# TIMEBENCH: A Comprehensive Evaluation of Temporal Reasoning Abilities in Large Language Models

**Anonymous ACL submission** 

## Abstract

Understanding time is a pivotal aspect of human cognition, essential for fully appreciat-002 003 ing the complexities of the world. Previous studies typically focus on specific aspects of 005 time, lacking a comprehensive temporal reasoning benchmark. To address this, we pro-007 pose TIMEBENCH, a comprehensive hierarchical temporal reasoning benchmark that covers a broad spectrum of temporal reasoning phenomena. TIMEBENCH provides a thorough 011 evaluation for investigating the temporal reasoning capabilities of large language models. 012 We conduct extensive experiments on GPT-4, LLaMA2, and other popular LLMs under various settings. Our experimental results indicate a significant performance gap between the stateof-the-art LLMs and humans, highlighting that there is still a considerable distance to cover in temporal reasoning. LLMs exhibit capability discrepancies across different reasoning tasks. Furthermore, we thoroughly analyze the impact 022 of multiple aspects on temporal reasoning and emphasize the associated challenges. We aspire 024 for TIMEBENCH to serve as a comprehensive benchmark, fostering research in temporal reasoning. Code and data will be released.

## 1 Introduction

027

034

040

**Time flies over us, but leaves its shadow behind.** Time continually moves forward, threading through the fabric of people's lives. Understanding time is a crucial part of human comprehension of the world. Envision the blossoming of flowers, and you'll associate it with the arrival of spring. The ponder within it encompasses the intricate interplay of world knowledge, causality, and event temporal relationships. Temporal reasoning, in contrast to reasoning of a singular nature, comes with inherent complexity, encompassing implicit arithmetic, logical implications, and world knowledge. It is a form of integrated reasoning built upon foundational reasoning like mathematical and logical



Figure 1: A brief overview of human and LLMs' performance on TimeBench. Human scores are annotated.

reasoning (Cobbe et al., 2021; Mishra et al., 2022; Yu et al., 2020). Recently, large language models (LLM) have demonstrated remarkable performance in complex reasoning (Hendrycks et al., 2021; Srivastava et al., 2022; Brown et al., 2020; Chowdhery et al., 2023; OpenAI, 2023; Touvron et al., 2023), but their performance in intricate temporal reasoning still lacks a definitive conclusion.

Recent research for temporal reasoning typically focuses only on a few aspects, such as temporal commonsense or temporal question answering (Zhou et al., 2019; Chen et al., 2021; Dhingra et al., 2022; Wang and Zhao, 2023). Due to the inherent complexity of temporal reasoning, it is challenging to accurately measure models' temporal reasoning capabilities based on limited aspects.

To address this issue, we propose TIMEBENCH, a comprehensive and hierarchical temporal reasoning benchmark, which is aligned with intrinsic realworld scenarios. Specifically, drawing inspiration from the human cognitive process of transitioning from abstraction and concreteness to integration (Barsalou et al., 2018), we categorize temporal reasoning into three levels: symbolic temporal reasoning, commonsense temporal reasoning, and event temporal reasoning. These levels respectively represent understanding abstract time expression,

grasping concrete world knowledge, and integrating and applying this knowledge in real-world sce-070 narios. It comprises 10 tasks with 16 sub-tasks, 071 covering a broad spectrum of temporal reasoning phenomena. Besides, prior work typically featured only a singular task form, overly simplistic and insufficient in capturing the model's performance. In contrast, we incorporate four distinct task forms, offering a more realistic simulation of challenges.

To quantify the temporal reasoning capabilities of contemporary LLMs, we extensively assess 079 widely-used LLMs, including closed-source models such as ChatGPT (Ouyang et al., 2022) and GPT-4 (OpenAI, 2023), as well as open-source like LLaMA2 (Touvron et al., 2023), Vicuna-1.5 (Chiang et al., 2023), Mistral (Jiang et al., 2023), Baichuan2 (Yang et al., 2023), ChatGLM3 (Zeng et al., 2023) and FLAN-T5 (Chung et al., 2022). We conduct experiments under zero-shot and fewshot settings, combining commonly used reasoning 880 techniques, chain-of-thought prompting (Kojima et al., 2022; Wei et al., 2022). The experimental results suggest that GPT-4 outperforms other models, showcasing strong temporal reasoning capabilities, as shown in Figure 1. Nevertheless, there is still a considerable gap with humans. On the contrary, open-source models show inferior performance in temporal reasoning, attributed to shortcomings 096 in abstract time understanding, temporal relations modeling, and a lack of temporal commonsense. In addition, we also observe that chain-of-thought prompting does not consistently enhance model 100 performance. These findings indicate that there is still significant room for improvement in mod-102 els' temporal reasoning capabilities. Moreover, we conduct an in-depth analysis of the obstacles encountered by models in temporal reasoning.

084

101

103

104

105

106

107

108

109

110

111

112

113

114

115

116

117

118

We aspire for temporal reasoning to garner increased attention within the research community. Our contributions can be summarized as follows:

- We introduce TIMEBENCH, a comprehensive and hierarchical temporal reasoning benchmark to quantify the temporal reasoning ability of LLMs.
- We conduct extensive experiments on a large number of LLMs, revealing a significant gap between even sota LLM and humans, indicating substantial research opportunities in this field.
- By conducting a thorough analysis, we uncover the dilemmas that models face in temporal reasoning and identify potential solutions.

#### DATE ARITH

**Q**: What is the time 2 year and 4 month before Mar, 1755 A: Nov, 1752

#### TIMEX NLI

Premise: On 28th May 1967, I graduated. Hypothesis: Before 23rd October 1920, I graduated. A: Contradiction

Table 1: Examples of symbolic temporal reasoning

119

120

121

122

123

124

125

126

127

128

129

130

131

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

#### 2 **TIMEBENCH Benchmark**

#### 2.1 Benchmark Design Principal

TIMEBENCH focuses on a comprehensive evaluation of the temporal reasoning capabilities of large language models in challenging and complex scenarios. To achieve this goal, we summarize the difficulties and challenges faced in temporal reasoning, categorize them into three levels, and meticulously design tasks that better align with complex real-world scenarios.

Just as the human cognitive process unfolds from foundational cognition and conceptual understanding to practical reasoning, we delineate temporal reasoning into three hierarchical levels. Specifically, TIMEBENCH categorizes temporal reasoning into symbolic, commonsense and event temporal reasoning, covering 10 datasets with a total of 16 subtasks. (1) Symbolic Temporal Reasoning focuses on comprehension of fundamental abstract temporal expressions. (2) Temporal Commonsense Reasoning emphasizes the mastery of temporal principles, concepts and world knowledge. (3) Event Temporal Reasoning concentrates on modeling the temporal relationships of events within authentic scenarios. Furthermore, to better align with real-world scenarios, we employ diverse question formats.

### 2.2 Difficulties and Challenges

We outline the necessary capabilities and challenges encountered from a human cognition perspective during temporal reasoning, and language models confront similar obstacles. Detailed information on TIMEBENCH and challenges involved in each task is presented in Table 7.

TimeX Understanding Time expressions (TimeX) denote words or phrases that convey information about time and represent the simplest and most basic units of expressing time, such

#### MCTACO

C: Ransome looks after her as well as for young Fern Simon, who has declared her love for him. O: How often do Ransome and Fern talk?

O: each century, once a day, once a century, every night

## TIMEDIAL

**Dialog:** ... Person1: Do you go to work by train every day Person2: Yes . I commute <MASK> a week by train...

O: five days, 25 days, a minute, six days

#### SITUATEDGEN

157

158

159

160

161

164

165

168

169

170

171

173

174

175

177

179

181

182

**Keywords**: *axis, one day, one month, Earth, Moon* **A**: <u>Earth</u> rotates on its <u>axis</u> once in <u>one day</u>. It takes <u>one month</u> for the <u>Moon</u> to rotate on its <u>axis</u>.

Table 2: Examples of commonsense temporal reasoning.

as *in August 2008*, which is fundamental in comprehending time.

Temporal Commonsense assesses the understanding of temporal world knowledge, including event order, event duration, event typical time, and event frequency, which is crucial for language models to comprehend daily scenarios.

**Event-Time Relations** evaluates the grounding ability to establish correspondence between events and time, enabling the model to comprehend the development and changes of events as they unfold over time.

**Event-Event Relations** not only involve eventtime grounding but also introduce multi-hop relative connections. Models with this capability can better handle temporal reasoning in complex scenarios involving multiple events.

**Implicit Temporal Reasoning** involves going beyond the surface of texts, engaging in deeper reasoning such as drawing upon temporal commonsense, identifying implicit temporal factors and discerning hidden temporal relationships among events. Implicit temporal reasoning is pivotal in complex real-world scenarios where events and time are intricately interwoven.

## 2.3 Symbolic Temporal Reasoning

To evaluate the language model's comprehension
of abstract time expressions, we utilize two symbolic reasoning tasks stripped of semantic content:
TimeX arithmetic and TimeX inference. Table 1
shows examples of symbolic temporal reasoning.

## TIMEQA

C: ... He worked in Utrecht for the firm of P Smits & de Wolf from 1864 to 1867 and then returned to ... Q: Where did Ludwig Mond work between Mar 1866 and Sep 1866? A: Utrecht

## MENATQA

C: ... After the French evacuated Egypt in 1801, Hurshid Pasha was named governor of Egypt in 1804.
Muhammad Ali had himself named governor of Egypt in May 1805 ...
Q: Which position did Hurshid Pasha hold from 1804 to 1806, if Hurshid Pasha tepped down as the governor of Egypt in 1808?
A: governor of Egypt

#### **TEMPREASON**

C: ... Peter Corke works for Queensland University of Technology from Jan, 2010 to Dec, 2022. Peter Corke works for Commonwealth Scientific from Jan, 1984 to Jan, 2009. ... Q: Which employer did Peter Corke work for before Queensland University of Technology? A: Commonwealth Scientific

Table 3: Examples of event temporal reasoning.

**TimeX Arithmetic (Tan et al., 2023)** assesses the model's grasp of abstract date calculation. When provided with a date, the model needs to accurately calculate the date a certain amount of time before or after the given date. 188

189

190

191

192

193

194

195

196

197

198

199

201

202

203

204

205

206

207

208

209

210

211

212

213

**TimeX NLI (Thukral et al., 2021)** focuses on the logical entailment relationships among abstract TimeX, including three aspects: order (s1), duration (s2), and duration with unit conversion (s3).

#### 2.4 Commonsense Temporal Reasoning

We measure the model's mastery of temporal common and world knowledge, along with its capacity for reasoning based on these insights. Table 2 presents examples of temporal commonsense reasoning in QA and generation forms.

MCTACO (Zhou et al., 2019) evaluates diverse commonsense knowledge from different aspects of events, including duration, frequency, order, stationary and typical event time.

**DurationQA (Virgo et al., 2022)** focuses specifically on temporal commonsense reasoning in the spectrum of event duration.

**TimeDial (Qin et al., 2021)** considers temporal commonsense reasoning in dialogue scenarios and involves various aspects of commonsense associated with duration, order, and world knowledge.

297

298

299

300

301

258

214 SituatedGen (Zhang and Wan, 2023) consid-215 ers generative commonsense reasoning in a con-216 strained text generation scenario. Given a set of 217 contrasting keywords, the model needs to choose 218 appropriate keywords for each sentence and gen-219 erate a pair of contrasting sentences that satisfy 220 temporal commonsense.

## 2.5 Event Temporal Reasoning

222

225

240

241

242

243

246

247

251

254

255

Event temporal reasoning assesses the model's understanding of relationships between events and time in real-world scenarios, as well as its ability to reasoning under certain temporal or event constraints. Examples are shown in Table 3.

TimeQA (Chen et al., 2021) requires the model
 to answer time-sensitive questions based on con text containing numerous time-involved facts. It is
 categorized into explicit reasoning and implicit rea soning based on time indicators (before, in, etc.).

MenatQA (Wei et al., 2023) introduces temporal factors to elicit implicit temporal reasoning, including time scope change, disruption of facts, and counterfactual questions, which provides a more in-depth assessment of implicit reasoning ability on event-time relations.

**TempReason (Tan et al., 2023)** removes irrelevant context and focuses on implicit temporal reasoning within structured facts, investigating the model's capability boundaries. It involves eventtime reasoning and event-event reasoning.

**TRACIE** (Zhou et al., 2021) evaluates the model's comprehension of temporal order between implicit events. The model needs to identify events implied in the context and then determine their chronological order.

#### 2.6 Task Format and Evaluation Metrics

We design four different task forms to evaluate the LLM's reasoning in various scenarios. Please refer to Appendix A.2 for specific task forms and corresponding evaluation metrics.

## 3 Methodology

We perform evaluations using the prompt-based approach, including standard prompting and chainof-thought prompting. Experiments are conducted under both zero-shot and few-shot settings. **Standard Prompting** We formulate specific instructions for each task. In the zero-shot setting, models follow the instructions to answer questions. In the few-shot setting, models are provided with several question-answer pairs as demonstrations and emulate those instances to answer questions.

$$prompt_{zs}^{io} = \{INST\}\{Q\}$$
(1)

$$prompt_{fs}^{io} = \{INST\}\{Q_1\}\{A_1\}..\{Q\}$$
 (2)

**Chain-of-Thought Prompting** The instructions of CoT are the same as standard prompting. In the zero-shot setting, following Zeroshot CoT (Kojima et al., 2022), we add a reasoning trigger *Let's think step by step* after questions to perform chainof-thought reasoning. In the few-shot setting, we manually annotate CoT demonstrations for each task to guide the model to reason step-by-step. The instructions and demonstrations can be found in the Appendix A.6.

$prompt_{zs}^{cot} = \{INST\}\{Q\}\{TRIG\}$	(3)
$prompt_{fs}^{cot} = \{INST\}\{Q_1\}\{R_1\}\{A_1\}\{Q\}$	(4)

# 4 Experimental Setup

## 4.1 Models

We evaluate several popular LLMs, including both open-source and closed-source models, with parameter sizes ranging from 6B to 175B. The complete list of models can be found in Appendix A.1.

## 4.2 Implementation Details

We access closed-source models through Azure OpenAI API 0613-version. For open-source models, we deploy them locally through FastAPI. We set the temperature to 0 for greedy decoding in all experiments. To improve answer extraction accuracy, we use a summarization trigger *Therefore, the answer is* to obtain final answers.

## **5** Experimental Results

## 5.1 Few-shot Results

Table 4 presents the experimental results under few-shot settings. GPT-4 achieves the best performance across three categories, while LLaMA2<sub>70b</sub> and GPT-3.5 rank in the second tier. However, there remains a substantial gap of 19.4% between the most powerful LLM and humans.

In symbolic temporal reasoning tasks, GPT-4 demonstrates exceptional performance. However,

		Syn	ıbolic	:	(	Comm	onser	ise			l	Event	Tem	poral				Over	rall	
Method	Ti1 <i>s1</i>	neXN s2	NLI s3	Arith	DQA	McT.	TiD.	SitGen	Tim Exp.	eQA Imp.	M Sco.	lenatQ <i>Ord</i> .	PA Ctf.	Ten L2	npR L3	TRACIE	Sym.	Comm.	Event	Avg.
Human	98.0	96.0	92.0	100.0	80.8	87.1	97.8	100.0	93.3	91.1	85.6	87.3	79.9	97.1	95.3	82.5	96.5	91.4	89.0	91.5
GPT-4 + FS CoT	85.3 <b>92.0</b>	73.3 <b>84.0</b>	53.3 <b>64.0</b>	100.0 100.0	<b>64.8</b> 55.1	<b>88.3</b> 72.3	<b>94.6</b> 93.4	<u>88.6</u> -	73.7 66.9	51.0 <b>52.8</b>	72.4 65.3	<b>54.8</b> 52.6	<b>28.7</b> 25.9	92.4 <b>96.9</b>	<b>95.9</b> 94.6	62.8 <b>66.4</b>	78.0 <b>85.0</b>	<b>84.1</b> 73.6	<b><u>66.5</u></b> <u>65.2</u>	73.7 72.1
GPT-3.5 + FS CoT	52.0 <b>51.6</b>	68.4 <b>71.8</b>	31.6 <b>36.6</b>	63.6 <u><b>84.4</b></u>	<b>67.7</b> 41.2	<b>71.2</b> 38.1	<b>76.4</b> 71.1	<u>79.1</u> -	66.1 <u>68.0</u>	<b>48.4</b> 47.0	<b>43.2</b> 42.5	<b>51.6</b> 41.7	17.9 <b>37.8</b>	84.7 <b>89.9</b>	<b>78.0</b> 76.6	<b>55.0</b> 50.2	53.9 <u>61.1</u>	<b>73.6</b> 50.1	55.6 <u>56.7</u>	<b>59.7</b> 56.6
LLaMA2 $^{\dagger}_{70b}$ + FS CoT	55.0 52.0	61.0 <b>73.0</b>	37.0 <u>39.0</u>	<b>82.0</b> 79.5	<b>67.4</b> 62.3	<b>85.3</b> 79.1	<b>82.7</b> 61.1	<u>74.9</u> -	<b><u>66.7</u></b> 64.3	<b>48.3</b> 43.0	<b>61.4</b> 57.7	42.5 <b>45.2</b>	33.8 <u>53.1</u>	85.2 <b>87.5</b>	<b>85.4</b> 81.6	61.0 <u>67.0</u>	58.8 <u>60.9</u>	77.6 67.5	60.5 <u>62.4</u>	<b><u>64.4</u></b> <u>63.0</u>
$\begin{array}{l} LLaMA2^{\dagger}_{13b} \\ + FS \ CoT \end{array}$	<b>50.0</b> 40.0	54.0 <b>61.0</b>	30.0 <b>37.0</b>	29.5 <b>52.0</b>	53.3 <b>59.3</b>	66.0 <b>68.8</b>	<b>55.6</b> 40.8	64.8 -	59.3 <b>59.4</b>	48.6 <b>49.1</b>	49.6 <u>58.4</u>	43.4 <b>43.8</b>	37.5 <b>44.1</b>	<b>78.7</b> 78.0	62.7 <b>68.2</b>	58.0 58.0	40.9 <b>47.5</b>	<b>59.9</b> 56.3	54.7 <b>57.4</b>	52.6 <b>54.5</b>
LLaMA2 <sup><math>\dagger</math></sup> <sub>7b</sub> + FS CoT	26.0 <b>37.0</b>	50.0 <b>52.0</b>	30.0 <b>36.0</b>	20.0 <b>25.5</b>	54.5 <b>56.9</b>	59.6 <b>67.0</b>	<b>45.2</b> 41.9	62.4	<b>54.4</b> 45.6	<b>45.3</b> 36.1	49.8 <b>50.9</b>	<b>41.9</b> 38.0	35.8 <u>57.3</u>	<b>64.0</b> 59.7	53.3 <b>57.7</b>	49.0 <b>50.0</b>	31.5 <b>37.6</b>	<b>55.4</b> 55.3	49.2 <b>49.4</b>	46.3 <b>47.4</b>
Baichuan $2^{\dagger}_{13b}$ + FS CoT	38.0 <b>50.0</b>	48.0 <b>56.0</b>	33.0 <b>34.0</b>	42.5 <b>47.0</b>	54.8 <b>62.0</b>	<b>73.0</b> 69.3	<b>45.7</b> 43.8	64.9 -	<b>59.4</b> 58.2	<b>54.2</b> 49.6	<b>52.7</b> 49.8	38.0 <b>40.1</b>	21.4 <b>45.6</b>	77.3 <b>81.3</b>	63.5 <b>65.6</b>	54.0 <u>60.0</u>	40.4 <b>46.8</b>	<b>59.6</b> 58.4	52.6 <b>56.3</b>	51.3 <b>54.2</b>
Baichuan $2^{\dagger}_{7b}$ + FS CoT	27.0 <b>30.0</b>	<b>66.0</b> 56.0	<b>41.0</b> 34.0	32.5 <b>34.0</b>	<b>59.8</b> 57.0	69.4 <b>69.5</b>	34.3 <b>44.5</b>	59.8 -	<b>53.8</b> 51.2	50.2 40.7	<b>49.6</b> 46.4	<b>38.5</b> 32.6	22.9 <b>46.3</b>	<b>65.9</b> 61.5	51.0 <b>64.1</b>	<b>55.0</b> 53.0	<b>41.6</b> 38.5	55.8 <b>57.0</b>	48.4 <b>49.5</b>	<b>48.5</b> 48.1
Mistral <sup>†</sup> + FS CoT	48.0 <u>57.0</u>	53.0 <b>63.0</b>	<b>38.0</b> 35.0	41.0 <b>54.0</b>	61.8 61.8	<b>76.2</b> 45.7	<b>61.8</b> 57.3	58.3	55.9 <b>60.4</b>	45.3 <b>46.2</b>	49.4 57.2	47.8 <b>47.9</b>	<b>45.5</b> 33.2	<b>76.7</b> 65.9	<b>74.8</b> 67.9	53.0 <b>57.0</b>	45.0 <b>52.3</b>	<b>64.5</b> 54.9	<b>56.1</b> 54.5	<b>55.4</b> 54.0
ChatGLM3 <sup><math>\dagger</math></sup> + FS CoT	<b>48.0</b> 47.0	<b>70.0</b> 68.0	32.0 32.0	35.0 <b>46.0</b>	51.8 <b>53.9</b>	62.6 <b>64.3</b>	55.0 <b>56.5</b>	61.6	<b>57.2</b> 52.5	<b>26.3</b> 24.5	<b>35.4</b> 35.0	<b>41.5</b> 40.2	22.5 22.5	76.4 <b>79.4</b>	55.9 <b>60.3</b>	<b>58.0</b> 54.0	46.3 <b>48.3</b>	57.8 <b>58.2</b>	<b>46.7</b> 46.1	<b>49.3</b> 49.1

Table 4: Experimental results under **few-shot** settings (standard prompting by default). <sup>†</sup> denotes the base model without alignment. Best results in each group are **bold** and global top-3 results are <u>underlined</u>. Figure 6 provides a horizontal comparison of the performance of all models. Full results in Appendix A.5.

other models exhibit a significant decline in comparison to GPT-4. In commonsense temporal reasoning tasks, GPT4 lags behind humans by only 8.0%, indicating its powerful internal knowledge reservoir. With the model scale shrinking, its knowledge reservoir also decreases gradually, leading to a decline in performance. Notably, there is a significant gap of 25.2% between LLMs and humans in event temporal reasoning, which suggests that LLMs encounter major challenges in modeling intricate event-time relationships.

## 5.2 Zero-shot Results

303

305

306

310

311

312

313

Experimental results of alignment models under 314 zero-shot settings are shown in Table 5. In zero-315 shot settings, GPT-4 and GPT-3.5 rank first and second respectively, and they significantly outper-317 form all open-source models by a large margin. It is noteworthy that open-source models exhibit 319 a larger performance decline compared to closedsource models when transitioning from few-shot to 321 zero-shot scenarios. GPT, Baichuan2 and LLaMA2 suffer drops of 5.6%, 14.6% and 27.2% respec-323 tively. We attribute this performance decline to 324 the quality of alignment. Restricted by their lim-325 ited instruction-following capability, open-source models struggle to fully unleash their performance



Figure 2: Performance gap with and without CoT prompting. The results are averaged from GPT-4, GPT-3.5, Baichuan2<sub>13b</sub>, LLaMA2<sub>70b</sub> and Mistral<sub>7b</sub>.

solely through instructions. Therefore, few-shot prompting is a better approach for stimulating their temporal reasoning abilities. 328

329

330

331

332

333

334

335

336

337

338

339

#### 5.3 Chain-of-Thought in Temporal Reasoning

Previous research has found that chain-of-thought prompting can enhance the model's reasoning ability (Wei et al., 2022; Kojima et al., 2022). *Does CoT prompting bring consistent improvement in temporal reasoning?* Due to the diversity of temporal reasoning, the above question has not yet been definitively answered. To investigate this, we select several popular LLMs and analyze their performance affected by chain-of-thought prompting.

		Syn	nbolic	2	(	Comm	onser	ise			I	Event	Tem	poral				Over	rall	
Method	Tii s1	meXN s2	NLI s3	Arith	DQA	McT.	TiD.	SitGen	Tim Exp.	eQA Imp.	M Sco.	enatQ Ord.	QA Ctf.	Ten L2	npR L3	TRACIE	Sym.	Comm.	Event	Avg.
Human	98.0	96.0	92.0	100.0	80.8	87.1	97.8	100.0	93.3	91.1	85.6	87.3	79.9	97.1	95.3	82.5	96.5	91.4	89.0	91.5
GPT-4 + CoT	78.6 <b>80.0</b>	76.0 <u>76.0</u>	50.7 <u>60.0</u>	<b>98.0</b> 92.0	<b>59.2</b> 58.1	80.0 <b>82.6</b>	<b>91.1</b> 89.3	<u>59.3</u>	60.6 <b>61.3</b>	<b>46.5</b> 41.2	<b>57.0</b> 54.6	57.0 <u>59.6</u>	<b>23.1</b> 22.6	95.3 <u>97.0</u>	<b>95.0</b> 94.5	$\frac{\textbf{64.8}}{58.0}$	75.8 <u>77.0</u>	72.4 <b>76.7</b>	<b><u>62.4</u></b> 61.1	68.3 <u>68.5</u>
GPT-3.5 + CoT	<b>45.4</b> 33.6	<b>67.6</b> 64.8	31.2 33.6	<b>97.0</b> 71.0	<b>50.5</b> 23.2	<b>68.6</b> 45.1	<b>69.1</b> 67.0	<u>62.3</u>	<b>70.8</b> 64.4	<b>35.4</b> 35.1	<b>40.9</b> 39.7	<b>43.9</b> 42.9	22.9 <b>26.3</b>	<b>81.2</b> 57.6	<b>73.8</b> 68.1	<b>57.4</b> 52.0	<b>60.3</b> 50.8	<b>62.6</b> 45.1	<b>53.3</b> 48.3	<b>57.4</b> 48.3
LLaMA2 <sub>70b</sub> + CoT	$\frac{44.0}{30.0}$	<b>47.0</b> 66.0	<b>32.0</b> 28.0	78.5 53.5	<b>59.2</b> 57.3	<b>68.9</b> 67.1	57.0 <b>58.6</b>	25.0	<b>40.8</b> 31.4	<b>40.6</b> 19.5	<b>18.9</b> 12.2	<b>16.6</b> 12.7	12.0 20.8	<b>63.5</b> 37.5	<b>54.5</b> 40.5	48.0 <b>51.0</b>	<b>50.4</b> 44.4	52.5 <u>61.0</u>	<b>36.8</b> 28.2	<b>44.1</b> 39.1
LLaMA2 <sub>13b</sub> + CoT	30.0 <b>36.0</b>	49.0 <b>50.0</b>	34.0 <b>38.0</b>	<b>22.5</b> 6.0	38.5 <b>39.2</b>	40.6 <b>51.7</b>	35.4 <b>36.9</b>	57.9 -	<b>61.9</b> 58.7	30.5 <b>38.9</b>	<b>46.1</b> 40.9	<b>36.1</b> 32.5	26.9 <b>33.6</b>	53.1 <b>58.0</b>	<b>69.4</b> 68.4	<b>49.0</b> 47.0	<b>33.9</b> 32.5	<b>43.1</b> 42.6	46.6 <b>47.3</b>	<b>42.6</b> 42.4
LLaMA2 <sub>7b</sub> + CoT	39.0 <b>44.0</b>	<b>53.0</b> 50.0	30.0 <b>33.0</b>	<b>13.0</b> 5.0	<b>39.3</b> 35.0	<b>41.0</b> 40.0	<b>6.3</b> 1.7	24.5	49.0 <b>49.9</b>	29.0 <b>31.6</b>	26.8 <b>31.4</b>	21.1 24.5	16.0 <b>17.8</b>	<b>63.9</b> 56.9	47.9 <b>48.1</b>	<b>49.0</b> 46.0	<b>33.8</b> 33.0	<b>27.8</b> 25.6	37.8 <b>38.3</b>	34.3 34.3
Baichuan $2_{13b}$ + CoT	<b>41.0</b> 40.0	<b>61.0</b> 57.0	<b>37.0</b> 31.0	<b>12.5</b> 10.0	<b>52.0</b> 44.6	<b>63.4</b> 61.9	57.7 <b>58.1</b>	52.2	<b>55.4</b> 41.5	34.6 <b>40.9</b>	48.8 <b>52.0</b>	$\frac{44.3}{38.5}$	39.5 <b>43.2</b>	57.4 62.8	61.4 <b>64.3</b>	49.0 <b>55.0</b>	<b>37.9</b> 34.5	<b>56.3</b> 54.9	48.8 <b>49.8</b>	<b>48.0</b> 46.7
Baichuan $2_{7b}$ + CoT	35.0 <b>38.0</b>	<b>50.0</b> 43.0	<b>37.0</b> 32.0	<b>4.5</b> 1.0	<b>47.9</b> 37.9	55.3 <b>58.0</b>	<b>54.3</b> 44.2	42.0	41.5 53.5	34.7 <b>38.8</b>	35.2 <b>39.9</b>	31.2 33.2	20.4 <b>29.3</b>	<b>43.4</b> 41.2	<b>47.7</b> 47.2	<b>55.0</b> 54.0	<b>31.6</b> 28.5	<b>49.9</b> 46.7	38.6 <b>42.1</b>	<b>39.7</b> 39.4
Vicuna1.5 $_{13b}$ + CoT	35.0 <b>42.0</b>	50.0 <b>51.0</b>	36.0 <b>37.0</b>	<b>15.0</b> 3.0	<b>39.2</b> 29.8	<b>59.1</b> 50.0	<b>34.2</b> 33.7	51.8	<b>60.4</b> 56.9	<b>37.0</b> 36.4	$\frac{\textbf{46.8}}{38.2}$	37.4 <b>37.7</b>	<b>23.2</b> 20.4	42.1 <b>49.0</b>	43.6 <b>49.1</b>	46.0 <b>51.0</b>	<b>34.0</b> 33.3	<b>46.1</b> 37.8	42.1 <b>42.3</b>	<b>41.1</b> 39.0
Vicuna1.5 <sub>7b</sub> + CoT	<b>37.0</b> 36.0	<b>58.0</b> 50.0	<b>43.0</b> 36.0	<b>5.0</b> 1.5	<b>40.4</b> 39.4	<b>52.5</b> 49.2	32.0 <b>36.2</b>	47.8	<b>47.1</b> 40.9	18.5 <b>24.6</b>	<b>35.7</b> 26.2	25.7 <b>28.5</b>	17.3 <b>25.0</b>	<b>33.0</b> 27.7	<b>46.8</b> 40.3	54.0 54.0	<b>35.8</b> 30.9	<b>43.2</b> 41.6	<b>34.8</b> 33.4	<b>37.1</b> 34.4
FLANT5 <sub>11b</sub> + $CoT$	53.0 <b>56.0</b>	63.0 <b>66.0</b>	43.0 <b>45.0</b>	0.0 0.0	<b>52.0</b> 49.7	<b>65.0</b> 63.4	<b>47.7</b> 42.7	49.5	61.7 <b>64.4</b>	26.8 28.2	33.6 <b>41.6</b>	<b>52.2</b> 50.2	21.8 <b>30.6</b>	<b>87.9</b> 79.5	<b>83.9</b> 68.9	<b>64.0</b> 55.0	39.8 <b>41.8</b>	<b>53.6</b> 51.9	<b>54.0</b> 52.3	<b>50.3</b> 49.4
Mistral <sub>7b</sub> + CoT	<b>47.0</b> 38.0	50.0 <b>56.0</b>	<b>43.0</b> 35.0	<b>26.5</b> 16.5	<b>49.8</b> 36.6	<b>58.8</b> 49.3	<b>23.2</b> 19.3	<u>58.3</u>	28.2 <b>31.3</b>	21.4 22.4	<b>24.3</b> 21.1	22.3 <b>24.9</b>	21.7 <b>25.6</b>	<b>39.6</b> 34.0	<b>31.6</b> 31.2	51.0 <u>61.0</u>	<b>41.6</b> 36.4	<b>47.5</b> 35.1	30.0 <b>31.4</b>	<b>37.3</b> 33.5
ChatGLM3 <sub>6b</sub> + CoT	<b>38.0</b> 27.0	<b>50.0</b> 49.0	34.0 <b>37.0</b>	<b>2.0</b> 0.0	<b>34.1</b> 24.8	<b>43.6</b> 37.1	<b>56.7</b> 44.8	38.9	41.2 <b>41.7</b>	<b>31.7</b> 25.4	33.8 <b>34.6</b>	26.0 <b>28.1</b>	32.2 <b>41.2</b>	<b>57.0</b> 44.5	<b>54.0</b> 52.0	<b>50.0</b> 48.0	<b>31.0</b> 28.3	<b>43.3</b> 35.6	<b>40.7</b> 39.4	<b>39.0</b> 35.7

Table 5: Experimental results under **zero-shot** settings (standart prompting by default). All models are alignment models (-chat or -instruct). Best results in each group are **bold**, global top-3 results are <u>underlined</u>.

**Chain-of-thought reasoning is not consistently effective.** As illustrated in Figure 2, introducing zero-shot CoT prompting results in consistent declines, with an overall decrease of 7.4%. In the fewshot scenario, CoT prompting also fails to yield consistent improvements, varying depending on the task. There is a 10.8% improvement in symbolic reasoning, while a significant decline of 15.2% in commonsense reasoning. In event temporal reasoning, there is a slight improvement of 1.3%. Next, we will conduct a more detailed analysis of the impact of CoT on specific tasks.

341

342

344

347

348

351

Impact of CoT prompting across tasks. In order to explore the impact of CoT on various tasks thoroughly, we delve into the performance changes of each model across specific tasks within each category, as illustrated in Figure 3. In the zeroshot setting, open-source models achieve a slight improvement in event temporal reasoning with chain-of-thought prompting, while in other cases, they face performance degradation. While in the few-shot setting, almost all models exhibit significant improvement in symbolic temporal reasoning, with a concurrent prevalent decline in commonsense temporal reasoning. We attribute this to the knowledge sensitivity inherent in commonsense reasoning, where step-by-step reasoning cannot compensate for the lack of knowledge. In event temporal reasoning, improvements mainly stem from datasets involving implicit multi-step reasoning (MenatQA and TempReason), indicating that CoT is more effective for complex questions. In summary, zero-shot CoT consistently has a negative impact on temporal reasoning. While in fewshot scenario, CoT has a positive impact on symbolic and complex tasks, while negatively affecting knowledge-sensitive tasks. 364

365

366

367

368

369

370

371

372

373

374

375

376

377

378

379

380

381

382

383

384

## **6** Analysis and Discussion<sup>1</sup>

## 6.1 Effect of Scaling

We investigated how models scale affects temporal reasoning capability. The trend is illustrated in Figure 7. With the scale increasing, there is a notable performance enhancement. When the parameter size expands from 7B to 13B, LLaMA2

<sup>&</sup>lt;sup>1</sup>We give an error analysis in Appendix A.7.



Figure 3:  $\Delta$ Score between the chain-of-thought prompting and direct I-O prompting. **Top:** zero-shot setting, **Bottom:** few-shot setting, **Left:** variation in each task, **Right:** averaged variation in the symbolic, common-sense, event, and overall tasks.

Model	Order	Duration	Freq.	Stationarity	Typical	Avg.
GPT-4	76.4↓	92.8↑	83.3↑	71.4↓	54.5↓	77.5
GPT-3.5	50.5↑	39.8↓	55.2↑	48.4↑	28.7↓	43.5
Baichuan $2^{\dagger}_{13b}$	40.5↓	51.8↑	43.7↑	46.2↑	29.8↓	42.5
$LLaMA2^{\dagger}_{70b}$	65.2↑	72.1↑	66.3↑	36.3↓	52.7↓	63.0
$Mistral_{7b}^{\dagger}$	27.0↓	44.4↑	58.3↑	38.5↓	38.3↓	42.5

Table 6: Results in each temporal commonsense aspect under few-shot setting. Models with  $\dagger$  are base models. Red  $\downarrow$  and Green  $\uparrow$  represent the performance is lower or higher than its average performance. Metric is EM.

and Baichuan2 show improvements of 13.0% and 10.5%, respectively. Furthermore, when LLaMA scales up to 70B, the trend of performance improvement continues without stopping. The overall improvement follows a log-linearity with scale. There are no significant performance differences among LLaMA2, Baichuan2 and ChatGLM3 under similar parameter specifications, while Mistral shows impressive prowess, outperforming 13B models with nearly half the parameters.

## 6.2 Challenges in Temporal Reasoning

391

395

400

401

402

403

404

405

**LLMs underperform in (multi-hop) symbolic reasoning** Except for GPT-4, the performance of all other models in symbol temporal reasoning is unsatisfactory. Besides, a noticeable decrease is observed in duration-conversion task compared to other atomic tasks (25% in GPT-4 and 27% in LLaMA2<sub>70b</sub>). This is because the durationconversion task (s3) necessitates a two-step reasoning process. It first unifies time units, and subsequently engages in numerical comparison, In contrast, other atomic tasks (s1, s2 and arith) can be completed with a single inference. In summary, LLMs perform poorly in symbolic temporal reasoning and exhibit more pronounced declines when encountering multi-step reasoning. 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

446

447

448

449

450

451

452

453

454

455

Mastery of commonsense knowledge varies in LLMs We analyze models' performance across various commonsense aspects, as shown in Table 6. We regard the model's average performance in commonsense reasoning tasks as the baseline. If the model outperforms the baseline in a specific aspect, it suggests greater proficiency in this type of knowledge, and vice versa. The findings indicate that LLMs generally demonstrate good knowledge of event duration and frequency, while their comprehension of event order and typical events is relatively weaker. The uneven mastery of commonsense knowledge significantly affects the model's reasoning performance, especially when dealing with complex questions that involve multiple types of knowledge. Retrieval-augmented reasoning offers a promising avenue to alleviate the model's knowledge scarcity.

LLMs exhibit poor implicit temporal reasoning capabilities. When comparing the explicit and implicit event temporal reasoning (TimeQAexplicit versus others), we observe a significant decrease in the model's performance in implicit reasoning. Additionally, on TRACIE with numerous implied events, most models only surpass a random baseline (50.0). Even GPT-4 achieves a mere 66.4% accuracy, which implies that the LLM struggles with modeling implicit temporal relationships. We consider it helpful to explicitly model the temporal relationships between events and time expressions, for instance constructing timelines or temporal graphs.

LLMs are good factual reasoners rather than factual extractors When humans engage in temporal reasoning, it generally involves two steps: first, extracting time-fact pairs from the context, and then performing fact-based reasoning. TempReason provides extracted facts for conducting fact-based reasoning. By comparing the model's performance in context-based (TimeQA) and factbased (TempReason) reasoning, we identify the bottleneck in event temporal reasoning. LLMs excel in TempReason, signifying their strong capability in fact-based reasoning. However, their performance in context-based reasoning is significantly weaker



Figure 4: Performance difference between base and alignment models under few-shot setting. Baichuan2 and LLaMA2 are aligned with SFT and RLHF. Vicuna, Mistral and ChatGLM3 are aligned with only SFT.

than the former. This suggests that errors occur during the process of extracting facts from the context.
We attribute this gap to the model's deficiency in factual extraction capabilities. Thus, we consider LLMs to be strong factual reasoners rather than factual extractors in event temporal reasoning.

## 6.3 Alignment Impairs Temporal Reasoning

In the experiments mentioned earlier (Table 5), we observe a sharp decline in zero-shot performance of alignment models. To investigate whether alignment is responsible for the decline in temporal reasoning, we conduct experiments on alignment models under few-shot settings. The full experimental results can be found in Table 9. Figure 4 illustrates the overall performance decline after alignment. With the exception of Baichuan2, rest of the models are severely impaired, with catastrophic decline of up to 22%. By manually analyzing error cases, we conclude three reasons: (1) Alignment reduces the model's usability, causing it to tend towards refusal to answer when confronted with knowledge-sensitive questions. (2) Alignment damages the model's in-context learning capability, resulting in situations where the model deviates from the demonstrations. Furthermore, we believe the lack of temporal reasoning-related training data in alignment exacerbates this, resulting in disparities between different reasoning capabilities (e.g. mathematical reasoning v.s. temporal reasoning).

## 7 Related Work

## 7.1 Temporal Reasoning

Numerous research efforts address diverse challenges in temporal reasoning. Early research mainly relies on TimeML (Pustejovsky et al., 2003), focusing TimeX extraction and temporal relation extraction (Verhagen et al., 2007, 2010; UzZaman et al., 2013; Llorens et al., 2015; Miller et al., 2015; Mathur et al., 2021; Vashishtha et al., 2019). The advent of pre-trained language models (PLMs) has brought about commonsense reasoning as a tool to explore the world knowledge in models (Zhou et al., 2019; Qin et al., 2021; Dhingra et al., 2022). Recently, much attention has shifted towards event temporal reasoning (Chen et al., 2021; Tan et al., 2023; Zhu et al., 2023; Son and Oh, 2023; Chu et al., 2023b). Besides, Wang and Zhao (2023) introduces a unified form for assessing the temporal understanding capability of language models. 492

493

494

495

496

497

498

499

500

501

502

503

504

505

507

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

Distinguished from other datasets and benchmarks, TIMEBENCH is multispectral, and closely aligned with intricate real-world scenarios, offering a comprehensive and hierarchical evaluation of LLM's temporal reasoning abilities.

#### 7.2 Large-scaled Language Model

In recent years, there has been rapid progress in the research of large-scale language models (LLM) (Zhao et al., 2023). They exhibit outstanding performance across a multitude of tasks without the need for fine-tuning (Brown et al., 2020; Kojima et al., 2022). Furthermore, they have achieved astonishing results in complex reasoning tasks, such as mathematical reasoning (Cobbe et al., 2021; Mishra et al., 2022) and logical reasoning (Yu et al., 2020; Liu et al., 2023). Moreover, some studies suggest that the chain-of-thought prompting can further enhance the model's capabilities in complex reasoning scenarios (Wei et al., 2022; Kojima et al., 2022; Chu et al., 2023a; Zhang et al., 2023).

## 8 Conclusion

Temporal reasoning entails inherent diversity and complexity. The lack of a comprehensive benchmark makes it challenging to quantify LLMs' temporal reasoning capabilities. In this work, we present TIMEBENCH, a comprehensive and hierarchical benchmark for LLM temporal reasoning, tailored to mirror temporal reasoning in complex real-world scenarios. We conduct extensive experiments on state-of-the-art LLMs to investigate their temporal reasoning capabilities. Our findings indicate a substantial gap between state-of-the-art LLM and human performance, emphasizing the need for further research in this area. Moreover, we provide a meticulous analysis and discussion, outlining the current challenges that models face and suggesting potential directions for improvement.

479

480

482

483

484

485

486

487

488

489

490

491

456

457

458

459

460

461

462

463

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

595

596

# 541 Limitations

TimeBench is a comprehensive benchmark to quan-542 tify the temporal reasoning capabilities of LLMs. 543 While we have taken various factors into account, 544 there are a few limitations. Firstly, the research language of TimeBench is solely English, and we are considering the possibility of exploring appro-547 priate data to establish benchmarks for temporal reasoning abilities in other languages. Secondly, 549 in our evaluations, we only employ zero-shot and 550 few-shot methods. In future work, we plan to incorporate fine-tuning methods to provide a more 552 comprehensive analysis. 553

## References

554

555

556

557

558

559

561

562

564

565

568

571

572

573 574

577

580

581

582

584

585

589

590

593

- Joshua Ainslie, James Lee-Thorp, Michiel de Jong, Yury Zemlyanskiy, Federico Lebrón, and Sumit Sanghai. 2023. Gqa: Training generalized multi-query transformer models from multi-head checkpoints. *arXiv preprint arXiv:2305.13245*.
- Satanjeev Banerjee and Alon Lavie. 2005. Meteor: An automatic metric for mt evaluation with improved correlation with human judgments. In *Proceedings of the acl workshop on intrinsic and extrinsic evaluation measures for machine translation and/or summarization*, pages 65–72.
- Lawrence W Barsalou, Léo Dutriaux, and Christoph Scheepers. 2018. Moving beyond the distinction between concrete and abstract concepts. *Philosophical Transactions of the Royal Society B: Biological Sciences*, 373(1752):20170144.
- Iz Beltagy, Matthew E Peters, and Arman Cohan. 2020. Longformer: The long-document transformer. *arXiv* preprint arXiv:2004.05150.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners. In Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual.
- Wenhu Chen, Xinyi Wang, and William Yang Wang. 2021. A dataset for answering time-sensitive questions. In Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks 1, NeurIPS Datasets and Benchmarks 2021, December 2021, virtual.

- Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. 2023. Vicuna: An opensource chatbot impressing gpt-4 with 90%\* chatgpt quality.
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh, Kensen Shi, Sasha Tsvyashchenko, Joshua Maynez, Abhishek Rao, Parker Barnes, Yi Tay, Noam Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan Du, Ben Hutchinson, Reiner Pope, James Bradbury, Jacob Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin, Toju Duke, Anselm Levskaya, Sanjay Ghemawat, Sunipa Dev, Henryk Michalewski, Xavier Garcia, Vedant Misra, Kevin Robinson, Liam Fedus, Denny Zhou, Daphne Ippolito, David Luan, Hyeontaek Lim, Barret Zoph, Alexander Spiridonov, Ryan Sepassi, David Dohan, Shivani Agrawal, Mark Omernick, Andrew M. Dai, Thanumalayan Sankaranarayana Pillai, Marie Pellat, Aitor Lewkowycz, Erica Moreira, Rewon Child, Oleksandr Polozov, Katherine Lee, Zongwei Zhou, Xuezhi Wang, Brennan Saeta, Mark Diaz, Orhan Firat, Michele Catasta, Jason Wei, Kathy Meier-Hellstern, Douglas Eck, Jeff Dean, Slav Petrov, and Noah Fiedel. 2023. Palm: Scaling language modeling with pathways. J. Mach. Learn. Res., 24:240:1-240:113.
- Zheng Chu, Jingchang Chen, Qianglong Chen, Weijiang Yu, Tao He, Haotian Wang, Weihua Peng, Ming Liu, Bing Qin, and Ting Liu. 2023a. A survey of chain of thought reasoning: Advances, frontiers and future. *CoRR*, abs/2309.15402.
- Zheng Chu, Zekun Wang, Jiafeng Liang, Ming Liu, and Bing Qin. 2023b. MTGER: multi-view temporal graph enhanced temporal reasoning over timeinvolved document. *CoRR*, abs/2311.04816.
- Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, Albert Webson, Shixiang Shane Gu, Zhuyun Dai, Mirac Suzgun, Xinyun Chen, Aakanksha Chowdhery, Alex Castro-Ros, Marie Pellat, Kevin Robinson, Dasha Valter, Sharan Narang, Gaurav Mishra, Adams Yu, Vincent Zhao, Yanping Huang, Andrew Dai, Hongkun Yu, Slav Petrov, Ed H. Chi, Jeff Dean, Jacob Devlin, Adam Roberts, Denny Zhou, Quoc V. Le, and Jason Wei. 2022. Scaling instruction-finetuned language models.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. Training verifiers to solve math word problems. *CoRR*, abs/2110.14168.
- Bhuwan Dhingra, Jeremy R. Cole, Julian Martin Eisenschlos, Daniel Gillick, Jacob Eisenstein, and

William W. Cohen. 2022. Time-aware language models as temporal knowledge bases. *Trans. Assoc. Comput. Linguistics*, 10:257–273.

Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. Measuring massive multitask language understanding. In 9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021. OpenReview.net.

655

657

673

674

675

677

699

703

706

709

710

- Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de Las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Lélio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2023. Mistral 7b. CoRR, abs/2310.06825.
  - Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. Large language models are zero-shot reasoners. In *NeurIPS*.
- Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*, pages 74–81.
- Hanmeng Liu, Zhiyang Teng, Ruoxi Ning, Jian Liu, Qiji Zhou, and Yue Zhang. 2023. Glore: Evaluating logical reasoning of large language models. *CoRR*, abs/2310.09107.
- Hector Llorens, Nathanael Chambers, Naushad Uz-Zaman, Nasrin Mostafazadeh, James F. Allen, and James Pustejovsky. 2015. Semeval-2015 task 5: QA tempeval - evaluating temporal information understanding with question answering. In *Proceedings* of the 9th International Workshop on Semantic Evaluation, SemEval@NAACL-HLT 2015, Denver, Colorado, USA, June 4-5, 2015, pages 792–800. The Association for Computer Linguistics.
- Puneet Mathur, Rajiv Jain, Franck Dernoncourt, Vlad Morariu, Quan Hung Tran, and Dinesh Manocha. 2021. TIMERS: Document-level temporal relation extraction. In Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers), pages 524–533, Online. Association for Computational Linguistics.
- Timothy Miller, Steven Bethard, Dmitriy Dligach, Chen Lin, and Guergana Savova. 2015. Extracting time expressions from clinical text. In *Proceedings of BioNLP 15*, pages 81–91, Beijing, China. Association for Computational Linguistics.
- Swaroop Mishra, Matthew Finlayson, Pan Lu, Leonard Tang, Sean Welleck, Chitta Baral, Tanmay Rajpurohit, Oyvind Tafjord, Ashish Sabharwal, Peter Clark, and Ashwin Kalyan. 2022. LILA: A unified benchmark for mathematical reasoning. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, EMNLP 2022, Abu Dhabi,*

*United Arab Emirates, December 7-11, 2022*, pages 5807–5832. Association for Computational Linguistics.

- OpenAI. 2023. GPT-4 technical report. *CoRR*, abs/2303.08774.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F. Christiano, Jan Leike, and Ryan Lowe. 2022. Training language models to follow instructions with human feedback. In *NeurIPS*.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the* 40th annual meeting of the Association for Computational Linguistics, pages 311–318.
- James Pustejovsky, José M. Castaño, Robert Ingria, Roser Saurí, Robert J. Gaizauskas, Andrea Setzer, Graham Katz, and Dragomir R. Radev. 2003. Timeml: Robust specification of event and temporal expressions in text. In New Directions in Question Answering, Papers from 2003 AAAI Spring Symposium, Stanford University, Stanford, CA, USA, pages 28–34. AAAI Press.
- Lianhui Qin, Aditya Gupta, Shyam Upadhyay, Luheng He, Yejin Choi, and Manaal Faruqui. 2021. TIME-DIAL: temporal commonsense reasoning in dialog. In Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing, ACL/IJCNLP 2021, (Volume 1: Long Papers), Virtual Event, August 1-6, 2021, pages 7066–7076. Association for Computational Linguistics.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *J. Mach. Learn. Res.*, 21:140:1–140:67.
- Jungbin Son and Alice Oh. 2023. Time-aware representation learning for time-sensitive question answering. *CoRR*, abs/2310.12585.
- Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam Fisch, Adam R. Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, Agnieszka Kluska, Aitor Lewkowycz, Akshat Agarwal, Alethea Power, Alex Ray, Alex Warstadt, Alexander W. Kocurek, Ali Safaya, Ali Tazarv, Alice Xiang, Alicia Parrish, Allen Nie, Aman Hussain, Amanda Askell, Amanda Dsouza, Ameet Rahane, Anantharaman S. Iyer, Anders Andreassen, Andrea Santilli, Andreas Stuhlmüller, Andrew M. Dai, Andrew La, Andrew K. Lampinen, Andy Zou, Angela Jiang, Angelica Chen, Anh Vuong, Animesh Gupta, Ana Gottardi, Antonio Norelli, Anu Venkatesh, Arash Gholamidavoodi,

- 770
- 779 780

- 783
- 790
- 791 792
- 794 795 796
- 799
- 801
- 802 803
- 804 805

810 811

- 812 813
- 814 815
- 816 817
- 818 819

820 821

822 823

824

826

827

- Arfa Tabassum, Arul Menezes, Arun Kirubarajan, Asher Mullokandov, Ashish Sabharwal, Austin Herrick, Avia Efrat, Aykut Erdem, Ayla Karakas, and et al. 2022. Beyond the imitation game: Quantifying and extrapolating the capabilities of language models. CoRR, abs/2206.04615.
- Qingyu Tan, Hwee Tou Ng, and Lidong Bing. 2023. Towards benchmarking and improving the temporal reasoning capability of large language models. In Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2023, Toronto, Canada, July 9-14, 2023, pages 14820-14835. Association for Computational Linguistics.
- Shivin Thukral, Kunal Kukreja, and Christian Kavouras. 2021. Probing language models for understanding of temporal expressions. In Proceedings of the Fourth BlackboxNLP Workshop on Analyzing and Interpreting Neural Networks for NLP, BlackboxNLP@EMNLP 2021, Punta Cana, Dominican Republic, November 11, 2021, pages 396-406. Association for Computational Linguistics.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. Llama 2: Open foundation and finetuned chat models.
  - Naushad UzZaman, Hector Llorens, Leon Derczynski, James F. Allen, Marc Verhagen, and James Pustejovsky. 2013. Semeval-2013 task 1: Tempeval-3: Evaluating time expressions, events, and temporal relations. In Proceedings of the 7th International Workshop on Semantic Evaluation, SemEval@NAACL-HLT 2013, Atlanta, Georgia, USA, June 14-15, 2013, pages 1-9. The Association for Computer Linguistics.
- Siddharth Vashishtha, Benjamin Van Durme, and Aaron Steven White. 2019. Fine-grained temporal relation extraction. In Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics, pages 2906–2919, Florence, Italy. Association for Computational Linguistics.

Ramakrishna Vedantam, C. Lawrence Zitnick, and Devi Parikh. 2015. Cider: Consensus-based image description evaluation. In IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2015, Boston, MA, USA, June 7-12, 2015, pages 4566-4575. IEEE Computer Society.

829

830

831

832

833

834

835

836

837

838

839

840

841

842

843

844

845

846

847

848

849

850

851

852

853

854

855

856

857

858

859

860

861

862

863

865

866

867

868

869

870

871

872

873

874

875

876

877

878

879

880

881

883

884

885

- Marc Verhagen, Robert J. Gaizauskas, Frank Schilder, Mark Hepple, Graham Katz, and James Pustejovsky. 2007. Semeval-2007 task 15: Tempeval temporal relation identification. In Proceedings of the 4th International Workshop on Semantic Evaluations, SemEval@ACL 2007, Prague, Czech Republic, June 23-24, 2007, pages 75-80. The Association for Computer Linguistics.
- Marc Verhagen, Roser Saurí, Tommaso Caselli, and James Pustejovsky. 2010. Semeval-2010 task 13: Tempeval-2. In Proceedings of the 5th International Workshop on Semantic Evaluation, SemEval@ACL 2010, Uppsala University, Uppsala, Sweden, July 15-16, 2010, pages 57-62. The Association for Computer Linguistics.
- Felix Giovanni Virgo, Fei Cheng, and Sadao Kurohashi. 2022. Improving event duration question answering by leveraging existing temporal information extraction data. In Proceedings of the Thirteenth Language Resources and Evaluation Conference, LREC 2022, Marseille, France, 20-25 June 2022, pages 4451-4457. European Language Resources Association.
- Yuqing Wang and Yun Zhao. 2023. TRAM: benchmarking temporal reasoning for large language models. CoRR, abs/2310.00835.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc V. Le, and Denny Zhou. 2022. Chain-of-thought prompting elicits reasoning in large language models. In NeurIPS.
- Yifan Wei, Yisong Su, Huanhuan Ma, Xiaoyan Yu, Fangyu Lei, Yuanzhe Zhang, Jun Zhao, and Kang Liu. 2023. Menatqa: A new dataset for testing the temporal comprehension and reasoning abilities of large language models. CoRR, abs/2310.05157.
- Aiyuan Yang, Bin Xiao, Bingning Wang, Borong Zhang, Ce Bian, Chao Yin, Chenxu Lv, Da Pan, Dian Wang, Dong Yan, Fan Yang, Fei Deng, Feng Wang, Feng Liu, Guangwei Ai, Guosheng Dong, Haizhou Zhao, Hang Xu, Haoze Sun, Hongda Zhang, Hui Liu, Jiaming Ji, Jian Xie, Juntao Dai, Kun Fang, Lei Su, Liang Song, Lifeng Liu, Liyun Ru, Luyao Ma, Mang Wang, Mickel Liu, MingAn Lin, Nuolan Nie, Peidong Guo, Ruiyang Sun, Tao Zhang, Tianpeng Li, Tianyu Li, Wei Cheng, Weipeng Chen, Xiangrong Zeng, Xiaochuan Wang, Xiaoxi Chen, Xin Men, Xin Yu, Xuehai Pan, Yanjun Shen, Yiding Wang, Yiyu Li, Youxin Jiang, Yuchen Gao, Yupeng Zhang, Zenan Zhou, and Zhiying Wu. 2023. Baichuan 2: Open large-scale language models. CoRR, abs/2309.10305.

Weihao Yu, Zihang Jiang, Yanfei Dong, and Jiashi Feng. 2020. Reclor: A reading comprehension dataset requiring logical reasoning. In 8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020. OpenReview.net.

889

892

898

900

901

902

903

904 905

906

907

908

909

910

911

912

913

914

915

916

917 918

919

921

923

925 926

927

928

930

931

936

937

939

- Aohan Zeng, Xiao Liu, Zhengxiao Du, Zihan Wang, Hanyu Lai, Ming Ding, Zhuoyi Yang, Yifan Xu, Wendi Zheng, Xiao Xia, Weng Lam Tam, Zixuan Ma, Yufei Xue, Jidong Zhai, Wenguang Chen, Zhiyuan Liu, Peng Zhang, Yuxiao Dong, and Jie Tang. 2023.
  GLM-130B: an open bilingual pre-trained model. In The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023. OpenReview.net.
- Yunxiang Zhang and Xiaojun Wan. 2023. Situatedgen: Incorporating geographical and temporal contexts into generative commonsense reasoning. *CoRR*, abs/2306.12552.
- Zhuosheng Zhang, Yao Yao, Aston Zhang, Xiangru Tang, Xinbei Ma, Zhiwei He, Yiming Wang, Mark Gerstein, Rui Wang, Gongshen Liu, and Hai Zhao. 2023. Igniting language intelligence: The hitchhiker's guide from chain-of-thought reasoning to language agents. *CoRR*, abs/2311.11797.
- Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, Yifan Du, Chen Yang, Yushuo Chen, Zhipeng Chen, Jinhao Jiang, Ruiyang Ren, Yifan Li, Xinyu Tang, Zikang Liu, Peiyu Liu, Jian-Yun Nie, and Ji-Rong Wen. 2023. A survey of large language models.
- Ben Zhou, Daniel Khashabi, Qiang Ning, and Dan Roth.
  2019. "going on a vacation" takes longer than "going for a walk": A study of temporal commonsense understanding. 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 2019, Hong Kong, China, November 3-7, 2019, pages 3361–3367. Association for Computational Linguistics.
- Ben Zhou, Kyle Richardson, Qiang Ning, Tushar Khot, Ashish Sabharwal, and Dan Roth. 2021. Temporal reasoning on implicit events from distant supervision. In Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2021, Online, June 6-11, 2021, pages 1361–1371. Association for Computational Linguistics.
- Xinyu Zhu, Cheng Yang, Bei Chen, Siheng Li, Jian-Guang Lou, and Yujiu Yang. 2023. Question answering as programming for solving time-sensitive questions. *CoRR*, abs/2305.14221.





Symbolic Common Sense Event Temporal

Figure 5: Sunburst figure of category, task and subtask in TIMEBENCH. The degree of arc indicates the ratio of data.

## A Appendix

#### A.1 Models

941

942

943

947

948

951

960

961

**ChatGPT-3.5/GPT-4 (Ouyang et al., 2022; OpenAI, 2023)** ChatGPT is a chat model aligned through SFT and RLHF based on GPT-3 (Brown et al., 2020). GPT-4 is an upgraded version of Chat-GPT with enhanced reasoning capabilities, making it the most powerful LLM. Unless otherwise stated, ChatGPT refers to *gpt-3.5-turbo-0613* and GPT-4 refers to *gpt-4-0613*.

Llama2/Vicuna-1.5 (Touvron et al., 2023; Chiang et al., 2023) LLaMA2 is an open foundation model trained on 2T tokens with efficient groupedquery attention (Ainslie et al., 2023). LLaMA2chat is the official aligned model with SFT and RLHF, and Vicuna-1.5 is an unofficial aligned model with SFT only.

**Baichuan2 (Yang et al., 2023)** is an open foundation model pre-trained on 2.6T tokens and outperforms LLaMA2 on several benchmarks. Baichuan2-chat is the official aligned model with SFT and RLHF.

Mistral (Jiang et al., 2023) is a 7B open foundation model incorporating efficient grouped-query attention (Ainslie et al., 2023) and sliding windows attention (Beltagy et al., 2020). It achieves the strongest performance among models of its size, even surpassing LLaMA2-13B. Mistral-instruct is the officially aligned model with SFT only. **ChatGLM3 (Zeng et al., 2023)** is an opensource bilingual model for Chinese and English, exhibiting competitive performance in models with sizes under 10B.

**FLAN-T5 (Chung et al., 2022)** is an opensource instruction-following model built on top of T5 (Raffel et al., 2020) through instruction finetuning.

## A.2 Task Formats and Evaluation Metrics

**Multi-choice & Multi-answer Questions** MCQ tasks require models to select the only correct answer from the options. However, this task form has shortcuts and may not truly reflect the model's capability. To address this, we use the MCMAQ form, where the model needs to select all possible correct answers from the options. In our task, each question includes four options, with two of them being correct. Task: MCTACO, DurationQA, TimeDial.

**Natural Language Inference** is the task of determining the logical relationship between two pieces of text. Specifically, given a premise and a hypothesis, the model needs to determine whether the hypothesis can be inferred from the premise and output entailment, contradiction, or neutral. Our tasks focus on the entailment in temporal domains. Tasks: TimeX-NLI, TRACIE.

**Free-form Reading Comprehension** requires models to answer questions based on the provided context, and the answer is free-form without predefined restrictions. Tasks: TimeQA, MenatQA, TempReason, Date Calculation.

**Constrained Text Generation** refers to the task of generating text under certain constraints. Our task is keyword-constrained text generation, where the model takes keywords as input and outputs sentences that include those keywords. Task: SituatedGen.

MetricsWe adopt the evaluation metrics from10previous work. Accuracy is used for NLI and date10calculation tasks.MAMCQ tasks are measured10using option-based EM and F1.FRC tasks (w/o10date calc.) are evaluated with token-based EM and10F1.For CTG tasks, we use composite generation10metrics, with more details in Appendix A.4.10



Figure 6: Performance comparison between state-of-the-art LLMs. **Up**: GPT-4/3.5 and alignment models under zero-shot setting. **Down**: GPT-4/3.5 and base models under few-shot setting.



Figure 7: Pareto Frontier of the model size and overall time reasoning performance. The x-axis (model size) is shown in the log-scale. Results show a log-linearity between parameter size and performance.

1016

1018

1019

1020

1021

1022

1023

1024

1025

1026

1030

## A.3 Benchmark Details

We present the subtasks for each dataset in TIMEBENCH, including data quantities, task formats, challenges faced by each task, and the capabilities required by models, as shown in Table 7.

## A.4 Metrics for CTG

Following SituatedGen (Zhang and Wan, 2023), we use BLEU-4 (Papineni et al., 2002), ME-TEOR (Banerjee and Lavie, 2005), ROUGE-L (Lin, 2004), CIDEr (Vedantam et al., 2015), and MATCH (Zhang and Wan, 2023) scores to metric the results of Constrained Text Generation.<sup>2</sup>

The overall score is calculated as the sum of the above scores. We set the weight of CIDEr to 1/10 for balance when summing.

$$S = BLEU-4 + METEOR + ROUGE-L$$
  
+ CIDER/10 + MATCH

As the overall score S does not represent a percentile, we proceeded to normalize the models' scores to align with humans' relative performance levels.

1032

1033

1034

1035

1036

1038

1039

1041

1043

1044

1047

1048

1049

## A.5 Full Results

The overall score is derived from the average of all corresponding metrics. For brevity, we omit some F1 scores in the main tables in the main text. Please refer to Table 9 for the full experimental results, and detailed metrics for SituatedGen can be found in Table 8.

## A.6 Prompts

The prompt instructions are showcased in Figure 9. The demonstrations can be found from Figure 10 to Figure 18.

## A.7 Error Analysis

We manually analyze 100 errors for each subtask of three categories on GPT-4, GPT-3.5 and LLaMa2-base<sub>70b</sub>, as shown in Figure 8.

Symbolic Reasoning We categorize symbolic 1050 reasoning errors into five groups: (a) Expression: 1051 The model provides an incorrect time calculation 1052 expression. (b) Computation: The model provides the correct time calculation expression, but there is 1054 a calculation error. (c) Conversion: The model has an error in the conversion of time units. (d) Com-1056 parison: The model has an error when comparing two time-expressions (or intervals). (e) Combina-1058 tion: The model encountered errors in the combi-1059 nation of multiple above operations. LLMs exhibit 1060 numerous computation, conversion, and compari-1061 son errors, which suggests a substantial deficiency 1062

<sup>&</sup>lt;sup>2</sup>We utilize pycocoevalcap package to calucate BLEU-4, METEOR, ROUGE-L, CIDEr.



Figure 8: Error analysis for Symbolic, Commonsense, and Event Temporal. We select 100 test samples from each subtask for GPT-4, GPT-3.5 and LLaMa2-base<sub>70b</sub>.

in their understanding of fundamental temporal expressions. Besides, more errors occurred in combination questions, indicating that time-related multistep reasoning remains a major challenge for current models.

1063

1064

1065

1066

1067

1084

1085

1086

1087

1088

1089

1090

1091

1092

**Commonsense Reasoning** We categorize the er-1068 rors of commonsense reasoning into two groups: 1069 (a) No Answer: The model fails to provide a final 1070 answer. (b) Reasoning Error: The model encoun-1071 ters reasoning errors, which can be subdivided into five types of knowledge-related errors. We observe that GPT series models have a higher No Answer 1074 rate, while LLaMA is always able to provide an-1075 swers. We believe that this phenomenon is caused 1076 by two factors. On the one hand, the model lacks relevant commonsense knowledge to answer this 1078 question, and on the other hand, RLHF makes the 1079 model choose to refuse to answer when encoun-1080 tering questions beyond the knowledge boundary. 1081 Retrieval-augmented reasoning can alleviate the 1082 problem of knowledge scarcity to a certain extent. 1083

**Event Temporal Reasoning** We categorize the errors of commonsense reasoning into four groups: (a) *No Answer*: The model is unable to find the answer in the context. (b) *Reasoning Error*: The model encounters reasoning errors. (c) *Hallucination*: The model's prediction does not exist in the context, known as hallucination reasoning. (d) *Metric*: The model's prediction is correct, but the metric is limited by the evaluation criteria.

Dataset	Format	#	Challenges
Symbolic			
TimeX Arith	FRC	4,000	TimeX Arithmetic
TimeX NLI	NLI	6,965	TimeX Causality
- Order	-	2,213	order
- Duration	-	2,332	duration
- Conversion	-	2,420	duration + time unit conversion
Commonsense			
MCTACO	MCMAQ	852	Temporal Commonsense
TimeDial	MCMAQ	1,446	Temporal Commonsense
DurationQA	MCMAQ	687	Event Duration
SituatedGen	CTG	115	Temporal Commonsense
Event			
TimeQA	FRC	1,000	Context-based Reasoning
- Explicit	-	500	explicit, event-time reasoning
- Implicit	-	500	implicit, event-time reasoning
MenatQA	FRC	1,599	Implicit, Context-based Reasoning
- Order	-	400	event-time reasoning
- Scope	-	400	event-time reasoning
- Counterfactual	-	400	event-time reasoning
TempReason	FRC	1,876	Implicit, Fact-based Reasoning
- l2 (e2t)	-	839	event-time reasoning
- <i>l3 (e2e)</i>	-	1,037	event-event reasoning
TRACIE	NLI	500	Implicit, Implied Event-Event Reasoning
In total		19,000	

Table 7: The task formats and challenges in TIMEBENCH.

Method	BLEU-4	METEOR	ROUGE-L	CIDEr	MATCH	Overall	Norm
Human	39.9	40.4	56.3	397	98.1	274.4	100.0
GPT-4	8.23	31.27	28.84	38.45	90.41	162.59	59.25
+ FS	28.64	<b>38.99</b>	<b>55.69</b>	<b>298.64</b>	90.11	243.29	<b>88.66</b>
GPT-3.5	13.38	30.12	35.91	125.41	<b>78.76</b>	170.70	62.21
+ FS	27.24	<b>33.77</b>	<b>51.18</b>	<b>282.75</b>	76.54	217.01	<b>79.08</b>
LLaMA2 <sub>70b</sub>	5.15	13.62	15.83	22.07	31.79	68.60	25.00
+ FS	19.10	29.09	41.74	171.36	65.29	172.35	62.81
LLaMA2 <sub>13b</sub>	4.66	21.43	20.80	17.72	61.62	110.28	40.19
+ $FS$	15.15	27.49	37.55	138.13	64.94	158.93	57.92
LLaMA $2_{7b}$	2.77	13.46	14.69	14.34	34.83	67.18	24.48
+ FS	6.90	15.82	21.77	52.99	33.81	83.60	30.47
Baichuan $2_{13b}$	8.33	25.86	30.07	82.63	70.63	143.15	52.17
+ FS	15.79	30.23	40.96	169.14	71.01	174.91	63.74
Baichuan $2_{7b}$	5.17	21.99	23.73	44.80	59.85	115.22	41.99
+ FS	15.06	23.45	32.29	137.94	52.04	136.64	49.79
Vicuna1.5 $_{13b}$	7.73	26.35	29.15	69.16	71.91	142.06	51.77
+ FS	6.85	18.66	25.99	92.96	46.19	106.99	38.99
Vicuna $1.5_{7b}$	6.29	24.34	26.91	46.90	68.84	131.07	47.77
+ FS	20.71	30.19	45.20	203.20	67.58	184.00	67.05
FLAN-T5	16.20	24.43	29.38	95.17	56.38	135.91	49.53
+ FS	12.88	30.38	36.27	92.20	76.44	165.19	60.20
Mistral <sub>7b</sub>	5.82	22.89	24.19	44.03	63.74	121.03	44.11
+ FS	18.96	29.02	43.15	185.61	63.24	172.93	63.02
ChatGLM3 <sub>6b</sub>	6.56	21.11	21.96	41.48	53.02	106.80	38.92
+ FS	10.53	24.17	33.44	124.50	56.94	137.53	50.12
LLaMA2 $^{\dagger}_{70b}$	22.34	33.03	50.93	243.31	74.96	205.59	74.92
LLaMA2 $^{\dagger}_{13b}$	17.54	29.44	45.21	200.14	65.64	177.84	64.81
LLaMA2 $^{\dagger}_{7b}$	17.49	28.33	45.24	202.08	59.98	171.25	62.41
Baichuan $2^{\dagger}_{13b}$	17.86	29.75	44.28	198.83	66.35	178.12	64.91
Baichuan $2^{\dagger}_{7b}$	15.30	27.54	41.80	171.59	62.40	164.20	59.84
$Mistral_{7b}^{\dagger}$	14.54	27.39	41.72	168.89	59.42	159.96	58.30
ChatGLM3 <sup><math>\dagger6b</math></sup>	17.11	29.35	40.74	156.49	66.18	169.02	61.60

Table 8: Full results of SituatedGen. Aligned models are under zero-shot setting by default. The top-3 results are **bold**. Methods with † are base models without alignment, under few-shot setting. We consider human performance as 100 points and normalize models' results accordingly.

			_					_							Ē							-			=	
	namke		_					_							Eve											
	TimeXNLI s1 s2 s3	Date Arit Acc	$\left  \frac{Dura}{EM} \right $	tionQA F1	McTAC( EM F1	D TimeDia EM F.	al SitGe I Norm	i E-EN	Tin E-FI	іеQА Н-ЕМ	H-FI	S-EM	S-FI	Mena 0-EM	tQA 0-FI (	C-EM	C-FI L	-EM L	lempRe 2-F1 L	ason 3-EM L	.3-FI	RACIE Acc	Sym.	Comm.	Event	Avg.
	98.0 96.0 92.0	100.0	64.0	80.8	75.8 87.	1 97.8 97.	.8 100.0	0.68   0	93.3	87.0	91.1	82.0	85.6	84.0	87.3	76.0	79.9	96.0	97.1	94.0	95.3	82.5	96.5	91.4	89.0	91.5
	78.6         76.0         50.7           80.0         76.0         60.0           85.3         73.3         53.3           92.0         84.0         64.0	<b>98.0</b> 92.0 100.0 100.0	35.0 35.0 <b>51.0</b> 42.0	59.2 58.1 <b>64.8</b> 55.1	61.2 80. 67.0 82. <b>77.6 88.</b> 68.0 72.	0 72.0 91 6 65.0 89 <b>3 85.0 94</b> 3 79.0 93	.1 59.3 .3 - .6 88.6 .4 -	50.0 50.0 59.2 48.0	60.6 61.3 <b>73.7</b> 66.9	40.4 33.0 <b>40.0</b> <b>44.4</b>	46.5 41.2 51.0 <b>52.8</b>	44.4 43.4 <b>59.6</b> 48.5	57.0 54.6 <b>72.4</b> 65.3	49.0 <b>53.0</b> 48.0 44.0	57.0 <b>59.6</b> 54.8 52.6	22.0 20.0 22.0	23.1 22.6 <b>28.7</b> 25.9	91.0 86.0 91.0	95.3 <b>97.0</b> 92.4 96.9	94.0 93.0 93.0 93.0	95.0 94.5 <b>95.9</b> 94.6	64.8 58.0 62.8 <b>66.4</b>	75.8 77.0 78.0 <b>85.0</b>	72.4 76.7 <b>84.1</b> 73.6	62.4 61.1 <b>66.5</b> 65.2	68.3 68.5 68.5 73.7 72.1
	45.4 67.6 31.2 33.6 64.8 33.6 52.0 68.4 31.6 51.6 71.8 36.6	<b>97.0</b> 71.0 63.6 84.4	19.2 12.4 <b>42.8</b> 20.8	50.5 23.2 <b>67.7</b> 41.2	34.1 68. 28.1 45. <b>43.5 71.</b> 21.4 38.	6 39.2 69 1 34.6 67 2 47.8 <b>76</b> 1 <b>48.3</b> 71	.1 62.3 .0 - .4 <b>79.1</b> .1 -	<b>60.5</b> 52.5 53.8 56.5	<b>70.8</b> 64.4 66.1 68.0	29.5 29.0 <b>37.9</b> 37.5	35.4 35.1 <b>48.4</b> 47.0	36.5 35.8 37.8 <b>38.1</b>	40.9 39.7 <b>43.2</b> 42.5	37.5 38.5 <b>43.5</b> 37.5	43.9 42.9 <b>51.6</b> 41.7	21.0 24.0 16.0 <b>33.0</b>	22.9 26.3 37.8	73.6 32.0 <b>36.2</b>	81.2 57.6 84.7 <b>89.9</b>	61.8 54.2 <b>70.0</b> 68.0	73.8 68.1 <b>78.0</b> 76.6	<b>57.4</b> 52.0 55.0 50.2	60.3 50.8 53.9 <b>61.1</b>	62.6 45.1 <b>73.6</b> 50.1	53.3 48.3 55.6 <b>56.7</b>	57.4 48.3 <b>59.7</b> 56.6
	44.0       47.0       32.0         30.0 <b>66.0</b> 28.0         49.0       42.0       38.0         54.0       63.0 <b>40.0</b>	<b>78.5</b> 53.5 62.0 69.5	<b>12.7</b> 8.0 1.3 8.0	59.2 57.3 <b>61.2</b> 55.2	<b>23.0 68.</b> 21.0 67. 13.0 66 21.5 62.	<b>9 10.0</b> 57 1 9.0 <b>58</b> 5 6.0 56 1 6.0 56	.0 25.0 .6 .6 62.8	28.0 17.0 <b>41.0</b> 36.6	40.8 31.4 <b>51.1</b> 50.9	31.0 13.0 16.0 <b>34.0</b>	40.6 19.5 20.0 <b>42.4</b>	8.0 5.0 8.0 <b>28.0</b>	18.9 12.2 16.4 <b>38.6</b>	11.0 8.0 17.0 <b>19.0</b>	16.6 12.7 19.9 <b>29.3</b>	9.0 18.0 18.0 18.0	12.0 20.8 18.7 <b>21.9</b>	50.0 12.0 34.0 77.0	63.5 37.5 52.2 <b>83.1</b>	39.0 20.0 31.0 <b>65.0</b>	54.5 40.5 41.1 7 <b>4.7</b>	48.0 51.0 51.0 <b>57.0</b>	50.4 44.4 47.8 <b>56.6</b>	52.5 61.0 <b>61.8</b> 57.9	36.8 28.2 33.8 <b>49.7</b>	44.1 39.1 <b>53.2</b>
	30.0 49.0 34.0 36.0 50.0 38.0 <b>43.0 57.0 60.0</b> 37.0 50.0 38.0 37.0 50.0 37.0 50.0 37.0 55.0 50.0 37.0 55.0 50.0 50.0 55.0 50.0 55.0 50.0 55.0 50.0 55.0 50.0 55.0 50.0 55.0 50.0 55.0 50.0 55.0 50.0 55.0 50.0 55.0 50.0 55.0 50.0 55.0 50.0 55.0 50.0 55.0 50.0 55.0 50.0 55.0 50.0 55.0 50.0 55.0 55.0 50.0 55.0	<b>22.5</b> 6.0 33.0	4.0 7.3 9.0 <b>12.0</b>	38.5 39.2 46.8 <b>49.5</b>	8.5 40. 8.6 51. 8.0 66.	6 10.0 35 7 10.0 36 6 15.0 62 6 8.0 44	.4 57.9 .9 . . <b>3 40.2</b>	46.0 45.0 24.0 35.0	<b>61.9</b> 58.7 34.2 46.0	21.0 <b>30.0</b> 17.0 21.0	30.5 <b>38.9</b> 18.4 25.4	28.0 20.0 11.0 <b>34.0</b>	46.1 40.9 25.9 <b>46.7</b>	<b>23.0</b> 18.0 5.0 23.0	36.1 32.5 14.6 <b>36.5</b>	18.0 21.0 7.0	26.9 33.6 33.3 16.5	13.0 13.0 54.0 <b>72.0</b>	53.1 58.0 68.1 <b>80.8</b>	55.0 <b>56.0</b> 50.0 54.0	<b>69.4</b> 68.4 64.8 66.2	49.0 47.0 47.0 <b>50.0</b>	33.9 32.5 <b>45.1</b> 43.8	43.1 42.6 <b>54.0</b> 46.5	46.6 <b>47.3</b> 38.3 46.0	42.6 42.4 43.9 <b>45.5</b>
I	39.0         53.0         30.0           44.0         50.0         33.0           44.0         60.0         34.0           38.0         51.0 <b>36.0</b>	<b>13.0</b> 5.0 11.0 14.5	2.7 2.7 4.0 <b>11.0</b>	39.3 35.0 <b>62.8</b> 42.8	4.0 41. 4.5 40. 8.0 64. <b>25.0 65.</b>	0 1.0 6. 0 1.0 1. 7 8.0 40 6 13.0 53.	3 24.5 7 - .0 <b>30.5</b>	<b>37.0</b> 27.0 36.0 36.0	49.0 49.9 50.8 <b>53.5</b>	14.0 17.0 20.0 <b>21.0</b>	29.0 31.6 29.4 <b>34.1</b>	7.0 <b>11.0</b> 5.0 1.0	26.8 <b>31.4</b> 22.3 13.6	8.0 <b>10.0</b> 3.0	21.1 24.5 18.0 11.2	<b>9.0</b> 7.0 5.0	16.0 17.8 <b>17.9</b> 14.0	<b>18.0</b> 14.0 22.0	<b>63.9</b> 56.9 36.3 46.7	<b>32.0</b> <b>32.0</b> 23.0 21.0	47.9 <b>48.1</b> 44.3 42.3	49.0 46.0 <b>53.0</b> 51.0	33.8 33.0 <b>37.3</b> 34.9	27.8 25.6 49.5 <b>53.9</b>	37.8 <b>38.3</b> 34.0 33.3	34.3 34.3 3 <b>8.7</b> 37.8
	41.0 <b>61.0</b> 37.0 40.0 57.0 31.0 43.0 59.0 40.0 <b>45.0</b> 54.0 <b>48.0</b>	<b>12.5</b> 10.0 42.5 47.0	4.0 3.3 <b>24.7</b> 10.7	52.0 44.6 <b>62.1</b> 44.4	18.5 63. 20.0 61. <b>27.5 70.</b> 27.0 68.	4 15.0 57 9 13.0 58 <b>2 18.0 58</b> 8 15.0 55	7 52.2 .1 - .9 63.7	45.0 36.0 47.0 43.0	55.4 41.5 60.7 57.8	29.0 <b>36.0</b> 35.0 27.0	34.6 40.9 <b>45.7</b> 36.7	31.0 <b>39.0</b> 37.0 38.0	48.8 <b>52.0</b> 51.9 49.8	<b>34.0</b> 27.0 31.0 34.0	<b>44.3</b> 38.5 41.5 40.7	30.0 29.0 <b>33.0</b>	39.5 4 <b>3.2</b> 31.8 43.0	40.0 16.0 7 <b>3.0</b> 72.8	57.4 62.8 <b>81.1</b> 80.4	45.0 46.0 43.0	61.4 <b>64.3</b> 59.4 60.2	49.0 55.0 48.0 44.0	37.9 34.5 46.1 <b>48.5</b>	56.3 54.9 <b>63.7</b> 56.1	48.8 49.8 <b>52.5</b> 51.6	46.7 53.7 51.7
	35.0         50.0         37.0           38.0         43.0         32.0           40.0         50.0         36.0           41.0         50.0         36.0	<b>4.5</b> 1.0 23.5	4.0           5.3           28.7           13.0	47.9 37.9 <b>59.4</b> 45.7	10.5 55. 13.0 58. <b>26.5 66</b> . 17.5 58.	3 15.0 54 0 15.0 44 <b>9 17.0 53</b> 1 7.0 39	.3 42.0 .2 - .0 49.8	26.0 41.0 <b>45.0</b> 36.0	41.5 53.5 60.7 51.2	20.0 28.0 <b>30.0</b> 29.0	34.7 38.8 <b>42.1</b> <b>43.0</b>	20.0 29.0 27.0 <b>42.0</b>	35.2 39.9 37.8 <b>52.5</b>	19.0 23.0 25.0	31.2 33.2 35.7 <b>39.3</b>	6.0 18.0 10.0 <b>20.0</b>	20.4 29.3 <b>31.0</b>	22.0 21.0 40.0	43.4 41.2 57.4 <b>70.1</b>	29.0 37.0 <b>39.0</b>	47.7 47.2 53.0 <b>60.2</b>	<b>55.0</b> 54.0 51.0 49.0	31.6 28.5 36.5 <b>37.6</b>	<b>49.9</b> 46.7 57.3 47.7	38.6 42.1 44.8 <b>49.5</b>	39.7 39.4 45.8
	35.0         50.0         36.0           42.0         51.0         37.0           48.0         57.0         38.0           38.0 <b>59.0 39.0</b>	<b>15.0</b> 3.0 39.5 39.5	8.0 1.3 7.3 <b>10.7</b>	39.2 29.8 33.6 <b>37.4</b>	21.5 59. 11.5 50. <b>27.5 57.</b> 14.0 45.	1 7.0 34 0 7.0 33 0 13.0 40 8 12.0 41.	.2 <b>51.8</b> .7 - .3 39.0	43.0 44.0 45.0 <b>47.0</b>	60.4 56.9 58.3 <b>59.5</b>	29.0 31.0 23.0 <b>27.0</b>	37.0 36.4 25.9 <b>30.7</b>	38.0 16.0 38.0 <b>39.0</b>	46.8 38.2 42.6 <b>48.1</b>	22.0 25.0 26.0 <b>31.0</b>	37.4 37.7 <b>41.4</b> 35.9	17.0 13.0 18.0 <b>26.0</b>	23.2 20.4 <b>31.2</b>	14.0 31.0 51.0	42.1 49.0 61.8 <b>77.5</b>	13.0 29.0 <b>53.0</b>	43.6 49.1 <b>42.6</b> <b>65.5</b>	46.0 51.0 52.0	34.0 33.3 43.4 <b>43.9</b>	<b>46.1</b> 37.8 42.5 41.6	42.1 42.3 43.6 <b>50.1</b>	41.1 39.0 46.7
	37.0         58.0         43.0         36.0         36.0         36.0         36.0         36.0         36.0         37.0         37.0         37.0         37.0         35.0 <th< td=""><td><b>5.0</b> 1.5 8.5 8.0</td><td>1.3 1.3 <b>3.3</b> 2.7</td><td>40.4 39.4 <b>44.6</b> 37.2</td><td><b>9.5 52.</b> 8.5 49. 5.5 42. 10.0 47.</td><td><b>5</b> 6.0 32 2 9.0 36 1 7.0 36 5 <b>10.0 41.</b></td><td>.0 47.8 .2 - .8 <b>67.1</b> .3 -</td><td><b>35.0</b> 30.0 24.0 31.0</td><td><b>47.1</b> 40.9 31.9 39.9</td><td>11.0 <b>14.0</b> 12.0 13.0</td><td>18.5 <b>24.6</b> 14.9 16.6</td><td><b>20.0</b> 16.0 15.0</td><td><b>35.7</b> 26.2 21.8 26.7</td><td>15.0 14.0 <b>20.0</b> 15.0</td><td>25.7 2<b>8.5</b> 27.5 23.8</td><td>12.0 12.0 <b>17.0</b> 16.0</td><td>17.3 <b>25.0</b> 22.2 23.1</td><td>14.0 9.0 5<b>5.0</b></td><td>33.0 27.7 34.3 <b>66.3</b></td><td>14.0 7.0 6.0 <b>32.0</b></td><td>46.8 40.3 32.2 <b>48.1</b></td><td><b>54.0</b> <b>54.0</b> <b>54.0</b> 43.0</td><td>35.8 30.9 <b>36.4</b> 33.0</td><td>43.2 41.6 <b>47.7</b> 42.0</td><td>34.8 33.4 29.9 <b>35.9</b></td><td><b>37.1</b> 34.4 35.5 36.4</td></th<>	<b>5.0</b> 1.5 8.5 8.0	1.3 1.3 <b>3.3</b> 2.7	40.4 39.4 <b>44.6</b> 37.2	<b>9.5 52.</b> 8.5 49. 5.5 42. 10.0 47.	<b>5</b> 6.0 32 2 9.0 36 1 7.0 36 5 <b>10.0 41.</b>	.0 47.8 .2 - .8 <b>67.1</b> .3 -	<b>35.0</b> 30.0 24.0 31.0	<b>47.1</b> 40.9 31.9 39.9	11.0 <b>14.0</b> 12.0 13.0	18.5 <b>24.6</b> 14.9 16.6	<b>20.0</b> 16.0 15.0	<b>35.7</b> 26.2 21.8 26.7	15.0 14.0 <b>20.0</b> 15.0	25.7 2 <b>8.5</b> 27.5 23.8	12.0 12.0 <b>17.0</b> 16.0	17.3 <b>25.0</b> 22.2 23.1	14.0 9.0 5 <b>5.0</b>	33.0 27.7 34.3 <b>66.3</b>	14.0 7.0 6.0 <b>32.0</b>	46.8 40.3 32.2 <b>48.1</b>	<b>54.0</b> <b>54.0</b> <b>54.0</b> 43.0	35.8 30.9 <b>36.4</b> 33.0	43.2 41.6 <b>47.7</b> 42.0	34.8 33.4 29.9 <b>35.9</b>	<b>37.1</b> 34.4 35.5 36.4
	53.0       63.0       43.0         56.0       66.0       45.0         53.0       65.0       43.0         54.0       68.0       46.0	0.0 0.0 3.5 3.5	4.0 4.7 4.0 4.0	<b>52.0</b> 49.7 50.2 50.7	14.0 65. 14.5 63. <b>15.5 65.</b> 13.0 64.	0 13.0 47 4 <b>13.0</b> 42 <b>8</b> 9.0 35 0 11.0 <b>43</b>	.7 49.5 .7 - .0 <b>60.2</b>	56.0 57.6 55.0 54.0	61.7 64.4 64.3 59.4	<b>24.0</b> 23.9 23.0 19.0	26.8 28.2 25.0 21.2	31.0 <b>39.0</b> 31.0 34.0	33.6 <b>41.6</b> 33.6 36.6	<b>48.0</b> 46.0 47.0 45.0	<b>52.2</b> 50.2 50.6 49.2	20.0 <b>28.0</b> 21.0 20.0	21.8 8 <b>30.6</b> 22.5 21.7	84.0 73.0 82.0	87.9 79.5 87.0 <b>93.8</b>	<b>78.0</b> 57.0 <b>78.0</b> 72.0	83.9 68.9 <b>84.5</b> 79.7	64.0 55.0 65.0 <b>66.0</b>	39.8 41.8 <b>42.9</b>	<b>53.6</b> 51.9 52.8 52.8	54.0 52.3 <b>54.1</b> 53.5	50.3 50.5
	47.0 50.0 <b>43.0</b> 38.0 56.0 35.0 <b>51.0</b> 57.0 35.0	<b>26.5</b> 16.5 18.0	11.3 13.3 12.0	49.8 36.6 <b>55.8</b>	15.0 58. 14.5 49. <b>25.5 71.</b>	8 6.0 23 3 8.0 19. 0 13.0 52.	.2 58.3 .3 - <b>9 63.0</b>	7.0 13.0 <b>24.0</b>	28.2 31.3 <b>43.5</b>	5.0 11.0 <b>12.0</b>	21.4 22.4 23.9	4.0 8.0 4.0	24.3 21.1 21.3	7.0 <b>14.0</b> 7.0	22.3 <b>24.9</b> 21.2	4.0 <b>12.0</b> 7.0	21.7 <b>25.6</b> 23.0	2.0 5.0 23.0	39.6 34.0 48.9	1.0 4.0 23.0	31.6 31.2 44.9	51.0 <b>61.0</b> 57.0	<b>41.6</b> 36.4 40.3	47.5 35.1 <b>60.7</b>	30.0 31.4 35.5	33.5 33.5 <b>43.0</b>

MenatQA         TempReason         TRACIE         Syn.         Comm.         Event Ave.           0 <b>25.3 9.0 24.4</b> 9.0 <b>24.3</b> 9.0 <b>78.6 42.0 57.4</b> 500 <b>37.9 43.4 39.6 39.6 39.6 39.6 39.6 39.6 39.6 39.6 39.6 39.6 39.6 39.6 39.7 39.6 39.7 39.6 39.7 39.6 39.7 39.6 39.7 39.6 39.7 39.6 39.7 39.6 39.7 39.6 39.7 39.6 39.7 39.6 39.7 39.7 39.7 39.7 39.7 39.7 39.7 39.6 39.7</b>
0         48.3         52.0         61.4         55.0         53.1         78.0         85.5         76.0         81.6         61.0         58.8         77.6         60.5         64.4         63.0           0         43.0         50.0         57.7         37.0         45.2         44.0         53.1         82.0         78.5         76.0         81.6         61.0         55.8         62.4         63.0           0         48.6         41.0         49.6         38.0         44.1         70.0         78.7         49.0         67.5         56.3         57.4         53.5           0         49.1         49.0         57.3         80.0         58.4         49.0         77.5         56.3         57.4         57.4         54.5           0         45.1         70.0         78.0         57.0         57.3         49.0         77.5         56.3         57.4         56.3         57.4         57.3         56.3         57.4         56.3         57.4         56.3         57.4         56.3         57.4         56.3         57.4         56.3         57.4         56.3         57.4         56.3         57.4         56.3         57.4         56.3         57.
$ \begin{array}{[cccccccccccccccccccccccccccccccccccc$
$ \begin{array}{c ccccccccccccccccccccccccccccccccccc$
0         54.2         42.0         52.1         31.0         38.0         13.0         21.4         68.0         77.3         50.0         63.5         54.0         40.4         59.6         52.6         51.3           0         49.6         39.0         49.8         34.0         40.1         37.0         45.6         73.0         81.3         46.0         65.6         60.0         46.8         58.4         56.3         54.2           0         49.7         38.0         49.6         35.0         81.3         46.0         65.9         34.0         51.0         55.0         49.4         48.5         43.1           0         40.7         38.0         46.4         26.0         32.5         48.0         57.0         49.4         48.5         43.1           0         40.7         38.0         45.7         45.3         60.0         65.9         51.0         51.0         57.0         49.5         48.4         48.5           0         45.3         46.0         65.9         58.0         67.0         74.8         56.3         54.1         54.5         54.1         54.0         54.5         54.0         54.5         54.0         54.5
0         50.2         41.0         49.6         33.0         38.5         18.0         22.9         49.0         65.9         34.0         51.0         55.0         41.6         55.8         48.4         48.5           0         40.7         38.0         46.4         26.0         32.6         41.7         46.3         46.0         61.5         43.8         64.1         53.0         43.5         57.0         49.5         48.4         48.5           0         45.3         37.0         61.7         46.3         46.0         61.5         43.8         63.0         70.7         49.4         57.0         49.5         58.1           0         45.3         37.0         47.9         28.0         45.7         64.0         74.8         53.0         45.7         54.9         54.5         54.9         54.5         54.0         54.5         54.0         54.5         54.0         54.5         54.0         54.5         54.0         54.5         54.0         54.5         54.0         54.5         54.0         54.5         54.0         54.5         54.0         54.5         54.0         54.5         54.0         54.5         54.0         54.5         54.0 <t< td=""></t<>
0         45.3         37.0         45.5         68.0         76.7         64.0         74.8         53.0         45.0         64.5         56.1         55.4           0         46.2         48.0         57.2         37.0         47.9         24.0         33.2         60.0         65.9         58.0         67.9         57.0         54.9         54.5         54.0           0         46.2         38.0         41.5         22.0         22.5         67.0         65.9         58.0         67.9         54.9         54.5         54.0           0         26.3         30.0         35.4         38.0         41.5         22.0         22.5         67.0         76.4         35.0         58.0         46.3         57.3         54.9         54.5         54.0           0         26.3         30.0         35.4         38.0         41.5         22.0         22.5         67.0         76.4         35.0         54.6         36.3         54.5         54.0           0         24.5         36.0         37.0         55.9         58.0         61.3         37.0         46.3         37.8         46.7         49.3           0         24.5
0         26.3         30.0         35.4         38.0         41.5         22.0         22.5         67.0         76.4         35.0         55.9         58.0         46.3         57.8         46.7         49.3           0         24.5         30.0         35.0         37.0         40.2         22.0         22.5         72.0         79.4         42.0         60.3         54.0         46.3         58.2         46.1         49.1           0         24.5         30.0         35.0         40.2         22.0         22.5         72.0         79.4         42.0         60.3         54.0         48.3         58.2         46.1         49.1

Table 9: Full results of TimeBench. Aligned models are under zero-shot setting by default. Methods with  $\ddagger$  are base models without alignment, under few-shot setting, thus incomparable with other methods. We consider human performance as 100 points and normalize models' results accordingly.

DURATIONQA, MCTACO
Answer the following question, select all the possible correct options, and each question has at least one correct option. Context: {} Question: {} Options: {} Answer:
TIMEDIAL
There is a two-person dialogue with several options. Choose all appropriate options to substitute the <mask> in the dialogue, and each question has at least one correct option. Dialogue: {} Options: {} Answer:</mask>
TRACIE
Read the following story and hypothesis, determine whether the hypothesis can be inferred from the story. You need to understand the implicit temporal relationships between events to make judgments. Story: {} Hypothesis: {} Options: A. Entailment B. Contradiction Answer:
SituatedGen
Generate a pair of contrastive sentences with the given set of keywords. Keywords: {}
DATE ARITHMETIC
Question: {}? Answer:
TIMEQA
I will give you a question with context. You need to answer my question based on the context. If you can infer the answer from the context, then output your answer. Otherwise, if there is no answer, output [unan- swerable]. Context: {} Question: {} Answer:
TEMPREASON
I will give you a question with context. You need to answer my question based on the context. Context: {} Question: {} Answer:
MENATQA
Get answers for the question based on the contxt, where answers derived from substrings in the context or categorized as [unanswerable]. Context: {} Question: {} Answer:
TIMEX-NLI
Read the following statements about time and determine if the hypothesis can be inferred from the premise. Premise: {} Hypothesis: {} Options: A. Entailment B. Contradiction C. Neutral Answer:

Figure 9: Zeroshot instructions and input formats.

CoT Demonstration of TIMEX-NLI (3-shot, order) Answer the following question, select all the possible correct options, and each question has at least one correct option. Premise: On Wednesday, they got married. Hypothesis: Before Friday, they got married. Options: A. Entailment B. Contradiction C. Neutral Answer: Wednesday is before Friday. As a result, we can infer that if something happens on Wednesday, it definitely happens before Friday. Therefore, the answer is A. Entailment. Premise: We went to Disneyland on Monday. Hypothesis: We went to Disneyland after Wednesday. Options: A. Entailment B. Contradiction C. Neutral Answer: Monday is before Wednesday. As a result, We can infer that if something happens on Monday, it definitely can not happen after Wednesday. Therefore, the answer is B. Contradiction. Premise: The failing company issued major layoffs after Tuesday. Hypothesis: The failing company issued major layoffs after Thursday. Options: A. Entailment B. Contradiction C. Neutral

**Options:** A. Entailment B. Contradiction C. Neutral **Answer:** Tuesday is before Thursday. If something happened after Tuesday, we cannot be certain whether it occurred after Thursday. Therefore, the answer is C. Neutral.

Figure 10: Chain-of-Thought demonstrations of TimeX-NLI (s1-order).

#### **CoT Demonstration of DATE ARITHMETIC (4-shot)**

**Question**: What is the time 4 year and 1 month after Apr, 2000? **Answer**: First, 4 years after 2000 is 2004. Next, 1 month after April is May. Therefore, 4 year and 1 month after Apr, 2000 is May, 2004.

**Question**: What is the time 3 year and 4 month before Jun, 1840? **Answer**: First, subtracting 3 years from 1840 gives 1837. Next, subtracting 4 months from June gives February. Therefore, 3 year and 4 month before Jun, 1840 is Feb, 1837.

**Question**: What is the time 7 year and 11 month after Feb, 1819? **Answer**: First, 7 years after 1819 is 1826. Next, 11 months after February is January of the next year. Therefore, 7 years and 11 months after Feb, 1819 is Jan, 1827.

**Question**: What is the time 6 year and 9 month before Jan, 1234? **Answer**: First, subtracting 6 years from 1234 gives 1228. Next, subtracting 9 months from January gives April of the previous year. Therefore, 6 year and 9 month before Jan, 1234 is Apr, 1227.

Figure 11: Chain-of-Thought demonstrations of Date Arithmetic.

## **CoT Demonstration of TRACIE (4-shot)**

Read the following story and hypothesis, determine whether the hypothesis can be inferred from the story. You need to understand the implicit temporal relationships between events to make judgments

.....

**Story**: Joe was a police officer. Joe was patrolling the streets of the city in his cruiser. Suddenly, Joe was alerted to a crime happening near him by dispatch. Joe responded to the scene and found a bank robber fleeing on foot. Joe arrested the criminal and was promoted.

Hypothesis: Joe put on his police uniform. starts after Joe arrest the criminal

Options: A. Entailment B. Contradiction

**Answer**: From the story we know Joe was patrolling. In the work state, Joe has already put on the police uniform. So we can infer that Joe put on his police uniform before arresting the criminal. This conflicts with hypothesis. Therefore, the answer is B. Contradiction.

Figure 12: Chain-of-Thought demonstrations of TRACIE.

#### **CoT Demonstration of DURATIONQA (4-shot)**

Answer the following question, select all the possible correct options, and each question has at least one correct option.

•••••

**Context**: actually i have an project on it so please give me as much as you have information about migratory birds in punjab

Question: How long did it take for them to have information about migratory birds in punjab?

Options: A. several months B. 12 weeks C. a few minutes D. almost instantly

**Answer**: This is a conversation scenario. In the conversation, providing relevant information about migratory birds in punjab to him is in real-time and takes very little time. Therefore, the answer is C. a few minutes, D. almost instantly.

Context: Hope she stops laying eggs because she will get really skinny !

Question: How long did it take for her to lay eggs?

Options: A. 1 week B. 22 hours C. 2 years D. 4 years

**Answer**: According to commonsense knowledge, the time it takes for birds to lay eggs typically varies from one day to several days. Therefore, the answer is A. 1 week, B. 22 hours.

Figure 13: Chain-of-Thought demonstrations of DurationQA.

## **CoT Demonstration of MCTACO (4-shot)**

Answer the following question, select all the possible correct options, and each question has at least one correct option.

.....

**Context**: She ordered the tastiest kind of each vegetable and the prettiest kind of each flower.

Question: How often does she order vegetables and flowers?

**Options**: A. once a second B. three days a week C. every 10 centuries D. once a week

**Answer**: According to commonsense knowledge, ordering vegetables and flowers typically happens on a regular basis, usually every few days. Therefore, the answer is B. three days a week, D. once a week.

**Context**: Wallace, 38, called Gastonia home from the age of 8 until she graduated from Hunter Huss High School in 1983.

Question: When did Wallace wake up for high school?

Options: A. at 6 am B. at 1 am C. 7:00 AM D. at 6 pm

**Answer**: According to commonsense knowledge, waking up for high school typically happens in the morning, usually between 6 AM and 8 AM. Therefore, the answer is A. at 6 am, C. 7:00 AM.

Figure 14: Chain-of-Thought demonstrations of MCTACO.

#### **CoT Demonstration of TIMEDIAL (4-shot)**

There is a two-person dialogue with several options.

Choose all appropriate options to substitute the <mask> in the dialogue, and each question has at least one correct option.

.....

#### Dialogue:

A:What schools have you attended ?

B: I finished Young Primary School in 1998, and entered Xi ' an Middle School that same September . I graduated from there in  $\langle MASK \rangle$ , and that September I entered Wuhan University, where I'm studying now.

A: How do you think the education you have received will contribute to your work in this company ?

B: I think I have a good understanding of fundamentals in the areas your company deals with , and I can go on from here to build up the specific skills and knowledge I need to do my job well .

A: Your graduation thesis was on Medical Application of Laser, right? What were your conclusions?

B: Yes . I did some work on that , and I found out some really interesting things about the conductivity of liquid helium . I was sure I had a great discovery until my teacher told me the same discovery already made twenty years ago . I think the most important thing , I learnt though , was the importance of keeping good records . **Options:** A. 1998 B. July of 2004 C. March of 2003 D. twenty years ago

**Answer**: Based on the dialogue, B entered middle school in Sep 1998. According to commonsense knowledge, it usually takes around 6 years from entering middle school to graduating from high school (and entering university). Adding 6 years to 1998 would be 2004, so the answer should be around the year 2004. Therefore, the answer is B. July of 2004, C. March of 2003.

Figure 15: Chain-of-Thought demonstrations of TimeDial.

## CoT Demonstration of TIMEQA, MENATQA (2-shot, implicit)

I will give you a question with context.

You need to answer my question based on the context.

If you can infer the answer from the context, then output your answer. Otherwise, if there is no answer, output [unanswerable]

.....

**Context:** Theo-Ben Gurirab Theo-Ben Gurirab (23 January 1938 Ž013 14 July 2018) was a Namibian politician who served in various senior government positions. He served as the second Prime Minister of Namibia from 28 August 2002 to 20 March 2005, following the demotion and subsequent resignation of Hage Geingob. Previously he was the countrys first Minister of Foreign Affairs from 1990 to 2002, and was President of the United Nations General Assembly from 1999 to 2000. He was Speaker of the National Assembly of Namibia from 2005 to 2015, when he was replaced by Peter Katjavivi. Gurirab ultimately resigned from politics in 2015. Death. Gurirab died at a Windhoek hospital on 14 July 2018 of natural causes. He is buried at Heroes Acre.

Question: Theo-Ben Gurirab took which position after Jan 2007?

**Answer**: Based on the context, we can summarize the following facts: Theo-Ben Gurirab served as second Prime Minister of Namibia from August 2002 to March 2005. Prior to that, he was the countrys first Minister of Foreign Affairs from 1990 to 2002 and and was President of the United Nations General Assembly from 1999 to 2000. From 2005 to 2015, he held the position of Speaker of the National Assembly of Namibia. He resigned from politics in 2015 and passed away in July 2018. According to the aforementioned facts, he took the position of Speaker of the National Assembly of Namibia in January 2007. Therefore, the answer is Speaker of the National Assembly of Namibia.

Figure 16: Chain-of-Thought demonstrations of TimeQA, MenatQA, implicit reasoning.

## CoT Demonstration of TEMPREASON (4-shot, event-time)

I will give you a question with context. You need to answer my question based on the context.

.....

Context (facts): Gian Piero Gasperini is the head coach of Atalanta B.C. from Jun, 2016 to Dec, 2022.
Edoardo Reja is the head coach of Atalanta B.C. from Mar, 2015 to Jun, 2016.
Stefano Colantuono is the head coach of Atalanta B.C. from Jun, 2010 to Mar, 2015.
Bortolo Mutti is the head coach of Atalanta B.C. from Jun, 2010 to Jun, 2010.
Emiliano Mondonico is the head coach of Atalanta B.C. from Jul, 1987 to Jun, 1990.
Marcello Lippi is the head coach of Atalanta B.C. from Jul, 1992 to Jun, 1993.
Angelo Gregucci is the head coach of Atalanta B.C. from Jul, 2009 to Sep, 2009.
Luigi Delneri is the head coach of Atalanta B.C. from Jul, 2007 to Jun, 2010.
Ottavio Bianchi is the head coach of Atalanta B.C. from Jul, 1981 to Jun, 1983.
Antonio Conte is the head coach of Atalanta B.C. from Jul, 1983 to Jun, 1983.
Valter Bonacina is the head coach of Atalanta B.C. from Jul, 1983 to Jun, 1987.
Valter Bonacina is the head coach of Atalanta B.C. from Jun, 2010 to Jan, 2010.
Question: Who was the head coach of the team Atalanta B.C. in Feb, 2016?
Answer: According to the context, Edoardo Reja was the head coach of Atalanta B.C. from Mar, 2015 to Jun, 2016.

Figure 17: Chain-of-Thought demonstrations of TempReason, event-time reasoning.

#### **CoT Demonstration of TEMPREASON (4-shot, event-event)**

I will give you a question with context.

You need to answer my question based on the context.

•••••

**Context (facts)**: Nicholas Macpherson holds the position of Member of the House of Lords from Oct, 2016 to Dec, 2022.

Nicholas Macpherson holds the position of Principal Private Secretary to the Chancellor of the Exchequer from Jan, 1993 to Jan, 1997.

Nicholas Macpherson holds the position of Permanent Secretary to the Treasury from Aug, 2005 to Jan, 2016.

Question: Which position did Nicholas Macpherson hold before Member of the House of Lords?

**Answer**: According to the context, Nicholas Macpherson holds the position of Permanent Secretary to the Treasury from Aug, 2005 to Jan, 2016. Afterthat, Nicholas Macpherson holds the position of Member of the House of Lords from Oct, 2016 to Dec, 2022. Nicholas Macpherson hold the position of Permanent Secretary to the Treasury before Member of the House of Lords. Therefore, the answer is Permanent Secretary to the Treasury."

Figure 18: Chain-of-Thought demonstrations of TempReason, event-event reasoning.