56.0	53.5 89.0	45.0 70.0	64.0 47.5	65.02 75.74
	EN-CoT	89.33	91.0	88.0	44.0	80.0	66.0	48.0	72.33
	Translate-EN XLT	78.00 88.00	85.0 88.0	83.0 86.0	48.0 70.0	76.5 93.5	64.0 60.0	54.5 58.0	69.86 77.64
DeepSeek-7B	Direct	47.33	24.0	35.0	14.0	22.0	41.0	17.5	28.69
	ZH-CoT EN-CoT	68.00 75.33	76.0 74.0	47.0 40.0	50.0 16.0	72.5 53.5	59.0 47.0	32.5 35.5	57.86 48.76
	Translate-EN	72.67	59.0	45.0	32.0	60.0	57.0	38.0	51.95
DeenSeek-67B	XLT	58.00	28.0	38.0	16.0	51.5	35.0	29.5	36.57
DeepSeek-07D	ZH-CoT	90.00	95.0	86.0	22.0	88.0	73.0	39.0	70.43
	EN-CoT Translata EN	61.33	96.0 87.0	76.0	30.0	90.5 81.0	71.0	35.0	65.69 72.02
	XLT	86.00	93.0	72.0	60.0	93.0	64.0	46.0	73.43
Qwen-7B	Direct	52.67	38.0	54.0	38.0	56.5	67.0	40.0	49.45
	EN-CoT	74.07	78.0 81.0	65.0	36.0	73.5	66.0	35.5	61.57
	Translate-EN	73.33	71.0	65.0	46.0	70.5	66.0	41.0	61.83
Qwen-14B	Direct	70.00	64.0 58.0	69.0 82.0	48.0	78.0	46.0	47.5	57.24 60.93
-	ZH-CoT	84.00	83.0	83.0	44.0	84.5	71.0	40.0	69.93
	EN-Col Translate-EN	86.67 86.67	82.0 72.0	81.0 85.0	44.0 48.0	79.5 78.0	66.0 64.0	42.5 48.5	68.81 68.88
	XLT	80.00	79.0	83.0	48.0	79.0	65.0	45.0	68.43
Qwen-72B	Direct ZH-CoT	88.00 94.00	63.0 92.0	85.0 93.0	56.0 64.0	83.5 93.0	78.0 71.0	65.5 63.5	74.14
	EN-CoT	90.00	92.0	86.0	60.0	92.5	66.0	58.0	77.79
	Translate-EN	91.33 92.67	87.0 70.0	89.0 91.0	54.0 66.0	81.5 91.5	63.0 66.0	64.0 50.5	75.69
	AL1	12.07	70.0	21.0	00.0	1.5	00.0	50.5	06.61

Table 9: Accuracy of reasoning tasks in the **global** commonsense domain of CHARM.

LLM	Prompt	Н	CA	LC	Ε	F	G	L	Avg.
GPT-3.5-1106	XLT	78.19	41.78	63.24	49.30	80.88	52.21	55.45	60.15
GPT-4-1106	XLT	91.49	58.22	86.76	73.24	86.27	71.68	67.27	76.42
LLaMA-2-7B	XLT	45.21	17.12	20.59	26.76	49.02	26.55	21.82	29.58
LLaMA-2-13B	XLT	50.00	38.36	39.71	36.62	57.84	46.02	30.91	42.78
LLaMA-2-70B	XLT	58.51	39.04	51.47	39.44	55.39	43.36	49.09	48.04
Vicuna-7B-v1.5-16k	XLT	48.40	31.51	33.82	38.03	50.98	44.25	32.73	39.96
Vicuna-13B-v1.5-16k	XLT	56.91	38.36	25.00	36.62	52.45	50.44	36.36	42.31
ChatGLM3-6B-32k	ZH-CoT	61.70	40.41	55.88	49.30	75.98	53.98	48.18	55.06
Baichuan2-7B	ZH-CoT	65.96	45.89	51.47	40.85	72.55	51.33	47.27	53.62
Baichuan2-13B	ZH-CoT	77.13	44.52	54.41	47.89	72.55	59.29	49.09	57.84
InternLM2-7B	ZH-CoT	76.60	33.56	61.76	49.30	75.49	49.56	45.45	55.96
InternLM2-20B	ZH-CoT	76.60	41.10	47.06	54.93	74.02	48.67	49.09	55.92
Yi-6B	ZH-CoT	62.23	43.84	54.41	38.03	75.00	44.25	36.36	50.59
Yi-34B	ZH-CoT	87.77	56.16	82.35	73.24	88.73	63.72	60.91	73.27
DeepSeek-7B	ZH-CoT	68.09	38.36	42.65	56.34	79.41	38.05	44.55	52.49
DeepSeek-67B	ZH-CoT	84.57	55.48	75.00	87.32	86.76	52.21	52.73	70.58
Qwen-7B	ZH-CoT	66.49	43.84	51.47	52.11	70.59	53.98	53.64	56.02
Qwen-14B	ZH-CoT	81.91	63.01	70.59	60.56	82.84	58.41	56.36	67.67
Qwen-72B	ZH-CoT	91.49	61.64	89.71	83.10	86.76	75.22	77.27	80.74

Table 10: Accuracy of reasoning questions on the 7 Chinese commonsense aspects of CHARM.

TIM	# Oniginal	Original Aga	# Datainad	Datained Ass
LLIVI	# Original	Original Acc.	# Retained	Retained Acc.
Qwen-72B	500	83.8	411	84.18
DeepSeek-67B	500	83.6	417	85.85
GPT-4	500	82.4	366	86.61
Yi-34B	500	82.2	372	88.71
Qwen-14B	500	76.2	333	81.98
InternLM2-20B	500	71.0	364	78.30
GPT-3.5	500	68.8	235	79.57
InternLM2-7B	500	68.2	328	80.79
DeepSeek-7B	500	67.4	351	75.21
Baichuan2-13B	500	66.6	329	76.29
Qwen-7B	500	62.8	332	65.36
Baichuan2-7B	500	62.6	354	64.97
ChatGLM3-6B	500	62.2	276	67.03
Yi-6B	500	61.4	311	69.77
LLaMA-2-70B	500	51.0	118	62.71
LLaMA-2-13B	500	49.2	90	58.89
Vicuna-13B	500	47.6	126	53.97
Vicuna-7B	500	44.6	92	46.74
LLaMA-2-7B	500	43.6	66	46.97

Table 11: Filtering Reasoning questions based on Mono-LLM-Memorization (FRMM) on the MRI tasks.

Commonsense Domain	Commonsense Aspect	Example of Entity	Example of Reasoning Question	# Question
Chinese commonsense Domain	History	朝代: 咸国、唐朝、宋朝 房支事件:赤壁之純、辛亥革命、北京奥运会 房支人物:李白、苏萩、成吉思汗 Dynasties: Warring States, Tang, Song Historical events: Battle of Red Cliffs, Xinhai Revolution, Beijing Olympics Historical figures: Li Bai, Su Shi, Genghis Khan	以下陈述是否包含时代错误,请遗择正确选项。一个接受了义务教育、具备基本常识的人会 如何选择? 刘邦在诸葛亮的辅佐下建立了汉朝。选项: (A) 是 (B) 否 Does the following statement contain historical errors? Please choose the correct option. How would a person who has received compulsory education and possesses basic knowledge choose? Liu Bang established the Han Dynasty with the assistance of Zhuge Liang. Option: (A) Yes (B) No	188
	Traditional Culture and Art	十二生肖:泉、牛、虎 艺木作品:《红枝梦》、《水浒传》、《三国演义》 发明: 指南针、大药、连纸木 Zodiac animals: Rat, Ox, Tiger Artistic works: "Dream of the Red Chamber", "All Men Are Brothers", "Romance of the Three Kingdoms" Invention: Compass, gunpowder, papermaking	小残在甲子年出生,他的表哥比他大5岁,那么他的表哥是在哪一年出生的? 透项: (A) こ卯 (B) 皮長(C) こ木 (D) 壬午 Xiaoqian was born in the year of Jiazi, and his cousin is 5 years older than him. So, in which year was his cousin born? Option: (A) Ji Mao (B) Geng Chen (C) Ji Wei (D) Ren Wu	146
	Daily life and Customs	生活方式: 高铁、网络购物、短视频 放舍: 设子、红烧肉、汤圆 节日: 端牛节、中秋节、重和节 Lifestyle: high-speed rail, online shopping, short videos Diet: dumplings, Braised pork belly, rice dumpling Festivals: Dragon Boat Festival, Mid-Autumn Festival, Double Ninth Festival	下列包決于約法程正确的是? 速項: (A) 得校子成改在手中、取适量校子指放在成的中央、探紧边缘、将校子成对析 (B) 将校子成这在手中、将校子成对析、 按紧边缘、取适量校子指放在成的中央 (C) 将校子方成在手中、再这量校子指放在成的中央、将校子皮对析、 按紧边缘 (D) 将校子方成在, 不能量效率, 新校子成放在中中、取适量校子指放在成的中央 (C) 将校子方成市, 探紧边缘、将校子方成在中中、取适量校子指放在成的中央 What is the correct process for making dumplings? Option: (A) Put the dumpling skin in your hand, take an appropriate amount of dumpling filling and place it in the center of the skin, place the edges tightly, and fold the dumpling skin in half (B) Put the dumpling skin in your hand, fold the dumpling skin in half (B) Put the dumpling skin in your hand, fold the dumpling skin in half (C) Put the dumpling skin in your hand, fold the dumpling skin in half (D) Put the dumpling skin in your hand, take an appropriate amount of dumpling filling and place it in the center of the skin, fold the dumpling skin in half, and pinch the edges tightly (D) Fold the dumpling skin in half, pinch the edges tightly, hold the dumpling skin in your hand, take an appropriate amount of dumpling align and place it in the center of the skin (D) fold the dumpling and place it in the center of the skin (D) Fold the dumpling skin in half, pinch the edges tightly, hold the dumpling skin in your hand, take an appropriate amount of dumpling filling and place it in the center of the skin (D) Fold the dumpling skin in half, pinch the edges tightly, hold the dumpling skin in your hand, take an appropriate amount of dumpling filling and place it in the center of the skin (D) Fold the dumpling skin in half, pinch the edges tightly, hold the dumpling skin in your hand, take an appropriate amount of dumpling filling and place it in the center of the skin (D) Fold the dumpling skin in place it in the center of the skin (D) Fold the dumpling filling and place it in the center of the skin (D) Fold the dumpling filling and place (D) Fold the dumpling filling place (D) Fold the dumpling filling place (D) Fold the dumplin	68
	Entertainment	电影: 《龙门飞甲》、《唐人街孩案》、《十画理庆》 音乐::《靖天》、《光释岁月》、《十年》 戏商: 京刹、得刹 Movies: "Dragon Gate Flying Armor", "Detective Chinatown", "Ambush from Ten Sides" Music: "Summ Day", "Glorious Years", "Ten Years" Traditional Chinese Opera: Peking Opera, Yu Opera	和这些电影《紅高梁》、《活着》、《大紅灯笼高高徒》、《英綽》有共同点的电影是? 造項: (A)《一个都不能少》(B)《让子弹飞》(C)《何飞正传》(D)《东邓西寿》 What movies have in common with these movies "Red Sorghum", "To Live", "Red Lantern Hanging High", and "Hero"? Option: (A) Not One Can Be Missing (B) Let the Bullets Fly (C) The True Story of Afei (D) Eastern Evil and Western Poison	71
	Public figures	公众人物: 刘翔、马龙、师市明 Public figures: Liu Xiang, Ma Long, Zou Shiming	下面的句子可信吗?"运动员刘翔三周半跳在冰面上画出了优美的弧线" 造項: (A) 可含 (B) 不可信 Is the following sentence credible? "Athlete Liu Xiang's three and a half jumps on the ice, drawing a beautiful arc" option: (A) Credible (B) Not credible	204
	Geography	城市:北京、上海、三亚 河流:长江、黄河、林江 省份:河北、河南、快西 Citiess Beijing, Shanghai, Sanya Rivers: Yangtze River, Yellow River, the Pearl River Provinces: Hebei, Henan, Shaanxi	i连句一: 郭尔多斯每天則市臺产煤炭 i运句二: 申国的河南布山东都是产煤大省 i荷问这两句话是什么关系? (A) 盧杏 (B) 矛盾 (C) 无关 Statement 1: Ordos and Datong are rich in coal mines Statement 2: Henan and Shandong in China are both major coal producing provinces May I ask what is the relationship between these two sentences? (A) Entailment (B) Contradiction (C) Unrelated	113
	Chinese language	成语: 生机物物、调虎高山、方大尚珊 诗词: "未天運叶无穷勞"、"无边落木黄膏下"、"千里共掉 弱" Idioms: vibrant, teasing tigers away from the mountains, dim lights Poems: "Endless blue lotus leaves reaching up to the sky", "Endless falling trees rustling down", "A thousand miles of shared beauty"	下列播绘一天时间变化的或话按照一天中时间的先后调序播序正确的是? 造項: (A) 晨光景微、泡田齐升、夕阳南下、呈月文晖 (B) 泡田序升、生月文晖、晨光景微、夕阳南下 (C) 差月文晖、晨光景微、池田东升、夕阳南下 (D) 夕阳南下、赵阳东升、夕阳南下 (D) 夕阳南下、赵阳东升、夕阳南下 (D) 夕阳南下、赵阳东升、夕阳南下 (A) The morning light is faint, the rising sun rises in the east, the setting sun sets in the west, and the stars and moon shine together (B) The rising sun rises in the east, the interplay of stars and moon shines, the morning light is faint, and the sunset sets in the west. The faint dawn, the rising sun in the east, and the setting sun in the west (D) The setting sun sets in the west, the rising sun rises in the east, the stars and moon shine together, and the morning light is faint.	110
Global commonsense domain	General knowledge worldwide	人物: 莎士比亚、贝多芬、拿成仑 生活: 做饭、穿衣、出行 地理: 四大洋、世界地图 Figures: Shakespeare, Beethoven, Napoleon General knowlege: Cooking, dressing, traveling Geography: The four major oceans, world map	以下除述是否包含时代错误,请选择正确遗项。一个接受了义务教育、具备基本常识的人会 如何选择? 貝多芬正在使用电子钢琴创作他的文响乐。选项: (A) 是 (B) 否 Does the following statement contain a chronological error? How would a person who has received compulsory education and has basic common sense choose? Beethoven is composing his symphony on an electronic piano. Options: (A) Yes (B) No	900

Figure 5: Entity and question examples of the commonsense aspects.

Task	Example of Reasoning Question (Chinese Domain)	Example of Reasoning Question (Global Domain)
Anachronisms Judgment	以下陈述是否包含时代错误,请选择正确选项。一个接受了义务教育、具 备基本常识的人会如何选择? 孙中山乘坐高铁从或昌前往南京。选项: (A) 是 (B) 否 Does the following statement contain historical errors? Please choose the correct option. How would a person who has received compulsory education and possesses basic knowledge choose? Sun Yat sen took the high-speed rail from Wuchang to Nanjing. Option: (A) Yes (B) No	 以下陈述是否包含时代错误,请选择正确选项。一个接受了义务教育、具备基本常识的人会如何选择? 古代玛雅文明的人们使用天文望远镜观测来制定 衣耕 日历。选项: (A) 是 (B) 否 Does the following statement contain historical errors? Please choose the correct option. How would a person who has received compulsory education and possesses basic knowledge choose? The people of the ancient Maya civilization used astronomical telescopes to observe and formulate agricultural calendars. Option: (A) Yes (B) No
Time Understanding	如果今天是小満,那么一个月后大约是什么节气? 遠項: (A)夏至(B)小署(C)大署(D) 立秋 If today is Xiaoman, what solar term will be approximately one month later? Options: (A) Summer Solstice (B) Minor Heat (C) Great Heat (D) Beginning of Autumn	请根据题目选择正确答案。今天是2007年的第一天。请问10天前的日期是多 少? 选项: (A)2006年12月22日 (B)2006年12月23日 (C)2007年02月24日 (D)2007年01月 31日 (E)1961年12月22日 (F)2006年12月21日 Please choose the correct answer based on the question. Today is the first day of 2007. May I ask what was the date 10 days ago? Option: (A) December 22, 2006 (B) December 23, 2006 (C) February 24, 2007 (D) January 31, 2007 (E) December 22, 1961 (F) December 21, 2006
Sequence Understanding	下列人物校时间先后顺序排序正确的是? 遠项: (A) 孙武、秦始皇、李白、袁世凯 (B) 秦始皇、袁世凯、孙武、李白 (C) 李白、孙武、秦始皇、袁世凯 (D) 孙武、秦始皇、袁世凯、李白 Which of the following characters is sorted correctly in chronological order? Options: (A) Sun Wu, Qin Shi Huang, Li Bai, Yuan Shikai (B) Qin Shi Huang, Yuan Shikai, Sun Wu, Li Bai (C) Li Bai, Sun Wu, Qin Shi Huang, Yuan Shikai (D) Sun Wu, Qin Shi Huang, Yuan Shikai (D) Sun Wu, Qin Shi Huang, Yuan Shikai	以下哪个列表按照人类发展历程排列正确? 遗项: (A) 现代社会, 铁器时代, 青铜时代, 石器时代 (B) 青铜时代, 石器时代, 铁器时代, 现代社会 (C) 石器时代, 青铜时代, 线器时代, 现代社会 (D) 铁器时代, 青铜时代, 現代社会, 石器时代 Which of the following lists is arranged correctly according to the history of human development? Option: (A) Modern society, Iron Age, Bronze Age, Stone Age (B) Bronze Age, Stone Age, Iron Age, Modern Society (C) Stone Age, Bronze Age, Iron Age, Modern Society (D) Iron Age, Bronze Age, Modern Society (D) Iron Age, Bronze Age, Modern Society, Stone Age
Movie and Music Recommendation	和这些歌曲《夜曲》、《本草纲目》、《听妈妈的话》、《七里香》有共 同点的歌曲是: 逸项: (A)《双节税》(B)《年少有为》(C)《浮夸》(D)《三人游》 The songs that share similarities with these songs "Nocturne", "Compendium of Materia Medica", "Listen to Mom's Words", and "Seven Miles Fragrance" are: options: (A) "Double knot Stick" (B) "Young and Promising" (C) "Exaggerate" (D) "Three person Tour"	寻找一部与《编嫱侠》、《变相怪杰》、《亡命夭涯》、《风月俏佳人》类 似的电影。选项: (A)《满城风雨》(B)《读情漩涡》(C)《狮子王》(D)《联社亚美利加》 Find a movie similar to "Batman", "The Mask", "The Fugitive", and "Pretty Woman". Options: (A) "The Front Page" (B) "Vertigo" (C) "The Lion King" (D) "Lamerica".
Sport Understanding	下面的句子可信吗? "运动员张怡宁大力扣篮" 选项: (A) 可信 (B) 不可信 Is the following sentence credible? "The athlete Zhang Yining dunks vigorously." Options: (A) Credible (B) Not credible	下面的句子可信吗? "科比·布莱恩特打板投篮得分" 选项: (A) 可信 (B) 不可信 Is the following sentence credible? Option for Kobe Bryant's rebounding and shooting scores: (A) Credible (B) Not credible
Natural Language Inference	 语句一:小明和家人在寒飯期间去三亚过年,发現酒店和旅游景点游客爆 満 満句二:三亚冬天的温度有20多度 请问这两句话是什么关系? (A) 蕴含 (B) 矛盾 (C) 元关 Statement 1: Xiao Ming and his family went to Sanya for the Chinese New Year during the winter vacation and found that hotels and tourist attractions were overcrowded. Statement 2: The temperature in Sanya during winter is over 20 degrees Celsius. What is the relationship between these two statements? (A) Entailment (B) Contradiction (C) Unrelated	语句一: 我们的朋友遍天下 语句二: 我们的朋友有很多。 请问这两句话是什么关系? (A) 蕴含 (B) 矛盾 (C) 无关 Sentence 1: We have friends all over the world. Sentence 2: We have many friends. What is the relationship between these two sentences? (A) Entailment (B) Contradiction (C) Unrelated
Reading Comprehension	在我国,中秋节是我国民间传统的五大节日之一,其核心的文化内涵是:"祝 愿社会和谐进步和家庭团圆幸福",但遗憾的是,如今商业化将中秋节演变成 为"月饼节",月饼越做越大,文化意义却越来越少.以下哪项是这段文字最有 可能支持的观点? (A) 传统文化不能作为经济资源加以利用 (B) 要挖掘和创新传统文化内涵,弘扬优秀民族文化 (C) 要充分挖掘传统节日蕴涵的巨大商机 (D) 商业活动应以传承民族文化为主要目的 In China, He Mid-Autumn Festival is one of the five traditional folk festivals in China. Its core cultural connotation is: "wish social harmony and progress and family reunion and happiness". Unfortunately, today's commercialization has turned the Mid-Autumn Festival into a "moon cake festival". The moon cakes are becoming bigger and bigger, but the cultural meaning is becoming less and less Which of the following is the viewpoint that this passage is most likely to support? (A) Traditional culture cannot be utilized as an economic resource. (B) We need to explore and innovate the connotation of traditional culture, promote excellent ethnic culture. (C) We need to fully tap into the huge business opportunities contained in traditional festivals. (D) Commercial activities should be primarily aimed at inheriting ethnic culture	研究表明,水污染的增长速度正在趋于平稳:今年造成的水污染量与去年几 乎相同。如果这种趋势持续下去,水污染问题将不再变得更加严重。 推理 是有问题的,因为它忽略了这样一种可能性 造项: (A)水污染的影响是累积的 (B)污染越未越严重 (C)水污染趋于平稳的趋势不会持续 (D)某些类型的水污染对于水生物没有明显影响 The study indicates that the growth rate of water pollution is stabilizing: the amount of water pollution caused this year is almost the same as last year. If this trend continues, the problem of water pollution will no longer become more serious. The reasoning is flawed because it overlooks the possibility that: (A) The impact of water pollution is cumulative (B) Pollution is becoming more severe (C) The trend of water pollution stabilizing will not continue (D) Some types of water pollution have no significant imp

Figure 6: Examples of the reasoning tasks in CHARM.

Task	Example of Memorization Question	Example of the Corresponding Reasoning Question	# Question
Anachronisms Judgment	华佗是中国哪个时期的人物? During which period in China was Hua Tuo a figure?	以下陈述是否包含时代错误, 请选择正确选项。一 个接受了义务教育、具备基本常识的人会如何选 择? 华佗使用麻沸散为病人手术。选项: (A) 是 (B) 否 Does the following statement contain an anachronism? Please choose the correct option. How would a person who has received compulsory education and possesses basic common sense choose? Hua Tuo used Mafeisan for surgery on patients. Options: (A) Yes (B) No	135
Time Understanding	清朝朝代对应的公元年份范围是什么? What is the range of AD years corresponding to the Qing Dynasty?	小刘在公元1912年出生,他的母亲比他大40岁,那 么他的母亲是在哪个朝代出生的?选项: (A) 清朝 (B)民国时期 (C) 元朝 (D) 明朝 Xiao Liu was born in AD 1912, and his mother was 40 years older than him. In which dynasty was his mother born? Options: (A) Qing Dynasty (B) Republic of China period (C) Yuan Dynasty (D) Ming Dynasty	83
Movie and Music Recommendation	《少年派的奇幻漂流》的主演有谁? Who are the main actors in "The Fantasy Drifting of the Youth School"?	和这些电影《鬼子来了》、《阳光灿烂的日子》、 《春桃》、《芙蓉镇》有共同点的电影是:选项: (A)《大佛普拉斯》 (B)《少年派的奇幻漂流》 (C)《让子弾飞》 (D)《大红灯笼高高挂》 The movie that has something in common with these films: "Devils on the Doorstep", "In the Heat of the Sun", "Spring Peach", and "Hibiscus Town" is: Options: (A) "The Great Buddha+" (B) "Life of Pi" (C) "Let the Bullets Fly" (D) "Raise the Red Lantern"	399
Sport Understanding	运动员王治郅从事哪项运动项目? Which sports does athlete Wang Zhizhi engage in?	下面的句子可信吗? "运动员王治邳水花压得很好"选 项: (A) 可信 (B) 不可信 Is the following sentence credible? "The athlete Wang Zhizhi is very good at splashing water." Options: (A) Credible (B) Not credible	127

Figure 7: Examples of the memorization tasks in CHARM.

LLM	Final Score
GPT-4	19.58
Qwen-72B	16.76
Yi-34B	16.62
DeepSeek-67B	15.63
Qwen-14B	10.19
InternLM2-7B	8.46
GPT-3.5	8.15
InternLM2-20B	7.40
DeepSeek-7B	3.42
Baichuan2-13B	3.04
ChatGLM3	-2.19
Baichuan2-7B	-2.26
Yi-6B	-3.91
Qwen-7B	-4.31
LLaMA-2-70B	-13.95
LLaMA-2-13B	-14.07
Vicuna-13B	-17.50
Vicuna-7B	-21.25
LLaMA-2-7B	-29.80

Table 12: Final results of the Memorization-Independent Battles among LLMs (MIB) on the MRI tasks.

Pr	ompt Strategy	Example
	Direct	Q:以下陈述是否包含时代错误,请选择正确选项。一个接受了义务教育、具备基本常识的人会如何选择?李白用钢笔写诗。选项:(A)是(B)否A:(A)
	ZH-Co T	Q:以下陈述是否包含时代错误,请选择正确选项。一个接受了义务教育、具备基本常识的人会如何选择?李白用钢笔写诗。选项:(A)是(B)否A:让我们一步一步来思考。这个陈述提到了"李白",他是中国唐朝时期的诗人。而陈述中提到的"钢笔"是现代设备,因此李白不可能使用钢笔写诗,该陈述包含时代错误。所以答案是(A)。
	EN-CoT	Q: 以下陈述是否包含时代错误,请选择正确选项。一个接受了义务教育、具备基本常识的人会如何选择? 李白用钢笔写诗。选项: (A) 是 (B) 否A: Let's think step by step. This statement mentions "Li Bai", a poet from the Tang Dynasty in China. The "pen" mentioned in the statement is a modern device, so it is impossible for Li Bai to write poetry with a pen. This statement contains errors from the times. So the answer is (A).
Т	`ranslate-EN	Q: Choose the correct option if the following statement contains an anachronism. How would a person with compulsory education and basic common sense choose?Li Bai wrote poetry with a fountain pen.Options:(A) Yes (B) No A: Let's think step by step.The statement mentions "Li Bai", a Chinese poet from the Tang Dynasty. The "fountain pen" mentioned in the statement is a modern device, so Li Bai could not have used a fountain pen to write his poems, and the statement contains an anachronism. The answer is (A).
	XLT	I want you to act as a commonsense reasoning expert for Chinese.Request:以下陈述是否包含时代错误,请选择正确选项。一个接受了义务教育、具备基本常识的人会如何选择? 李白用钢笔写诗。选项: (A) 是 (B) 否 You should retell the request in English. You should do the answer step by step to choose the right answer.You should step-by-step answer the request. You should tell me the answer in this format 'So the answer is'. Request: How would a typical person answer each of the following statements whether it contains an anachronism? Li Bai writes poetry with a pen. Option:(A) Yes (B) No Step-by-step answer: 1.This statement mentions "Li Bai", a poet from the Tang Dynasty in China. 2.The pen mentioned in the statement is a modern device. 3. so, it is impossible for Li Bai to write poetry with a pen. This statement contains errors from the times. So the answer is (A).

Figure 8: Examples of prompt strategies.



Figure 9: Accuracy of reasoning and memorization on the 4 MRI tasks.



Figure 10: Averaged accuracy of the intermediate checkpoint models throughout the LLM pretraining across the 4 *MRI* tasks.



Figure 11: Results of the Memorization-Independent Battles among LLMs (MIB) on the MRI tasks.



Figure 12: Examples of the 3 types of memorization-independent reasoning errors of LLMs