FORMULAREASONING: A DATASET FOR FORMULA-BASED NUMERICAL REASONING

Anonymous authors

Paper under double-blind review

Abstract

The application of formulas is a fundamental ability of humans when addressing numerical reasoning problems. However, existing numerical reasoning datasets seldom indicate explicitly the formulas employed during the reasoning steps. To bridge this gap, we construct a dataset for formula-based numerical reasoning called FormulaReasoning, which consists of 5,420 reasoning-based questions. We employ it to conduct evaluations of LLMs with size ranging from 7B to over 100B parameters utilizing zero-shot and few-shot chain-of-thought methods, and we further explore using retrieval-augmented LLMs provided with an external formula database associated with our dataset. We also experiment with supervised methods where we divide the reasoning process into formula generation, parameter extraction, and numerical calculation, and perform data augmentation. Our empirical findings underscore the significant potential for improvement in existing models when applied to our challenging, formula-driven FormulaReasoning.

023

004

010 011

012

013

014

015

016

017

018

019

021

1 INTRODUCTION

025 026

Numerical reasoning constitutes one of the significant forms within natural language reason-027 ing (Frieder et al., 2023). The study of numerical reasoning has seen substantial progress in recent 028 years, largely driven by the development of LLMs (OpenAI, 2023; Touvron et al., 2023; Li et al., 029 2023c) and specialized datasets (Wang et al., 2017; Dua et al., 2019; Amini et al., 2019; Cobbe et al., 2021a). Current datasets for numerical reasoning typically include simple, commonsense nu-031 merical questions that do not reflect the complexity of real-world problems. These datasets have not fully addressed the interpretability issue in numerical reasoning, as they often rely on implicit 033 commonsense knowledge without explicit guidance knowledge during the reasoning process. This 034 issue becomes particularly evident when LLMs meet hallucination (Frieder et al., 2023; Bang et al., 035 2023). Consequently, one might naturally ask "What knowledge could be used to guide numerical reasoning process?". Formulas exactly represent such knowledge that has been largely overlooked 036 037 in previous research but is frequently utilized in real-life applications.

Take a question from the GSM8K (Cobbe et al., 2021a) as an example: "A robe takes 2 bolts of blue fiber and half that much white fiber. How many bolts in total does it take?". This example only requires the use of implicit *commonsense mathematical knowledge* to solve without domain-specific formula. However, in our FormulaReasoning dataset, we require *domain-specific formulas* to guide the numerical reasoning process, such as the formula used to calculate the heat absorption of an object.

Recently, Liu et al., 2023 constructed two formula-based datasets, Math23K-F and MAWPS-F. However, the formulas in these datasets primarily consist of commonsense formulas (such as total_amount = unit_amount \times total_number), and only 33.5% and 38.4% of the questions in these datasets, respectively, require the use of formulas.

To fill this gap, we constructed a dataset for numerical reasoning that requires the use of formulas called FormulaReasoning. We annotated formulas for each question in FormulaReasoning. An example of FormulaReasoning is shown in Figure 1.¹ The formula-based feature makes Formula-Reasoning a more challenging dataset for developing systems that can tackle real-world numerical

⁰⁵²

¹Please note that FormulaReasoning is originally in Chinese. For the convenience of understanding, we translated Chinese into English in all the examples presented in this paper.

54 Question

060

061

062

063

064

065

067

068

069 070 071

073

075 076 077

078

079

080 081

082

083

084 085

090

There is a electric water heater, after 50kg of water is loaded into its tank, the water is heated from 20°C to 60°C by electricity. It is known that the specific heat capacity of water is C_water = 4.2×10^3J/(kg*C).
If the total electrical energy consumed during the heating process is 1×10^7J, what is the thermal efficiency of the water heater?

Explanation (Reasoning Steps)

Calculating the degree of temperature increase in water: [Degree of water temperature increase] = [Final temperature] - [Initial temperature] = $60 \ ^\circ C - 20 \ ^\circ C = 40 \ ^\circ C$. The degree of water temperature increase = $40 \ ^\circ C$. The heat absorbed by water is given by: [Heat absorbed by water] = [Mass of water] * [Specific heat capacity of water] * [Degree of water temperature increase] = $50 \ \text{kg} * 4.2 * 10^3 \ \text{J}/(\text{kg} \cdot \text{C}) * 40 \ ^\circ \text{C} = 840000 \ \text{J}$. The thermal efficiency of the water heater can be obtained from: [Thermal efficiency of the water heater] = [Heat absorbed by water] / [Total electrical energy consumed] * 100% = $8400000 \ \text{J} / (1 * 10^7 \ \text{J}) * 100\%$ = 84%. The thermal efficiency of the water heater = 84%.

Parameter Table

Parameter	Symbol	Value	Unit
Degree of water temperature increase	Δt	40	°C
Final temperature	t _{final}	20	°C
		•••	
Heat absorbed by water	$Q_{absorbed}$	8400000	J
Mass of water	m _{water}	50	kg

Figure 1: An example taken from FormulaReasoning. Numerical values (including units) given in the question and obtained from intermediate steps are highlighted in red and purple, respectively. Formulas and their elements are in blue.

reasoning problems. Indeed, in fields such as mathematics and physics, formulas serve as an important vessel for representing domain knowledge. However, existing datasets scarcely consider explicit incorporation of formulas into numerical reasoning.

Table 1: Statistics of Math23-F, MAWPS-F, GSM8K, MATH and our FormulaReasoning.

Dataset	Math23K-F	MAWPS-F	GSM8K	MATH	FormulaReasoning
# questions	23,162	2,373	8,792	12,500	5,420
# formulas (and variants)	51 (131)	18 (46)	0 (0)	0 (0)	272 (824)
# questions requiring formula (proportion)	7,750 (33.46%)	911 (38.39%)	N/A	N/A	5,420 (100%)
Avg. # reasoning steps	1.16	1.01	3.59	Not Provided	2.37

092 We collected questions requiring formula-based numerical reasoning from Chinese junior high 093 school physics examinations. With the combined efforts of manual annotation and assistance from 094 LLMs, we annotated each question with an explanation text, a final answer, and a set of relevant 095 formulas (including formula structures, parameter names, symbols, numerical values, and units) 096 and built a consolidated formula database. The formula database functions as an external knowl-097 edge base, which can be used to evaluate retrieval-based/augmented systems. In Table 1, we com-098 pare FormulaReasoning with two existing formula-based datasets and the well-known GSM8K and 099 MATH (Hendrycks et al., 2021). In comparison to Math23K-F and MAWPS-F, FormulaReasoning contains a larger number of formulas (272), whereas the other two datasets contain 51 and 18 for-100 mulas. Additionally, all questions in FormulaReasoning require the use of formulas. The higher 101 average number of reasoning steps (2.37 vs. 1.16/1.01) implies that FormulaReasoning is more 102 challenging and better suited for evaluating existing models as a multi-step formula-based reasoning 103 task. 104

We used FormulaReasoning to evaluate LLMs ranging from 7B to >100B parameters, as well as fine-tuned models such as Qwen-1.8B (Bai et al., 2023) and ChatGLM3-6B (Zeng et al., 2022) with a proposed Chain-of-Thought supervised fine-tuned method and a data augmentation method. We also trained an encoder for formula retrieval and experimented with retrieval-augmented generative

108 models. Our empirical findings show that the best existing models only achieve an accuracy of 109 around 84%, lagging behind an accuracy 92% of humans, indicating that there is still significant 110 room for exploration in formula-based numerical reasoning.

111 Our contributions are summarized as follows: 112

- We construct a formula-based numerical reasoning dataset FormulaReasoning, with finegrained annotations for each question. As a formular knowledge-guided numerical reasoning dataset, it can be applied to tasks involving trustworthy and verifiable reasoning.
- We conduct evaluations on LLMs of various sizes, supervised fine-tuned models, and retrieval-augmented generative models. The experimental results establish a strong baseline for future research and also indicate that the task remains unresolved.

The dataset and code is currently available on anonymous GitHub https://anonymous. 4open.science/r/FormulaReasoning.

2 **RELATED WORK**

113

114

115

116

117

118

119 120

121

122 123

124

126

125 2.1 NUMERICAL REASONING DATASETS

127 Numerical reasoning is one of the fundamental capabilities of natural language reasoning. The 128 study of numerical reasoning in natural language has existed for several years. Numerous datasets, 129 such as DROP (Dua et al., 2019), GSM8K (Cobbe et al., 2021b), TSQA (Li et al., 2021) and 130 MATH (Hendrycks et al., 2021), have introduced natural language numerical reasoning. Another 131 line of research focusing on numerical reasoning in natural language is math word problem (MWP). 132 MWP tasks typically provide a short passage (i.e., a question) and require the generation of an arithmetic expression that can compute an answer. Representative datasets include MAWPS (Koncel-133 Kedziorski et al., 2016), Math23K (Wang et al., 2017), MathQA (Amini et al., 2019), etc. Several 134 works focus on specialized domains where some of the questions in their datasets require numer-135 ical reasoning. Examples include GeoSQA (Huang et al., 2019), which focuses on the geography 136 domain, the STEM (Drori et al., 2023) dataset and the ScienceQA (Lu et al., 2022) which covers 137 multiple disciplines in science and technology. The distinguishing feature of our FormulaReasoning 138 is that the numerical reasoning questions within these datasets lack explicitly labeled formulas. 139

The recently introduced datasets (Liu et al., 2023) Math23K-F and MAWPS-F require formulas for 140 only 33.5% and 38.4% of the questions, respectively, and the formulas within these datasets are 141 all simple commonsense formulas (e.g., total_cost = unit_cost \times total_number). By contrast, our 142 FormulaReasoning dataset collects questions from junior high school physics examinations, with 143 every question accompanied by formulas. In addition, we also annotated a *formula database* for 144 FormulaReasoning that can serve as an external knowledge base, used to assess retrieval-augmented 145 systems.

146

147 2.2 NUMERICAL REASONING METHODS 148

149 The methods for solving numerical reasoning have evolved from statistical approaches (Hosseini 150 et al., 2014; Kushman et al., 2014) to those based on rules and templates (Shi et al., 2015; Wang et al., 151 2019) and further to methods based on deep learning models (Gupta et al., 2019; Chen et al., 2022; 152 Kim et al., 2022; Li et al., 2023a). In the past two years, with the rapid development of LLMs, LLMs have demonstrated strong capabilities in resolving numerical reasoning questions. Consequently, 153 several methods aimed at enhancing the reasoning abilities of LLMs have been proposed, including 154 the notable Chain of Thoughts (CoTs) method (Wei et al., 2022), along with many subsequent variant 155 approaches (Kojima et al., 2022; Wang et al., 2022; Zhou et al., 2022; Li et al., 2023b). 156

157 We established representative existing methods as baselines for FormulaReasoning, including 158 zero/few-shot CoTs prompting methods to LLMs ranging from 7B to over 100B parameters. We 159 trained a specialized formula retriever for retrieving formulas and explored retrieval-enhanced numerical reasoning. We also divided the reasoning process into formula generation, parameter ex-160 traction, and numerical calculation, and used data augmentation to enhance fine-tuned models with 161 fewer than 7B parameters.

162 3 DATASET CONSTRUCTION

164 We collected raw questions from Chinese junior high school physics examinations from 2015 to the 165 present. We had a total of five postgraduate volunteer students, and they all hold a bachelor's degree 166 in science and engineering. We then annotated the reasoning steps and corresponding formulas for 167 each question. This process involved a combination of manual annotation and the assistance of 168 LLMs to improve the efficiency of annotation. Each question is associated with an explanation of the reasoning steps in natural language with a symbolic representation of the reasoning steps using 170 formulas, including the values and units for all the parameters within the formulas. Finally, we compiled all the formulas and we merged those expressing the same meaning to create a formula 171 database. We describe this process to construct FormulaReasoning in detail below. 172

173 174 3.1 PREPROCESSING

We crawled 18,433 junior high school physics examination questions in China from 2015 to the present from public sources, including only those with free-text answers and excluding multiple-choice and true/false questions. Each raw question contains a *question text* and an *explanation text that includes the reasoning steps*. We eliminated questions requiring diagrams.

Subsequently, we filtered the questions by assessing the presence of numerical values within the explanation and confirming that the final answer was numerical. Utilizing a regular expression-based approach, we extracted the *final numerical answer*, including its unit, from the explanation. We found that for 487 questions, the regular expressions did not return results, so we manually annotated the positions of their answers in the text explanations. Following the preprocessing phase, we compiled an initial dataset comprising 6,306 questions.

Table 2: Original explanation and explanation with normalized formulas (highlighted in blue).

Original explanation.

The change in water temperature is 60 - 20 = 40 °C. Therefore, the heat absorbed by the water is $Q_{absorbed} = 50 \text{ kg} \times 4.2 \times 10^3 \text{ J/(kg.°C)} \times 40 \text{ °C} = 8.4 \times 10^6 \text{ J}$. Given that the total electrical energy consumed in the heating process is $1 \times 10^7 \text{ J}$, the thermal efficiency of the water heater can be calculated using the formula for the efficiency of a heat engine: $\eta = Q_{absorbed}/W_{absorbed} \times 100\% = (8.4 \times 10^6 \text{ J})/(1.0 \times 10^7 \text{ J}) \times 100\% = 84\%$. Answer: If it is known that the total electrical energy consumed during the heating process is 1×10^7 , the thermal efficiency of the water heater is 84%.

Explanation with normalized formulas.

1. Calculating the temperature increase in water: [Degree of water temperature increase] = [Final temperature] - [Initial temperature] = $60 \degree C - 20 \degree C = 40 \degree C$. The degree of water temperature increase = $40 \degree C$. 2. Calculating the heat absorbed by water: [Heat absorbed by water] = [Mass of water] × [Specific heat capacity of water] × [Degree of water temperature increase] = $50 \text{ kg} \times 4.2 \times 10^3 \text{ J/(kg} \degree C) \times 40 \degree C = 8400000 \text{ J}$. 3. The thermal efficiency of the water heater can be obtained from: [Thermal efficiency of the water heater]

= [Heat absorbed by water] / [Total electrical energy consumed] $\times 100\%$ = 8400000 J / (1 $\times 10^7$ J) * 100% = 84%. The thermal efficiency of the water heater = 84%. Answer = 84%

202 203

187

188

189

190

191

192

193 194

195

196

197

199

200

201

204 205

3.2 FORMULA NORMALIZATION

206 We found that the reasoning steps (i.e. the explanation) in the obtained raw dataset lacked a normal-207 ized format and were expressed quite casually. Some formulas mixed parameter names (e.g., "mass 208 of water") and symbols (e.g., " m_{water} "), while others simply provided calculations in numerical 209 form without parameter names or symbols. In order to ensure that all explanations adopted a nor-210 malized form of formulas, we normalized the formula annotations in the explanations. An example 211 can be found in Table 2. In this process, we need to *identify the formulas used within the original* 212 explanations and to correct any formatting issues. Manually undertaking such tasks would require 213 significant effort. However, since the process is not open-ended, but rather structured and verifiable, we could automatically, e.g., using a LLM, extract formulas from the explanations, calculate 214 each step, and compare the result with the given answer to ensure the accuracy of this normalization 215 process.

Specifically, to enhance the efficiency of the annotation, we adopted a coarse-to-fine annotation approach with the help of a LLM². We first prompted the LLM in a few-shot manner to generate accurate explanations of the reasoning process. Then, we used few-shot prompts to guide the LLM in correcting minor errors within the normalized explanations, including formatting errors in formula annotations and inaccuracies in the parameters used during computations. Both prompts can be found in Appendix A.1.1. Next, we will provide a detailed description of this process.

Initially, we introduced the question along with its original explanation and the corresponding answer to guide the LLM through few-shot prompting to revise the original explanation. We observed that the ability of the LLM to revise explanations towards normalized explanations remained satisfactory. To assess the correctness of the revised explanations, we extracted formulas from these explanations and then computed the answer using the numbat tool³. In addition to providing explanations, we also required the LLM to present the values, symbols, and units of each parameter in the formulas in the form of a table. An example is shown in Figure 1.

229 At this stage, we checked the correctness of the formula format in the explanations by automatic 230 rules, including whether there were omissions in parameter names, parameter symbols, or corre-231 sponding units, and these issues were all correctable. Therefore, if our program detected that the 232 LLM had not successfully generated an accurate normalized explanation, we used few-shot prompt-233 ing to identify and correct these specific errors. More details can be found in Appendix A.1.1. We observed that the questions which remained incorrect despite multiple attempts by the LLM were 234 of notably poor quality, including missing important reasoning steps, unclear question formulation, 235 and so on. Some examples of these questions can be found in Appendix A.1.2. These questions 236 were removed from our dataset. Following this step, our dataset contains a remaining total of 5,420 237 questions. 238

239 240

3.3 FORMULA DATABASE CONSTRUCTION

241 Our next step was to construct a unified formula 242 database for the entire dataset. Given that pa-243 rameters in the same formula can be expressed 244 differently across various problem contexts, for 245 instance, the two formulas "[weight of water] = 246 [mass of water] * [gravitational acceleration]" 247 and "[weight] = [mass] * [gravitational acceleration]" both calculate the weight of an object, 248 we need to merge these formulas into a single 249 representation. 250

Table 3: Changes in the number of formulas after each merging step.

Step	# Formulas
Before merging	12,906
After symbolic rules based merging	1,163
After semantic-based merging	439
After manual review and error correction	272

We divided the construction process of the formula database into three steps: 1) Merge the formulas through symbolic rules. 2) Merge the formulas through semantic-based method. 3) Manual review and error correction. In Table 3, we present the initial number of formulas and the remaining number of formulas after each step.

256 **Symbolic rules based merging.** In this step, we merged formulas through symbolic rules. Specif-257 ically, this was achieved by comparing the structure of the formulas and the symbols. Take the following as an example of judging whether two formulas have the same structure: the formulas 258 " $f_1: a_1 = (b_1 + c_1)/d_1$ ", " $f_2: a_2 = (b_2 + c_2)/d_2$ " and " $f_3: b_1 = a_1 * d_1 - c_1$ " have the same 259 structure because f_2 can be derived from f_1 by renaming parameters, and f_3 can be obtained from f_1 260 by transformation. Moreover, in physics, certain physical quantities are conventionally represented 261 by specific symbols. For example, the mass of an object is often denoted by "m" and the density 262 of an object is frequently represented by the symbol " ρ ". Subscripts are then used to distinguish 263 which specific object a physical quantity refers to, such as " ρ_{water} " for the density of water. For 264 any two formulas, we first computed all the transformations of each formula to obtain a set of all 265 its variants. Then, we compared the formula structures in the two sets to determine if two formulas 266 were structurally equivalent. If they shared the same structure, we then compared whether their

²⁶⁷ 268 269

²During dataset construction, we accessed Qwen-max via API (https://help.aliyun.com/zh/dashscope/developerreference/quick-start). Qwen-max is a LLM with over 100B parameters and a strong capability in Chinese.

³https://numbat.dev. Numbat is designed for scientific computations with support for physical units.

270 symbols, with subscripts removed, were identical. If they were, we considered these two formulas 271 to be mergeable. When merging, we retained the parameter with the shorter length from the two. 272 After merging based on symbolic rules, we reduced the number of formulas in the formula database 273 from 12,906 to 1,163.

274

275 **Semantic-based merging.** In the symbolic rules based merging process, the semantic information 276 of the parameter names was neglected. This led us to perform merges grounded on the semantics 277 of the parameter names. For instance, two formulas that were not merged during the symbolic fusion stage, "[density] = [mass] / [volume]" and "[density of water] = [mass of water] / [volume 278 of water]", can actually be merged. We would carry out the merging of these two formulas based 279 on the semantic information of the parameter names (for example, "density" and "density of water" 280 are semantically similar). Specifically, for formulas with identical structures, we tokenized each 281 pair of corresponding parameters to create two sets of words⁴. When the two sets overlapped, the 282 parameters were considered to have semantic connection, and the formulas became candidates for 283 merging. Utilizing this approach, we identified a set of pairs of potentially mergeable formulas 284 and then consulted the LLM for a thorough evaluation of each pair. The prompts can be found in 285 Appendix A.1.3. After this step, the number of formulas in the formula database was reduced to 286 439.

288 Manual review and error correction. Upon completing the aforementioned merging process, we 289 manually inspected the correctness of the results, rectified instances where errors occurred during merging, and manually merged formulas that were overlooked by the LLM. In this process, there 290 were two human volunteers cross-validating the results of manual review and annotation. Finally, we obtained a formula database consisting of 272 formulas. 292

293 294

291

287

4 **EXPERIMENTS SETUP**

295

296 In this section, we explore several methods for handling the questions within FormulaReasoning, 297 including prompting LLMs using zero-shot and few-shot chain-of-thought (CoT, Wei et al., 2022; 298 Kojima et al., 2022), and training a formula retriever to retrieve formulas to be incorporated into 299 LLM prompts. Additionally, we employed two approaches to enhancing the reasoning abilities of fine-tuned models with fewer than 7B parameters. The first approach involved dividing the reasoning 300 process into distinct steps: formula generation, parameter extraction, and numerical calculation. The 301 second approach leveraged data augmentation to improve the models' reasoning ability. 302

4.1 DATASET SPLIT 304

305 We divided FormulaReasoning into into subsets for training, id (in-distribution) test, and ood (out-306 of-distribution) test, comprising 4,608, 421 and 391 questions, respectively. We required that all 307 formulas in the id test must appear in the training set, whereas in the ood test, each question involves 308 at least one formula that has not been seen in the training set. This division is designed to evaluate 309 the generalizability of fine-tuned models on formulas that they have not previously encountered.

310 311

319

320

323

303

- 4.2 EVALUATED METHODS 312
- 313 4.2.1 HUMAN PERFORMANCE 314

We recruited 108 students from a high school, with each student being assigned 7–8 questions. Each 315 student was given 40 minutes to complete these questions. These questions were used as part of their 316 in-class exercises, and at the end, each student received a gift. The final statistics were collected to 317 evaluate human performance, which was consented by all the students. 318

4.2.2 LLMs

321 Following Kojima et al., 2022, we incorporated the phrase "Let's think step by step" into the zero-322 shot prompt to guide LLMs in generating the reasoning steps. For the few-shot setting, we randomly

⁴We used jieba: https://github.com/fxsjy/jieba.

sampled five questions from the training set to serve as examples for in-context learning. Each
 example includes the question text and the reasoning steps (i.e., the explanation). Examples of the
 prompts can be found in Appendix A.2.2.

We conducted experiments on GPT-40, GPT-4-turbo, GPT-3.5-turbo, GLM-4-plus, GLM-4flash (GLM et al., 2024), and Qwen-max. We also evaluated on Qwen2.5-7B/14B (Yang et al., 2024) and Llama3.1-8B (Meta, 2024).

4.2.3 FORMULA RETRIEVER

333 We trained a formula retriever on the training set. Specifically, we encoded each question using the 334 Chinese-BERT-wwm-base (Devlin et al., 2019; Cui et al., 2021) model to obtain the CLS vector of the question. Each formula in the formula database was represented by a randomly initialized 335 vector. During training, we calculated the cosine score between the question vector and the formula 336 vector. The retriever was then trained with in-batch negatives and contrastive learning loss (Gao 337 et al., 2021). Subsequently, for each question in the id test, we retrieved the top five formulas with 338 the highest scores and included them in the prompt to observe the change in the performance of the 339 LLM when provided with relevant formulas. More details can be found in Appendix A.2.3. 340

341
 342
 4.2.4 SUPERVISED FINE-TUNED MODELS

We found that directly prompting models possessing fewer than 7B parameters failed to produce satisfactory outcomes (for example, ChatGLM3-6B attained merely 8.99 points in a zero-shot setting). Therefore, we conducted supervised fine-tuning of models with fewer than 7B parameters, yet discerned that, dissimilar to larger models (such as GLM-4-plus), smaller models did not exhibit proficient performance in numerical extraction and calculation. In order to augment the reasoning capabilities of smaller models, we explored two approaches for improvement. We conducted experiments on Qwen-1.8B (Bai et al., 2023) and ChatGLM3-6B (Zeng et al., 2022).

350

362 363

Chain-of-Thought Supervised Fine-Tuning (CoT-SFT) We decomposed the reasoning process into several steps. First, we instructed the model to generate the formulas required to solve the question. Subsequently, the parameter names within the formulas were extracted, allowing the model to retrieve the corresponding values and units from the context. Next, the formulas and the associated parameter values were provided to a calculator to obtain the final result. This approach relieved the model from numerical calculation, allowing it to concentrate on the reasoning aspect.

357 Data Augmentation (DA) We augmented the training dataset with the assistance of larger models.
 358 Firstly, we utilized a few-shot approach to prompt a LLM (Qwen-max) to generate new question 359 answer pairs. The correctness of the computation process generated by the LLM was meticulously
 360 verified using a calculator. Subsequently, the formulas generated by the model were extracted and
 361 normalized. More details could be found in Appendix A.2.1.

4.3 METRIC

We utilized numbat to evaluate the predictions generated by the model against the gold-standard answers. A prediction is deemed correct if the relative error (prediction - gold) / gold is less than 1%. We employed *accuracy*, which is the proportion of questions answered correctly, as our metric.

368369 4.4 IMPLEMENTATION DETAILS

We accessed to GPT-40 (gpt-40-2024-08-06 version), GPT-4-turbo (gpt-4-1106-preview version), GPT-3.5-turbo (gpt-3.5-turbo-1106 version)⁵, GLM-4-plus, GLM-4-flash⁶, Qwen-max and Qwen2.5-7B/14B⁷ through API calls with the default hyper-parameters. For Llama3.1, we conducted experiments on NVIDIA V100-32G GPUs. These LLMs generated using nucleus sampling with top_p=0.8. Models that require fine-tuning were experimented on NVIDIA V100 GPUs with

377 ⁶https://open.bigmodel.cn/

⁵https://platform.openai.com/docs

⁷https://help.aliyun.com/zh/dashscope/developer-reference/quick-start

378 Huggingface Transformers and Pytorch 2.0. For Qwen-1.8B, we used a learning rate of 1e-5 and a 379 batch size of 32, and tested the model after training for 10 epochs. For ChatGLM3-6B, we fine-tuned 380 with LoRA (Hu et al., 2021) with r=8, alpha=32 and learning rate of 5e-5, batch size of 1. The max 381 input length and output length are both set to 512. We utilized nucleus sampling with top_p=0.8 for 382 generation. In the case of CoT-SFT, which directly outputted formulas along with corresponding parameter values and units, if the generation output contained formatting errors, we allowed the small 383 model to retry up to 5 times until a correctly formatted output was generated. Training Qwen-1.8B, 384 ChatGLM3-6B models required 12 and 24 hours respectively. 385

386 387

388 389

390

391

392 393

394

- **5** EXPERIMENTS RESULTS
- 5.1 HUMAN PERFORMANCE

In FormulaReasoning, humans achieved impressive performance, with a score of 93.49 on the id test, 90.47 on the ood test, and an average score of 92.03.

5.2 RESULTS OF LLMS

395 396 397

Table 4: Results of LLMs with zero-shot and few-shot prompting.

Model	Size	ze	ro-shot Co	т	fe	w-shot Co	Т
WIOdel	5120	id test	ood test	Avg.	id test	ood test	Avg.
GPT-40	unknown	77.20	72.38	74.88	76.01	73.66	74.88
GPT-4-turbo	unknown	70.07	72.89	71.43	71.50	77.49	74.38
GPT-3.5-turbo	unknown	26.13	25.58	25.87	32.07	29.92	31.03
GLM-4-plus	>100B	84.32	81.07	82.76	82.90	85.68	84.24
GLM-4-flash	unknown	71.50	71.87	71.68	61.76	67.01	64.29
Qwen-max	>100B	57.24	60.10	58.62	55.82	61.38	58.50
Qwen2.5	14B	61.28	64.71	62.93	61.28	65.22	63.18
Qwen2.5	7B	42.04	43.73	42.38	59.62	65.73	62.56
Llama3.1	8B	13.06	9.74	11.46	9.74	9.72	9.73
Human	-	93.49	90.47	92.03	93.49	90.47	92.03

409 410 411

The evaluation results on LLMs are shown in Table 4. GLM-4-plus exhibited the best performance 412 in both zero-shot and few-shot settings, surpassing the second-ranked GPT-40 by an average of 7.88 413 points in zero-shot setting and 9.36 in few-shot setting. Among models with size not exceeding 414 20B, Qwen2.5-14B demonstrated commendable performance in both zero-shot and few-shot set-415 tings. The subpar performance of Llama3.1 might be due to its pre-training data being primarily in 416 English. After incorporating few-shot examples, GPT-4-turbo, GPT-3.5-turbo, GLM-4-plus and 417 Owen2.5 demonstrated performance improvements, ranging from 0.25 to 20.18. However, similar 418 performance changes were not observed on other LLMs. Surprisingly, the open-source Qwen2.5-419 14B model outperformed the closed-source Qwen-max model⁸.

Human performance surpassed the performance of the flagship model GLM-4-plus with zero-shot setting and few-shot setting by margins of 9.27 and 7.79 points, respectively. Such results demon-strated that there remained a substantial gap between the current capabilities of state-of-the-art LLMs and human performance. This was even more pronounced when considering smaller-scale models.
These findings underscored *the challenging nature of FormulaReasoning as an unresolved dataset*, and that there was significant room for improvement in LLMs as they struggled to match human levels of reasoning.

We also compared the chain of thought (CoT) and program of thought (PoT, Chen et al., 2023)
methods, with the results presented in Appendix A.2.4. The results indicated that CoT and PoT demonstrated varying performances between different models and under different settings.

⁸We have not yet found clear information indicating whether the closed-source Qwen-max is also based on version 2.5.

432 5.3 RESULTS OF LLMS WITH FORMULA RETRIEVER

434 The results of LLMs utilizing the formula re-435 triever are shown in Table 5. We found that the 436 impact on performance varied among different 437 LLMs when incorporating retrieved formulas 438 into prompts. We observed a positive enhance-439 ment on Qwen2.5-7B, with score increments of 440 10.92 and 4.04 with zero-shot and few-shot, re-441 spectively, on the id test. However, we found that the performance was essentially on par on 442 the GLM-4-flash. Specifically, we found that 443 the top 5 retrieved formulas often included ir-444 relevant ones, as the number of formulas re-445 quired varies for each problem. The presence 446

Table 5: Results of LLMs with Formula Retriever on the id test.

Model	zero-shot	few-shot
GLM-4-flash	71.50	61.76
+ formula retriever	70.55	62.95
Qwen2.5-7B	42.04	59.62
+ formula retriever	52.96	63.66

of these extraneous formulas affected the model's performance, indicating that there is considerable
 room for further research in retrieving from a formula database.

5.4 RESULTS OF SUPERVISED FINE-TUNED MODELS

453 Table 6 shows the results for the supervised fine-tuned models, with and without CoT-SFT 454 and DA, which were detailed in Section 4.2.4. 455 In most settings, both models achieved higher 456 scores on the id test than the ood test, yet they 457 still exhibited considerable performance on the 458 ood test. This indicates that 1) the ood formu-459 las indeed challenged model performance and 460 2) the models still demonstrated a certain level 461 of generalizability. We hope that the division of 462 id test and ood test will be helpful for assessing 463 the generalization ability of fine-tuned models in future work. 464

Table 6: Results of supervised fine-tuned models on FormulaReasoning.

Model	Size	id test	ood test	Avg.
Qwen-1.8B	1.8B	55.91	44.58	50.25
+ DA		56.16	45.32	50.74
+ CoT-SFT		73.65	74.38	74.00
ChatGLM3-6B	6B	52.95	40.64	47.02
+ DA		53.44	45.32	49.53
+ CoT-SFT		74.63	73.89	74.23

465 It was noteworthy that with CoT-SFT, Qwen-1.8B and ChatGLM3-6B, with a mere parameter count 466 of 1.8B and 6B, respectively, achieved performance comparable to GPT-40 (though such a com-467 parison may not be entirely fair). This indicated that the incorporation of CoT-SFT and the use of 468 calculators could significantly enhance the reasoning capabilities of small models. Our findings re-469 vealed that focusing on reasoning with CoT while delegating numerical calculation to a calculator 470 could enhance the performance of small models, given their limited calculating capability. The assistance of LLMs for data augmentation could also enhance smaller models' reasoning capability. 471 This finding provides valuable insights for future deployment of numerical reasoning systems with 472 small models. 473

474

449 450

451 452

475 476

5.5 CASE STUDY AND ERROR ANALYSIS

477 478

We sampled 50 error cases from the id test (few-shot setting) of GPT-3.5-turbo and manually cat-479 egorized the types and proportions of errors. We divided the error types into two main categories: 480 formula errors and calculation errors. Formula errors encompass inappropriate formulas and omit-481 ted formulas, while calculation errors primarily involve inaccuracies in numerical calculation and 482 unit errors. We found that 38% of errors were caused by incorrect formulas, while the remaining 483 62% were attributable to calculation errors. We provide one example for each of the two types of errors listed in Figure 2. It could be observed that FormulaReasoning poses challenges to existing 484 models in terms of formula application and numerical calculation (including unit calculation and 485 arithmetic calculation).

486 The maximum power of a tank engine is 7.2×10^{5} W. During a certain mission, the tank traveled 3000m on a straight road at a constant speed using its maximum power for 300 seconds, consuming 15kg of diesel fuel. (q_diesel = 4.0×10^{7} J/kg) Q: The efficiency of the tank's engine. 487 488 489 Output of GPT-3.5-turbo (wrong formula highlighted in red) Solution: 1. Calculate the time it takes for the tank to travel: ... omitted ... 490 3. Calculate the engine efficiency: ... omitted 491 Gold Explanation (correct formula highlighted in green) 492 Solution: 1. Calculate the heat released by the complete combustion of diesel fuel: ... omitted ... 3. Calculate the efficiency of the tank's engine: [[Tank engine efficiency] = [Work done by the complete combustion of diesel fuel] × 100% Equation ... omitted ... 493 494 495 (a) An error case caused by wrong formulas. 496 Given that the calorific value of natural gas is 3.8×10^7 J/m^3 and the specific heat capacity of water is 4.2×10^3 J/(kg·°C). If the efficiency of a natural gas stove in heating water is 50%, then at standard atmospheric pressure, how many degrees Celsius can 10L of natural gas, burning completely, heat 3.0kg of water initially at 60°C? 497 498 499 Output of GPT-3.5-turbo (wrong calculation highlighted in red. Note that 1 L = 10⁻{-3} m⁻³) 500 Solution: 1. Calculate the heat content of natural gas: [Heat content of natural gas] = [Calorific value of natural gas] × [Volume of natural gas] Calculation = <mark>3.8×10^7 J/m^3 × 10 L = 3.8×10^8 J H</mark>eat content of natural gas = 3.8×10^8 J *... omitted ...* 501 Gold Explanation (correct calculation highlighted in green) Solution: 1. Calculate the heat released from the complete burning of natural gas: [Heat released from complete combustion of natural gas] = [Volume of natural gas] × [Calorific value of natural gas] Calculation = 10 L × 3.8×10^7 J/m³ = 3.8 × 10^5 J ... omitted ... 504 505 (b) An error case caused by wrong calculation. 506 507 Figure 2: Error cases. 508 509

6 CONCLUSION AND LIMITATIONS

511 512

510

We introduced FormulaReasoning, a dataset for formula-based numerical reasoning. We annotated 513 the reasoning steps with formulas for each question with both manual and LLM-assisted efforts. Fur-514 thermore, we constructed a formula database after merging formulas with similar meanings, serving 515 as an external knowledge base for subsequent retrieval-based/augmented approaches. We evalu-516 ated FormulaReasoning across various sizes of LLMs, supervised fine-tuned models, and retrieval-517 augmented LLMs, demonstrating its challenging nature as an unresolved task. Our findings indicate 518 substantial room for improvement of existing models on formula-based numerical reasoning, thus 519 motivating future research efforts.

In the future work, we plan to utilize the formula knowledge from FormulaReasoning to improve the numerical reasoning capabilities of LLMs. Possible approaches include enhancing reasoning abilities through knowledge-driven methods, preference learning methods based on formula feedback.

One limitation of this work is that our evaluation results reported in the paper were obtained from the 524 original Chinese version of FormulaReasoning. We have employed a combination of LLM-based 525 translation and manual review to release an English version of FormulaReasoning. Currently, we 526 provide a preview English version in our GitHub repository, and we will release the official English 527 version of FormulaReasoning after completing the manual review process. Another limitation is 528 that, our dataset is limited to the domain of physics. Although junior high school physics is not 529 overly complex and can be understood by most people which would benefit evaluation efforts, it is 530 still possible to explore formula-based question answering data in other domains. 531

532

533 REFERENCES

Aida Amini, Saadia Gabriel, Shanchuan Lin, Rik Koncel-Kedziorski, Yejin Choi, and Hannaneh Hajishirzi. MathQA: Towards interpretable math word problem solving with operation-based formalisms. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pp. 2357–2367, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1245. URL https://aclanthology.org/N19-1245. 540 Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenhang Ge, 541 Yu Han, Fei Huang, Binyuan Hui, Luo Ji, Mei Li, Junyang Lin, Runji Lin, Dayiheng Liu, Gao 542 Liu, Chengqiang Lu, K. Lu, Jianxin Ma, Rui Men, Xingzhang Ren, Xuancheng Ren, Chuanqi 543 Tan, Sinan Tan, Jianhong Tu, Peng Wang, Shijie Wang, Wei Wang, Shengguang Wu, Ben-544 feng Xu, Jin Xu, An Yang, Hao Yang, Jian Yang, Jian Yang, Shusheng Yang, Shusheng Yang, Bowen Yu, Yu Bowen, Hongyi Yuan, Zheng Yuan, Jianwei Zhang, Xing Zhang, Yichang Zhang, 545 Zhenru Zhang, Chang Zhou, Jingren Zhou, Xiaohuan Zhou, and Tianhang Zhu. Qwen techni-546 cal report. ArXiv, abs/2309.16609, 2023. URL https://api.semanticscholar.org/ 547 CorpusID:263134555. 548

Yejin Bang, Samuel Cahyawijaya, Nayeon Lee, Wenliang Dai, Dan Su, Bryan Wilie, Holy Lovenia, Ziwei Ji, Tiezheng Yu, Willy Chung, Quyet V. Do, Yan Xu, and Pascale Fung. A multitask, multilingual, multimodal evaluation of chatgpt on reasoning, hallucination, and interactivity. *ArXiv*, abs/2302.04023, 2023. URL https://api.semanticscholar.org/ CorpusID:256662612.

 Jiayi Chen, Xiao-Yu Guo, Yuan-Fang Li, and Gholamreza Haffari. Teaching neural module networks to do arithmetic. In *Proceedings of the 29th International Conference on Computational Linguistics*, pp. 1502–1510, Gyeongju, Republic of Korea, October 2022. International Committee on Computational Linguistics. URL https://aclanthology.org/2022.coling-1.
 129.

- Wenhu Chen, Xueguang Ma, Xinyi Wang, and William W. Cohen. Program of thoughts prompting: Disentangling computation from reasoning for numerical reasoning tasks. *Trans. Mach. Learn. Res.*, 2023, 2023. URL https://openreview.net/forum?id=YfZ4ZPt8zd.
- 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. Training verifiers to solve math word problems. *ArXiv*, abs/2110.14168, 2021a. URL https://api.semanticscholar.org/CorpusID:239998651.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser,
 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to
 solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021b.
- 570 Yiming Cui, Wanxiang Che, Ting Liu, Bing Qin, and Ziqing Yang. Pre-training with whole word masking for chinese bert. *IEEE Transactions on Audio, Speech and Language Processing*, 2021. doi: 10.1109/TASLP.2021.3124365. URL https://ieeexplore.ieee.org/ document/9599397.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pp. 4171–4186, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1423. URL https://aclanthology.org/N19-1423.
- Bosheng Ding, Chengwei Qin, Ruochen Zhao, Tianze Luo, Xinze Li, Guizhen Chen, Wenhan Xia, Junjie Hu, Anh Tuan Luu, and Shafiq Joty. Data Augmentation using LLMs: Data Perspectives, Learning Paradigms and Challenges. *arXiv e-prints*, art. arXiv:2403.02990, March 2024. doi: 10.48550/arXiv.2403.02990.
- Iddo Drori, Sarah Zhang, Zad Chin, Reece Shuttleworth, Albert Lu, Linda Chen, Bereket Birbo,
 Michele He, Pedro Lantigua, Sunny Tran, et al. A dataset for learning university stem courses
 at scale and generating questions at a human level. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pp. 15921–15929, 2023.
- Dheeru Dua, Yizhong Wang, Pradeep Dasigi, Gabriel Stanovsky, Sameer Singh, and Matt Gardner.
 DROP: A reading comprehension benchmark requiring discrete reasoning over paragraphs. In
 Proceedings of the 2019 Conference of the North American Chapter of the Association for Com- putational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), pp.
 2368–2378, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1246. URL https://aclanthology.org/N19-1246.

- Simon Frieder, Luca Pinchetti, Ryan-Rhys Griffiths, Tommaso Salvatori, Thomas Lukasiewicz, Philipp Christian Petersen, Alexis Chevalier, and J J Berner. Mathematical capabilities of chatgpt. ArXiv, abs/2301.13867, 2023. URL https://api.semanticscholar.org/ CorpusID:256415984.
- Tianyu Gao, Xingcheng Yao, and Danqi Chen. SimCSE: Simple contrastive learning of sentence embeddings. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pp. 6894–6910, Online and Punta Cana, Dominican Republic, November 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.emnlp-main.552. URL https://aclanthology.org/2021.emnlp-main.552.
- Team GLM, Aohan Zeng, Bin Xu, Bowen Wang, Chenhui Zhang, Da Yin, Diego Rojas, Guanyu
 Feng, Hanlin Zhao, Hanyu Lai, et al. Chatglm: A family of large language models from glm-130b
 to glm-4 all tools. *arXiv preprint arXiv:2406.12793*, 2024.
- Nitish Gupta, Kevin Lin, Dan Roth, Sameer Singh, and Matt Gardner. Neural module networks for reasoning over text. In *International Conference on Learning Representations*, 2019.
- ⁶¹⁰ Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn
 ⁶¹¹ Song, and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset.
 ⁶¹² In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks* ⁶¹³ Track (Round 2), 2021.
- Mohammad Javad Hosseini, Hannaneh Hajishirzi, Oren Etzioni, and Nate Kushman. Learning to solve arithmetic word problems with verb categorization. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 523–533, Doha, Qatar, October 2014. Association for Computational Linguistics. doi: 10.3115/v1/D14-1058. URL https://aclanthology.org/D14-1058.
- Edward J Hu, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen,
 et al. Lora: Low-rank adaptation of large language models. In *International Conference on Learning Representations*, 2021.
- Zixian Huang, Yulin Shen, Xiao Li, Gong Cheng, Lin Zhou, Xinyu Dai, Yuzhong Qu, et al. Geosqa: A benchmark for scenario-based question answering in the geography domain at high school level. 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)*, pp. 5866–5871, 2019.
- Jeonghwan Kim, Junmo Kang, Kyung-min Kim, Giwon Hong, and Sung-Hyon Myaeng. Exploiting numerical-contextual knowledge to improve numerical reasoning in question answering. In *Findings of the Association for Computational Linguistics: NAACL 2022*, pp. 1811–1821, Seattle, United States, July 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.
 findings-naacl.138. URL https://aclanthology.org/2022.findings-naacl. 138.
- Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large
 language models are zero-shot reasoners. *Advances in neural information processing systems*, 35:22199–22213, 2022.
- Rik Koncel-Kedziorski, Subhro Roy, Aida Amini, Nate Kushman, and Hannaneh Hajishirzi.
 MAWPS: A math word problem repository. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 1152–1157, San Diego, California, June 2016. Association for Computational Linguistics. doi: 10.18653/v1/N16-1136. URL https://aclanthology.org/N16-1136.
- Nate Kushman, Yoav Artzi, Luke Zettlemoyer, and Regina Barzilay. Learning to automatically solve algebra word problems. In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 271–281, Baltimore, Maryland, June 2014. Association for Computational Linguistics. doi: 10.3115/v1/P14-1026. URL https://aclanthology.org/P14-1026.

- Kiao Li, Yawei Sun, and Gong Cheng. Tsqa: tabular scenario based question answering. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pp. 13297–13305, 2021.
- Xiao Li, Yin Zhu, Sichen Liu, Jiangzhou Ju, Yuzhong Qu, and Gong Cheng. Dyrren: A dynamic retriever-reranker-generator model for numerical reasoning over tabular and textual data. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pp. 13139–13147, 2023a.
- Yifei Li, Zeqi Lin, Shizhuo Zhang, Qiang Fu, Bei Chen, Jian-Guang Lou, and Weizhu Chen. Making
 language models better reasoners with step-aware verifier. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 5315–
 5333, Toronto, Canada, July 2023b. Association for Computational Linguistics. doi: 10.18653/
 v1/2023.acl-long.291. URL https://aclanthology.org/2023.acl-long.291.
- Yuan-Fang Li, Sébastien Bubeck, Ronen Eldan, Allison Del Giorno, Suriya Gunasekar, and Yin Tat Lee. Textbooks are all you need ii: phi-1.5 technical report. *ArXiv*, abs/2309.05463, 2023c. URL https://api.semanticscholar.org/CorpusID:261696657.
- Jia-Yin Liu, Zhenya Huang, Zhiyuan Ma, Qi Liu, Enhong Chen, Tianhuang Su, and Haifeng Liu.
 Guiding mathematical reasoning via mastering commonsense formula knowledge. Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, 2023. URL https://api.semanticscholar.org/CorpusID:260499357.
- Pan Lu, Swaroop Mishra, Tanglin Xia, Liang Qiu, Kai-Wei Chang, Song-Chun Zhu, 667 Oyvind Tafjord, Peter Clark, and Ashwin Kalyan. Learn to explain: Multimodal rea-668 soning via thought chains for science question answering. In S. Koyejo, S. Mo-669 hamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (eds.), Advances in Neural 670 Information Processing Systems, volume 35, pp. 2507-2521. Curran Associates, Inc., 671 2022. URL https://proceedings.neurips.cc/paper_files/paper/2022/ 672 file/11332b6b6cf4485b84afadb1352d3a9a-Paper-Conference.pdf. 673
- 674 Meta. Meta llama 3, 2024. URL https://llama.meta.com/llama3/.
- 675
 OpenAI. Gpt-4 technical report. ArXiv, 2023. URL https://api.semanticscholar.org/ CorpusID:257532815.
- Shuming Shi, Yuehui Wang, Chin-Yew Lin, Xiaojiang Liu, and Yong Rui. Automatically solving number word problems by semantic parsing and reasoning. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pp. 1132–1142, Lisbon, Portugal, September 2015. Association for Computational Linguistics. doi: 10.18653/v1/D15-1135. URL https://aclanthology.org/D15–1135.
- Kashun Shum, Shizhe Diao, and Tong Zhang. Automatic prompt augmentation and selection with
 chain-of-thought from labeled data. In Houda Bouamor, Juan Pino, and Kalika Bali (eds.), *Find- ings of the Association for Computational Linguistics: EMNLP 2023*, pp. 12113–12139, Singapore, December 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.
 findings-emnlp.811. URL https://aclanthology.org/2023.findings-emnlp.
 811.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Niko lay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open founda tion and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- Lei Wang, Dongxiang Zhang, Jipeng Zhang, Xing Xu, Lianli Gao, Bing Tian Dai, and Heng Tao
 Shen. Template-based math word problem solvers with recursive neural networks. In *Proceedings* of the AAAI Conference on Artificial Intelligence, volume 33, pp. 7144–7151, 2019.
- Kuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V Le, Ed H Chi, Sharan Narang, Aakanksha
 Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language
 models. In *The Eleventh International Conference on Learning Representations*, 2022.
- Yan Wang, Xiaojiang Liu, and Shuming Shi. Deep neural solver for math word problems. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pp. 845–854, Copenhagen, Denmark, September 2017. Association for Computational Linguistics. doi: 10.18653/v1/D17-1088. URL https://aclanthology.org/D17-1088.

- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in Neural Information Processing Systems*, 35:24824–24837, 2022.
- Chenxi Whitehouse, Monojit Choudhury, and Alham Fikri Aji. LLM-powered data augmentation for enhanced cross-lingual performance. In *The 2023 Conference on Empirical Methods in Natural Language Processing*, 2023. URL https://openreview.net/forum?id= wWFWwyXElN.
- An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, et al. Qwen2 technical report. arXiv preprint arXiv:2407.10671, 2024.
- Aohan Zeng, Xiao Liu, Zhengxiao Du, Zihan Wang, Hanyu Lai, Ming Ding, Zhuoyi Yang, Yifan Xu, Wendi Zheng, Xiao Xia, et al. Glm-130b: An open bilingual pre-trained model. In *The Eleventh International Conference on Learning Representations*, 2022.
- Chujie Zheng, Sahand Sabour, Jiaxin Wen, Zheng Zhang, and Minlie Huang. AugESC: Dialogue augmentation with large language models for emotional support conversation. In *Findings of the Association for Computational Linguistics: ACL 2023*, pp. 1552–1568, Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.findings-acl.99. URL https://aclanthology.org/2023.findings-acl.99.
 - Denny Zhou, Nathanael Schärli, Le Hou, Jason Wei, Nathan Scales, Xuezhi Wang, Dale Schuurmans, Claire Cui, Olivier Bousquet, Quoc V Le, et al. Least-to-most prompting enables complex reasoning in large language models. In *The Eleventh International Conference on Learning Representations*, 2022.
- 726 727 728 729

722

723

724

725

A APPENDIX

- 730 A.1 DATASET CONSTRUCTION
- 732 A.1.1 PROMPTS IN FORMULA NORMALIZATION

The process of formula normalization is delineated into three distinct stages: the generation of natu-734 ral language explanations, the extraction of the associated parameters from the explanations, and the 735 subsequent error correction phase. The initial two stages are illustrated in Figures 3 and 4. The third 736 stage is further splited into three specific error categories, each addressed by a dedicated prompt: 737 input errors, where the parameters mentioned in the explanation are absent from the question; cal-738 culation errors, which occur when the calculator reports an error during the computation process; 739 and output errors, where the final computed answer is incorrect. We provide an example here fo-740 cusing on prompts for correcting calculation errors, while prompts for the other two error types can 741 be found in our code submission. The prompts designed to correct calculation errors are depicted in Figure 5. The entire normalization procedure employs a 6-shot prompting, an instance of which is 742 provided herein for illustrative purposes. 743

- 744 745
- A.1.2 EXAMPLES OF DELETED QUESTIONS

The questions which remained incorrect despite multiple attempts by the LLM were of notably poor
 quality, including missing important reasoning steps, wrong reference answer, and so on. Here is an
 example of these questions in Figure 6.

749 750

A.1.3 SEMANTIC-BASED MERGING FOR FORMULA DATABASE CONSTRUCTION

Semantic-based merging primarily employs the LLM to comprehend formulas, ascertain if two formulas are semantically equivalent, and subsequently determine whether they can be merged into a
single formula. The prompt for this procedure is illustrated in Figure 7. This approach ensures that
the nuanced meanings embedded within formulas are accurately captured and evaluated for potential
merging, thereby enhancing the quality of formula database.

756 A.2 EXPERIMENTS

758 A.2.1 DATA AUGMENTATION (DA) FOR FORMULAREASONING

759 760 There have been several studies utilizing large language models (LLMs) for data augmentation (Ding 761 et al., 2024). The data generated in these related works (Zheng et al., 2023; Whitehouse et al., 2023) 762 primarily focus on daily conversations or sentiment analysis and do not require rigorous numerical 763 calculations. Some research on data augmentation involving numerical calculations (Shum et al., 764 2023) employs LLMs to generate solutions to questions to aid in training, rather than creating com-765 plete questions. In contrast to these approaches, our work generates complete questions that involve 766 numerical calculations (particularly formula calculations), along with automatic improvement and 766 selection to ensure data quality.

In order to enhance the capabilities of models, we use LLM to generate more data for fine-tuning.
 We divide the process of data generation into the following several steps.

First, we randomly generated 17,000 prompts. Each prompt was obtained by stacking five questionanswer pairs sampled form training set. At the end of the prompt, LLM was required to generate the sixth question-answer pair. Second, we normalized the generated formulas. Except for the absence of manual review, the remaining steps were consistent with those in Section 3.2. At last, we unitized the calculator to check whether the calculation process in the data generated by the LLM is correct, and discarded the generated data with incorrect calculation processes. After the above steps, we finally retained more than 2500 questions.

777 We found that mixing the newly generated data into the original training set did not always bring 778 positive improvement, perhaps because the newly generated data has not undergone manual re-779 view. We found that randomly selecting a small portion of the newly generated data can enable 780 the model to have performance improvement. We set several different mixing ratios selected from 781 {5%, 10%, 15%, 20%, 2%, 30%, 35%, 40%}. We fine-tuned each model using the augmented data 782 set. After training for a fixed number of steps (150k and 200k), we selected the checkpoints with the 783 smallest loss among models of different mixing ratios.

- 784 A.2.2 ZERO-SHOT AND FEW-SHOT PROMPTS
- Zero-shot and few-shot prompts are shown in Figure 8.
- 788 A.2.3 FORMULA RETRIEVER

796

797

801 802 803

Let the number of formulas in the formula database be N. During training, we randomly initialized a matrix $\mathbf{F} \in \mathbb{R}^{N \times d}$, where d is the hidden size and the *i*-th row in \mathbf{F} represented the initial representation of the *i*-th formula in formula database. We denoted a batch of questions with a batch size of B as $Q = \{q_1, q_2, ..., q_B\}$. The indices of the gold-standard formulas corresponding to these Bquestions were denoted as $L = \{l_1, l_2, \cdots, l_B\}$ (i.e. the label of q_i is l_i , where $1 \le i \le B$).

5 BERT was utilized to encode each question,

 $\mathbf{h}_{cls}^{i}, \mathbf{h}_{1}^{i}, \dots = \mathtt{BERT}(q_{i}), 1 \le i \le B.$ $\tag{1}$

Subsequently, we took the CLS vector \mathbf{h}_{cls}^i as the representation for the *i*-th question.

800 We utilized in-batch negatives and contrastive learning loss,

$$\mathcal{L} = -\frac{1}{B} \sum_{1 \le i \le B} \log \frac{\exp(\cos(\mathbf{h}_{cls}^{i}, \mathbf{F}_{l_{i}}))}{\sum_{1 \le j \le B} \exp(\cos(\mathbf{h}_{cls}^{i}, \mathbf{F}_{l_{j}}))}.$$
(2)

Each question might correspond to multiple correct formulas, and we ensured that the same question did not appear twice in the same batch when loading the data. Based on the implementation of Chinese-BERT-wwm-base, we tested the retrieval performance on the id test set and found that Recall@5 reached 97.69%.

Models were evaluated with top-5 retrieved formulas. Prompts can be found in Appendix A.2.5. We utilized zero-shot CoTs.

Model		zero-shot			few-shot	
WIOUCI	id test	ood test	Avg.	id test	ood test	Avg.
GPT-40 (CoT)	77.20	72.38	74.88	76.01	73.66	74.88
GPT-40 (PoT)	80.76	73.91	77.46	81.47	82.61	82.02
GLM-4-plus (CoT)	84.32	81.07	82.76	82.90	85.68	84.24
GLM-4-plus (PoT)	84.08	78.51	81.40	86.70	84.91	85.84
Human	93.49	90.47	92.03	93.49	90.47	92.03

Table 7: Results of LLMs with zero-shot and few-shot chain of thought (CoT) and program of
 thought (PoT).

A.2.4 COMPARISON OF COT AND POT PROMPTS

Results are shown in Table 7. In the PoT approach, we utilized a Python interpreter to execute the code and obtain the final results. We found that the performance comparison between CoT and PoT varies across models. GPT-40 consistently demonstrated superior performance with PoT across all settings, achieving improvements of 2.58 points on average in the zero-shot setting and 7.14 points on average in the few-shot setting. In contrast, GLM-4-plus showed an average decline of 1.36 points in the zero-shot setting but showed an average improvement of 1.60 points in the few-shot setting. The finding might be related to the code capabilities of the models.

A.2.5 PROMPTS FOR LLMS WITH FORMULA RETRIEVER

We added the formulas before each question in the few-shot setting. For the examples sampled from the training set, gold-standard formulas were added before each question. For the final question from the test set in both zero-shot and few-shot prompts, we included the top 5 retrieved formulas. The prompts are shown in Figure 9.

866 867 868 870 871 872 873 874 875 876 877 878 **Prompt actually used English translation** 879 我需要你修改问题原有的解析,给出规范格式的新解析,要求 I need you to modify the original explanation of the question and 加下 rovide a new explanation with the following requirements: 1.请逐步地进行思考,如果有公式组合的部分需要一步步地拆分 1. Please think step by step. If there has formula combination, you 成基本公式进行求解 2.公式中的计算符号,如"+"、"-"、"×"、"/"和"^"不能省略 need to decompose the combination into basic formulas step by step 2. Calculation symbols such as "+", "-", "×", "/" and "^" in formulas 3.公式需要同时给出符号和有具体含义的两种形式,然后代入 cannot be omitted. 883 数值计算得出答案 3. The formula needs to be given in both symbolic and concrete 4.涉及到单位换算的部分需要展示出来具体过程 forms. After that, you need to substitute into the numerical 5.使用"[]"标注公式中的变量,其中科学计数法形式的数字 calculation to obtain the answer. "a×10'b"以及复杂的单位,需要使用"()"标注 6.如果有latex格式的公式,比如"\frac{Q_吸}{Q_放}"需要改成 4. The part related to unit conversion needs to show the specific 885 proces 正常算式的形式: "Q_吸/Q_放" 5. Use "[]" to label variables in formula, "()" for numbers like 7.解析开始如果给出"(解析可能存在错误)"表明给出的解析中 "a×10^b" in scientific notation and for complex units. 887 可能存在错误 6. If you have a latex formula such as $\frac{Q_i}{Q_i} \{Q_out\}$, change 888 下面是一些例子: it to a normal formula: "Q in /Q out". 7. Start with "(Explanation may be wrong)" indicates that there input: 889 问题: 在某次测试中该车匀速行驶时受到的阻力是车重的 0.02 may be an error in the given explanation. 倍,在 0.5h 内行驶了 28km 的路程,求:行驶的平均速度为多 890 Here are some examples: 小 input: 891 解析:解:车辆在0.5h内行驶了 28km 的路程,所以行驶的平 Ouestion: In a test, the resistance of the car at a constant speed is 均速度为: V=s/t=28km/0.5h=56km/h。 892 0.02 times the weight of the car, and the car traveled a distance of 28km in 0.5h. What is the average speed? output: 893 Explanation: The car has traveled a distance of 28km in 0.5h. The 解 原始解析比较清晰,下面将其转换为规范格式 average speed of the vehicle is V=s/t=28km/0.5h=56km/h. 894 第一步,先计算平均速度: output: 895 v=s/t Solution: [平均速度]=[路程]/[时间] The raw parsing is clear, so let's convert it to a canonical format: 896 算式=28km/0.5h=56km/h 1.we compute the average speed: 平均速度=56km/h v=s/t [Average speed]=[distance]/[time] 答案=56km/h expression=28km/0.5h=56km/h ...omitted... 下面是一个新的问题 Average speed =56km/h 问题: {{question}} Answer=56km/h 900 解析: {{explanation}} ...omitted... There's a new question: 901 Question: {{question}} 902 Explanation: {{explanation}} 903 Figure 3: Prompt of the formula normalization stage 1. 904 905 906 907 908 909

911 912

910

- 913
- 914
- 915
- 916
- 917

969 970 971

919 920 921 922 923 **Prompt actually used** English translation 924 我需要你帮我从给定的问题和解析中提取出解决问题使用到的 I need you to help me extract the formula and parameters table for 925 公式,以及公式对应的参数表 the formula from given question and explanation: 926 1.请逐步地进行思考,先对解析进行分析, 然后生成提取结果 1. Think step by step, analyze the explanation first, and then 如果有公式组合的部分需要一步步地拆分成基本公式进行求解 generate the extracted results. If there is a combination of formulas, 927 2.公式中所有的计算符号,如"+"、"-"、"×"、"/"和"^"不能省略 3.公式中的每个变量需要使用"[]"标注出来,而且变量需要使用 the combination needs to be split into basic formulas step by step. 2. All calculation symbols such as "+", "-", "×", "/" and "^" in the 928 有意义的文字标识,尽量避免直接使用数值 formula cannot be omitted. 929 4.如果有latex格式的公式,比如" $frac{Q_W}{Q_b}$ "需要改成 3. Each variable in the formula needs to be labeled with "[]", and the 正常算式的形式: "[Q_吸]/[Q_放]", 算式中的单位换算部分不 930 variable needs to be identified with meaningful text instead of 属于公式,不需要被提取 numbers 931 5.参数表中的参数是公式中使用到的参数(参数名称要与公式 4. If a latex formula such as $\frac{Q_i}{Q_i} \{Q_ou\}$ needs to be 中的参数一致),表格包括:概念、符号、数值、单位,使用"|" changed to a normal formula: [Q_in]/[Q_out]. The unit conversion 932 分割单元格 does not need to be extracted. 933 6.参数表中的数值和单位来自于问题本身以及解析计算的中间 结果,如果参数进行了单位换算,参数表要给出原始的参数形 5. The parameters table come from the parameters in formula (the parameter name should be consistent with the parameters in the 934 式 (没有讲行单位换算) formula), the table include: concept, symbol, numeric, unit, using 935 7.参数表中的参数如果没有单位或者是表示某个物体的数量, cell division " 比如"3盏灯"、"转了8圈"、"4个人"等等,那么单位写入"null" 6. The numeric and unit in the parameter table come from the 936 8.提取结果中给出的公式和参数表中不需要任何额外的解释信 problem itself and the intermediate results of analytical calculation. 937 If the parameters are converted into different units, the parameter 下面是一些例子: table should give the original parameter form (without unit 938 conversion). input: 问题: 在某次测试中该车匀速行驶时受到的阻力是车重的 0.02 7. If the parameter in the parameters table has no units or represents 939 the amount of an object, such as "3 lights", "8 revolutions", "4 倍, 少? 在 0.5h 内行驶了 28km 的路程, 求: 行驶的平均速度为多 940 people", etc., then the units are written as "null" -解析:解 8. Apart from formula and parameter table, no additional 941 原始解析比较清晰,下面将其转换为规范格式: information is required in the extraction results. 942 第一步,先计算平均速度: Here are some examples: input: 943 [平均速度]=[路程]/[时间] Question: In a test, the resistance of the car at a constant speed is 944 算式=28km/0.5h=56km/h 0.02 times the weight of the car, and the car traveled a distance of 平均速度=56km/h 28km in 0.5h. What is the average speed? 945 答案=56km/h Explanation 1.we compute the average speed: 946 output: 提取分析 v=s/t 947 1.解析中的公式格式符合要求,计算符号没有省略,变量使用 [Average speed]=[distance]/[time] "[]"标注,而且变量使用有意义的文字标识 2.计算过程中没有进行单位换算,参数可以直接提取 expression=28km/0.5h=56km/h 948 Average speed =56km/h 949 Answer=56km/h 下面给出提取结果(单位换算部分不属于公式,不需要被提取; 参数表中不需要额外解释信息) output: 950 提取结果 Extraction analysis: 1.先计算平均速度,根据路程和时间的关系: 1. The formula format in the analysis meets the requirements. The 951 [平均速度]=[路程]/[时间] calculation symbols are not omitted. Variables are labeled with "[]", 952 参数表: and variables are expressed with meaningful text. 概念| 符号| 数值| 单位 2. No unit conversion was performed during the computation, and 953 路程| s| 28| km parameters can be directly extracted. 954 时间| t| 0.5| h Below is the extraction result (the unit conversion part does not 平均速度| v| 56| km/h belong to the formula and does not need to be extracted; no 955 omitted additional explanatory information is required in the parameter 下面是一个新的问题 956 table). 问题: {{question}} Extraction result: 957 解析: {{explanation}} 1. First calculate the average speed, based on the relationship between distance and time 958 [average speed]=[distance]/[time] 959 Parameter table Concept | Symbol | Numeric | Unit 960 distance | s | 28 | km 961 time | t | 0.5 | h average speed | v | 56 | km/h ...omitted... 962 963 There's a new question Question: {{question}} 964 Explanation: {{explanation}} 965 966 Figure 4: Prompt of the formula normalization stage 2. 967 968

Prompt actually used	English translation
我需要你帮助我纠正解析中的错误,我会给出问题和错误信息.	I need your help to correct the error in the explanation. I will
下面是错误纠正的要求:	provide the question and error information. The following are the
1.你需要先进行错误分析,分析如何修改米纠止错误,然后给 出错误纠正部分,纠正解析中的错误	requirements for error correction: 1. You need to first conduct error analysis, analyze how to modify to
2.错误纠正部分不需要任何额外解释信息、错误纠正部分的格	correct the error, and then provide the error correction to correct the
式为:"内容: 修改前的内容->修改后的内容", 增加内容时"修改 前的内容"为null 删除内容时"修改后的内容"为null	error in the explanation. 2 The error correction section does not require any additional
3.问题缺失参数:如果问题中没有缺失的参数,那么向题目中	explanatory information. The format of the error correction section
增加缺失的参数;如果问题中的参数与缺失参数的含义相同但格 式不同 修改题目中的参数与缺失参数相同	is: "Content: Pre modified Content ->Modified Content". When adding content "Pre modified Content" is null and when deleting
4.算式错误:算式存在错误需要对公式和错误的参数进行修改,	content, "Modified Content" is null.
如果算式中存在"[参数]"或"null",需要补齐缺失的参数;如果 参数没有问题可能需要对公式进行修改	3. Missing parameters in the question: If there are no missing
5.公式的格式为"[待求解参数]=[参数1](+ - × /)[参数2]";参数	question; If the parameters in the question have the same meaning
表的格式为:"概念 符号 数值 单位",比如"水的沸点是100℃", 表示为"水的沸点上;沸山001℃"	as the missing parameters but different formats, modify the
表示为 小的沸点 [1_沸] 100 [0	 Expression error: The formula and incorrect parameters need to
下面是一些例子:	be modified. If there is "[parameter]" or "null" in the expression, the
mput: 问题:假设13.0r烟煤在煤炉中完全燃烧,放出的热量部分被水	missing parameters need to be filled in; If there are no issues with the parameters, it may be necessary to modify the formula.
吸收,可以使4×10^5kg的水从20℃升高到100℃,求水吸收的	5. The format of the formula is "[parameter to be
热量为多少J[c_水=4.2×10^3J / (kg·℃)] 结误信息	solved]=[parameter 1] (+ - \times /) [Parameter 2]"; The format of the parameter table is: "concept symbol numeric unit" for example
算式错误:1.计算水升高的温度差:	"The boiling point of water is 100 °C", which is represented as "the
公式: [水升高的温度差]=[末温]-[初温] 質式-[末現] [初現]	boiling point of water t_boiling 100 °C"
异式气不温气初温」 问题缺失参数:水升高的温度差=80 °C;	input:
output: 维문分析	question: Assuming that 13.0 tons of bituminous coal is completely humad in a coal furnace, the heat released is partially absorbed by
1.根据错误信息:算式存在错误,而且算式中存在"[参数]"的情	water, which can make 4×10^{5} kg of water is raised from 20 °C to
况: "[末温]"、"[初温]", 需要对参数表增加缺失的参数	100 °C, how much heat is absorbed by the water ?
依据镇长镇长信息, [术血]-[初血], 从题日中可以找到相关 文本"从20℃升高到100℃", 按照要求的参数格式表示为:	Error message:
初温 t_0 20 ℃	Expression error: 1. Calculate the temperature difference of water
末篇[4]100] C 这样参数表增加缺失的参数后,代入1. 计算水升高的温度差的	temperature] - [initial temperature]
公式可以得到: 第二十-(100) %() (20) %()-%0 %()	expression=[final temperature] - [initial temperature]
算式=((100) ℃)-(20) ℃)=80 ℃ 水升高的温度差=80 ℃	rise=80 °C;
2.根据错误信息,问题缺失参数,由于分析1中纠正算式后计算	output:
得到了"水井局的温度差=80°C",所以问题个冉缺矢参数,个 需要讲行修改	Error analysis: 1. According to the error message: there is an error in the formula.
	and there are cases of "[parameter]" in the formula: "[final
错误纠止: 参数表:_mull->初温 t_0 20 ℃	temperature]" and "[initial temperature]". It is necessary to add then to the parameter table
参数表: null->末温 t 100 ℃	According to the error message, "[Last Temperature] - [Initial
<i>omitted</i> 下面早—个新的问题:	Temperature]", the relevant text "Increase from 20 °C to 100 °C"
问题: {{question}}	it is represented as:
错误: {{error}}	Initial temperature $ t_0 20 $ °C
	After adding missing parameters to the parameter table, substitute $T_{1}^{(1)}$
	them into the formula for calculating the temperature difference can
	be obtained as follows: expression= $((100) \ ^{\circ}C)-((20) \ ^{\circ}C)=80 \ ^{\circ}C$
	temperature difference of water rise=80 °C
	2. According to the error message, the question is missing parameters. After analyzing the correction equation in 1 step it was
	calculated that "the temperature difference of water rise=80 °C", so
	the question is no longer missing parameters and does not need to
	Error correction:
	Parameter table: null ->Initial temperature $ t_0 20 ^{\circ}C$
	Parameter table: null ->final temperature t 100 °C
	There's a new question:
	Question: {{ question }}
	EITOT: {{error}}

Question:	1 37 1 0 '		1 . 1. 01.
As shown in the figure, f	the Xuelong 2 scien	tific research icebreaker d	lesigned in China.
resistance experienced h	w the icebreaker is	$\frac{1000}{2}$	alculate the
propulsion power of the	icebreaker at this ti	me.	alculate the
Reference answer: 2×10	0^7 W		
Formula:			
[thrust]=[resistance]			
[propulsion power]=[thr	ust]×[constant spee	d]	
Parameter table:			
Parameter	symbol	value	unit
resistance	f	2×10^7	Ν
ship speed	v	1	m/s
Explanation:	•		
1.Calculate thrust:			
thrust=resistance= 2×10^{4}	^7N		
2.Calculate propulsion p	ower:		
propulsion power=thrus	t×constant speed=2	×10^7N×constant speed(a	cannot find value)
Error:			
1. The parameter "resis	tance" in the questi	on is in the incorrect form	nat.
2. "constant speed" cou	ld not be located in	the parameter table.	
F .			
F1g	ure 6: An examp	le of deleted questions	5.
F1g	ure 6: An examp	le of deleted questions	5.
F1g	ure 6: An examp	le of deleted questions	5.
F1g	ure 6: An examp	le of deleted questions	5.
F1g	ure 6: An examp	le of deleted questions	5.
F1g	ure 6: An examp	le of deleted questions	5.
F1g	ure 6: An examp	le of deleted questions	5.
F1g	ure 6: An examp	le of deleted questions	5.
F1g	ure 6: An examp	le of deleted questions	5.
F1g	ure 6: An examp	le of deleted questions	5.
F1g	ure 6: An examp	le of deleted questions	S.
Fig	ure 6: An examp	le of deleted questions	s. Jation
F1g Prompt actu 下面我会给出两个公式,每个 构成,II中的表示参数。	ure 6: An examp ally used -公式由参数和运算符号	le of deleted questions English trans	S. Iation w. Each formula consists o ols. The text in ∏ represen
Prompt actu 下面我会给出两个公式,每个构成,[]中的表示参数。 你需要判断我给出的两个公式。不同一副不同一一个人们。	ure 6: An examp ally used 公式由参数和运算符号 行对应参数表达含义是	le of deleted questions English trans I will give two formulas belor parameters and operation symb parameter.	S. Iation w. Each formula consists o ols. The text in [] represen
Fig Prompt actu 下面我会给出两个公式,每个 构成,□中的表示参数。 你需要判断我给出的两个公式 否相同,是否是同一个公式: 如果含义不相同,不是同一个	ure 6: An examp ally used 公式由参数和运算符号 C中对应参数表达含义是 -公式,只需要回答不是;	le of deleted questions English trans I will give two formulas belor parameters and operation symb parameter. You need to judge whether the co two formulas I gave have the san	S. lation w. Each formula consists o ols. The text in [] represen prresponding parameters in the me meaning and whether they
Fig Prompt actu 下面我会给出两个公式,每个 构成,[]中的表示参数。 你需要判断我给出的两个公式 否相同,是否是同一个公式: 如果含义不相同,不是同一个 如果各个参数含义相同,最合一个 如果各个参数含义相同,是一一个 如果含义不相同,不是同一个	ure 6: An examp ally used 公式由参数和运算符号 公式,只需要回答不是;)一个公式,则需要短给出 (六句达表根本表示会教的	le of deleted questions English trans I will give two formulas belor parameters and operation symb parameter. You need to judge whether the cc two formulas I gave have the sau are the same formula: If the meaning is different, and t	S. Jation w. Each formula consists o rols. The text in [] represen presponding parameters in the me meaning and whether they hey are not the same formula
Fig Prompt actu 下面我会给出两个公式,每个 构成,但中的表示参数。 你需要判断我给出的两个公式 否相同,是否是同一个公式: 如果含义不相同,不是同一个 如果各个参数之相同,不是同一个 如果各个参动这人相同,不是同一个 如果各个参动这人相同,不是同一个 如果各个参动话,并且给出一个言 对应关系,每个单元格内容是	ure 6: An examp ally used 公式由参数和运算符号 公式,只需要回答不是; 一个公式,则需要的当本是; 一个公式,则需要的当本是; 一个参数,前两行填写 过一个参数,前面行填写	le of deleted questions English trans I will give two formulas belor parameters and operation symb parameter. You need to judge whether the cc two formulas I gave have the san are the same formula: If the meaning is different, and t just answer no;	S. Ilation w. Each formula consists o ols. The text in [] represen prresponding parameters in the me meaning and whether they hey are not the same formula
Fig Prompt actu 下面我会给出两个公式,每个 构成,[]中的表示参数。 你需要判断我给出的两个公式 否相同,是否是同一个公式: 如果各个参数含义相同,是居 是终的公式,并且给出一个三 对应关系,每个单元格内容是 两个公式的参数,第三行填写 下面是公式1:	ure 6: An examp ally used 公式由参数和运算符号 公式,由参数表达含义是 公式,则需要回答出 一个公式,则需要给出 一个参数,前两行填写 统一后的公式参数。	le of deleted questions English trans I will give two formulas belor parameter. You need to judge whether the co two formulas I gave have the sau are the same formula: If the meaning is different, and t just answer no; If each pair of parameters have th the same formula, the final form	S. Ilation w. Each formula consists o ols. The text in [] represen prresponding parameters in the me meaning and whether they hey are not the same formula he same meaning, and they are ula needs to be given, and i
Fig Prompt actu 下面我会给出两个公式,每个 构成,□中的表示参数。 你需要判断我会出的两个公式 否相同,是否是同一个公式: 如果各个参数含义相同,是同 最终的公式,并且给出一个三 对应关系,每个单元格内容是 两个公式的参数,第三行填写 下面是公式1: {公式1) 下面是人士2:	ure 6: An examp ally used 公式由参数和运算符号 公式,只需要回答不是;)一个公式,则需要给出 行的表格来表示参数的 (一个参数,前两行填写 3统一后的公式参数。	le of deleted questions English trans I will give two formulas belor parameter. You need to judge whether the co two formulas I gave have the sau are the same formula: If the meaning is different, and t just answer no; If each pair of parameters have th the same formula, the final form three-row table needs to be given	S. Jation w. Each formula consists o ols. The text in [] represen prresponding parameters in the me meaning and whether they hey are not the same formula he same meaning, and they are ula needs to be given, and it to indicate the corresponding The actuat of onk of they
Fig Prompt actu 下面我会给出两个公式,每个 构成,□中的表示参数。 你需要判断我给出的两个公式 否相同,是否是同一个公式: 如果各个参数含义相同,是后 最终的公式,并且给出一个三 对应关系,每个单元格内容是 两个公式的参数,第三行填写 下面是公式1: {公式1} 下面是公式2: {公式2}	ure 6: An examp ally used 公式由参数和运算符号 计中对应参数表达含义是 公式,只需要回答不是;)—个公式,则需要给出 行的表格来表示参数的 —个参数,前两行填写 3统一后的公式参数。	le of deleted questions English trans I will give two formulas belor parameters and operation symb parameter. You need to judge whether the cc two formulas I gave have the sau are the same formula: If the meaning is different, and t just answer no; If each pair of parameters have th the same formula, the final form three-row table needs to be given relationship between the parameter parameter, and the first two rows	S. Jation W. Each formula consists o ols. The text in [] represen orresponding parameters in th me meaning and whether they hey are not the same formula he same meaning, and they ar ula needs to be given, and to indicate the corresponding rs. The content of each cell is: are filled with two formulas
Fig Prompt actu 下面我会给出两个公式,每个 构成,[]中的表示参数。 你需要判断我给出的两个公式 否相同,是否是同一个公式: 如果各义不相同,不是同一个 如果各义不相同,不是同一个 如果各义不相同,不是同一个 如果各义不相同,不是同一个 如果各义不相同,不是同一个 如果各义不相同,不是同一个 如果各义不相同,不是同一个 如果各人和我们,不是而一个 如果各人和我们,不是而一个 如果各人和我们,不是而是同一 不过来了。 下面是公式1: (公式 1) 下面是公式2: (公式 2) 通过表达含义判断,是否是同	ure 6: An examp ally used 公式由参数和运算符号 计对应参数表达含义是 公式,只需要回答不是;)一个公式,则需要给出 行的表格来表示参数的 一个参数,前两行填写 3统一后的公式参数。	le of deleted questions English trans I will give two formulas belor parameters and operation symb parameter. You need to judge whether the cc two formulas I gave have the sau are the same formula: If the meaning is different, and t just answer no; If each pair of parameters have th the same formula the final form three some formula, the final form three some formula the first two rows Parameters, fill in the unified form Here is formula 1:	S. S. S. S. S. S. S. S. S. S.
Fig Prompt actu 下面我会给出两个公式,每个 构成,[P中的表示参数。 你需要判断我给出的两个公式 否相同,是否是同一个公式: 如果各个参数含义相同,是同一个 如果各个参数含义相同,是同一 對应关系,每个单元格内容是 两个公式的参数,第三行填写 下面是公式1: (公式1) 下面是公式2: (公式2) 通过表达含义判断,是否是同	ure 6: An examp ally used 公式由参数和运算符号 公式,只需要回答不是: 一个公式,则需要给出 行的表格来表示参数的 一个参数,前两行填写 统一后的公式参数。	le of deleted questions English trans I will give two formulas belov parameters and operation symb parameter. You need to judge whether the cc two formulas I gave have the san are the same formula: If the meaning is different, and t just answer no; If each pair of parameters have tf the same formula, the final form three-row table needs to be given relationship between the parameter parameters, fill in the unified form Here is formula 1; You will be the comparison of the same formula 1; You will be the comparison of the same formula 1; You will be the comparison of the same formula 1; You will be the comparison of the same formula 1; You will be the comparison of the same formula 1; You will be the comparison of the same formula 1; You will be the comparison of the same formula 1; You will be the comparison of the same formula 1; You will be the comparison of the same formula 1; You will be the comparison of the same formula 1; You will be the comparison of the same formula 1; You will be the same formula 1; You	S.
Fig Prompt actu 下面我会给出两个公式,每个 构成,[]中的表示参数。 你需要判断我给出的两个公式 否相同,是否是同一个公式: 如果含义不相同,不是同一个 如果含义不相同,不是同一个 如果含义不相同,不是同一个 如果含义不相同,不是同一个 就应关系,每个单元格内容是 两个公式的参数,第三行填写 下面是公式1: (公式 1) 下面是公式2: (公式 2) 通过表达含义判断,是否是同	ure 6: An examp ally used 公式由参数和运算符号 公式,只需要回答不是; 一个公式,则需要给出 行的表格来表示参数的 一个参数。前两行填写 统一后的公式参数。	le of deleted questions English trans I will give two formulas belor parameters and operation symb parameter. You need to judge whether the cc two formulas I gave have the sau are the same formula: If the meaning is different, and t just answer no; If each pair of parameters have tf the same formula the final form three-row table needs to be given relationship between the parameter parameters, fill in the unified form Here is formula 1: {formula 1} Here is formula 2: {formula 2}	S. Itation w. Each formula consists o ols. The text in [] represen orresponding parameters in the me meaning and whether they hey are not the same formula he same meaning, and they are ula needs to be given, and in to indicate the corresponding in the corresponding is are filled with two formulas una parameters in the third row
Fig Prompt actu 下面我会给出两个公式,每个 构成,□中的表示参数。 你需要判断我给出的两个公式 否相同,是否是同一个公式: 如果含义不相同,不是同一个 如果含义不相同,无是否是同一个 如果含义不相同,不是同一个 如果含义不相同,不是同一个 如果含义不相同,不是同一个 如果含义和的,子子。 而是公式2: 《公式1》 而是公式2: 《公式2) 通过表达含义判断,是否是同	ure 6: An examp ally used 公式由参数和运算符号 公式,只需要回答不是; 一个公式,则需要经出 行的表格来表示实数的 一个参数,前两行填写 统一后的公式参数。	le of deleted questions English trans I will give two formulas belov parameters and operation symb parameter. You need to judge whether the co two formulas I gave have the sau are the same formula: If the meaning is different, and t just answer no; If each pair of parameters have th the same formula, the final form three-row table needs to be given relationship between the parameter parameters, fill in the unifed form Here is formula 1; Here is formula 2: {formula 2; Judge whether they are the same formula 2: }	S. Ilation w. Each formula consists o ols. The text in [] represen orresponding parameters in the me meaning and whether they hey are not the same formula he same meaning, and they are ula needs to be given, and i to indicate the corresponding rs. The content of each cell is is a are filled with two formulas nula parameters in the third row comula by their meanings:
Fig Prompt actu 下面我会给出两个公式,每个 构成,0中的表示参数。 你需要判断我给出的两个公式 否相同,是否是同一个公式: 如果含义和相同,不是同一个 如果含个参数名义相同,是后一 最终的公式,并且给出的两个公式 为应关系,每个单元格内容是 两个公式的参数,第三行填至 下面是公式1: {公式 1} 下面是公式2: {公式 2} 通过表达含义判断,是否是同	ure 6: An examp ally used 公式由参数和运算符号 公式由参数和运算符号 公式,只需要回答不是; 一个公式,则需要给出 行的表格来表示参数的 一个参数,前两行填写 3统一后的公式参数。	le of deleted questions English trans I will give two formulas below parameters and operation symb parameter. You need to judge whether the co two formulas I gave have the sau are the same formula: If the meaning is different, and t just answer no; If each pair of parameters have th the same formula, the final form three-row table needs to be given relationship between the parameter parameter, fill in the unified form Here is formula 1: {formula 1} Here is formula 2: {formula 2} Judge whether they are the same formula	S. •Iation w. Each formula consists o ols. The text in [] represen presponding parameters in the me meaning and whether they hey are not the same formula he same meaning, and they ar ula needs to be given, and i to indicate the corresponding rs. The content of each cell is is a are filled with two formulas uula parameters in the third row formula by their meanings:
Fig Prompt actu 下面我会给出两个公式、每个 构成、[中的表示参数、 你需要判断我给出的两个公式、 如果含义不相同,不是同一个 如果含义不相同,不是同一个 如果含义不相同,不是同一个 如果含义不相同,不是同一个 如果含义不相同,不是同一个 如果含义不相同,不是同一个 如果含义和同,不是同一个 如果含义和同,不是同一个 如果含义和同,不是同一个 如果含义和同,不是同一个 如果含义和同,不是同一个 如果含义和同,不是同一个 如果含义和同,不是同一个 如果含义和同,无是是可。 《如子》 第二章 "你就是你了。" 》 第二章 "你就是你了。" 》 第二章 "你们" 》 第二章 "你们"	ure 6: An examp ally used 公式由参数和运算符号 公式由参数和运算符号 公式,只需要回答不是;)一个公式,则需要给出 ;行的表格来表示参数的 一个公式,可需要给出 ;行的表格来表示参数的 一个公式;)一个公式:)一个公式:	le of deleted questions English trans I will give two formulas belor parameters and operation symb parameters and operation symb If each pair of parameters have th the same formula 1 gave have the same formula 1: {formula 1} Here is formula 2: {formula 2} Judge whether they are the same formula semantic-based mergin	S. Jation w. Each formula consists o ols. The text in [] represen orresponding parameters in the me meaning and whether they hey are not the same formula he same meaning, and they ar nula needs to be given, and he same meaning, and they ar nula needs to be given, and is are filled with two formulas nula parameters in the third row cornula by their meanings: ng.
Prompt actu இன்னு பிர்க்கு இன்னு பிர்க்கு <	ure 6: An examp ally used 公式由参数和运算符号 公式由参数和运算符号 公式,只需要回答不是;)—个公式,则需要回答不是;)—个公式,则需要给出 之行的表格来表示参数的 是一个参数,前两行填写 3统一后的公式参数。 3—一个公式:	le of deleted questions English trans I will give two formulas belor parameters and operation symb parameter. You need to judge whether the cy two formulas I gave have the sau are the same formula: If the meaning is different, and t just answer no; If each pair of parameters have th the same formula, the final form three-row table needs to be given relationship between the parameter parameters, fill in the unified form Here is formula 1: {formula 1} Here is formula 2: {formula 2} Judge whether they are the same form	S. Ilation w. Each formula consists o ols. The text in [] represen orresponding parameters in the me meaning and whether they hey are not the same formula he same meaning, and they ar ula needs to be given, and i; hey are not the same formula so are filled with two formulas ula parameters in the third row formula by their meanings: ang.
Fig Prompt actu 下面我会给出两个公式、每个构成、[]中的表示参数、(新二人前一次) 你需要判断我会出的两个公式 否相同、是否是同一个公式: 如果各个条数含义相同,不是同一个 如果各个参数含义相同,不是同一个 如果各个参数含义相同,不是同一个 如果各个参数含义相同,不是同一个 如果各个参数含义相同,不是同一个 如果各个参数含义相同,不是同一个 如果各个参数含义相同,不是同一个 如是名人参数含义相同,是后 最终的公式,并且给出一个当 对应关系、每个单元格内容是 不同是公式1: 《公式1) 下面是公式2: 通过表达含义判断,是否是同一	ure 6: An examp ally used 公式由参数和运算符号 公式,只需要回答不是; 一个公式,则需要的当本是; 一个公式,前两行填写 统一后的公式参数。 一个公式: 一个公式:	le of deleted questions English trans I will give two formulas belor parameters and operation symb parameter. You need to judge whether the cc two formulas I gave have the sau are the same formula: If the meaning is different, and t just answer no; If each pair of parameters have ti the same formula 1: If each pair of parameters have ti the same formula the final form three-row table needs to be given relationship between the parameter parameter, and the first two rows Parameter, and the first two rows Parameter, and the first two rows Parameter, and It he sint formula 1: {formula 1} Here is formula 2: {formula 2} Judge whether they are the same fi	S. Ilation w. Each formula consists o ols. The text in [] represen orresponding parameters in the me meaning and whether they hey are not the same formula he same meaning, and they arr he same meaning, and they arr to indicate the corresponding rs. The content of each cell is : a are filled with two formulas nula parameters in the third row ormula by their meanings: Thg.
Fig Prompt actu 下面我会给出两个公式、每个构成、[]中的表示参数、 你需要判断我给出的两个公式 否相同,是否是同一个公式: 如果各个多数含义相同,不是同一个 如果各个参数含义相同,不是同一个 如果各个参数含义相同,不是同一个 如果各个参数含义相同,不是同一个 如果各个参数含义相同,不是同一个 如果各个参数含义相同,不是同一个 如果各个参数含义相同,不是同一个 如果各个参数含义相同,不是同一个 如果各个参数含义相同,不是同一个 如果各个参数含义相同,不是同一个 如是名人参数含义相同,不是同一个 如是公式:: 《公式 1) 下面是公式2: 通过表达含义判断,是否是同 好了。 好了。 通过表达含义判断,是否是同一个 Figure	ure 6: An examp ally used 公式由参数和运算符号 公式,只需要回答不是; 一个公式,则需要给出 后的表裙来表示参数的 一个参数,前两行填写 统一后的公式参数。 一个公式:	le of deleted questions English trans I will give two formulas belor parameters and operation symb parameter. You need to judge whether the cc two formulas I gave have the sau are the same formula: If the meaning is different, and t just answer no; If each pair of parameters have th the same formula it fif each pair of parameters have th the same formula the final form three-row table needs to be given relationship between the parameter parameter, and the first two rows Parameters, fill in the unified form Here is formula 1: {formula 1} Here is formula 2: {formula 2} Judge whether they are the same f	S. S. S. S. S. S. S. S. S. S.

1081 1082 1083 Prompt actually used **English translation** 1084 -个初中物理题目,根据问题给出计算的过程 这是-This is a junior high school physics question. Based on the given 1086 让我们一步一步地地思考,在最后用"###"作为开始 question, provide the calculation process and let's think step by 给出最终答案(一个数字)和答案的单位。 step. Finally, use "###" to start giving the final answer (a number) 1087 and the unit of the answer. 1088 Question: {{问题}} Answer: Question: {{question}} 1089 Answer: 1090 1091 (a) Zero-shot prompt for LLMs. **Prompt actually used English translation** 1093 这是--个初中物理题目,根据问题给出计算的过程, This is a junior high school physics question. Based on the given 1094 用公式表示。 question, provide the calculation process. 1095 Question: {{样例1问题}} Question: {{question of example 1}} Answer: {{样例1解析}} Answer: {{explanation of example 1}} 1097 1098 ...omitted... ...omitted... 1099 Question: {{问题}} Question: {{question}} 1100 Answer: Answer: 1101 1102 (b) Few-shot prompt for LLMs. 1103 Figure 8: Zero-shot and few-shot prompts for LLMs. 1104 1105 1106 1107 1108 1109 1110 1111 1112 Prompt actually used **English translation** 1113 这是一个初中物理题目,根据问题给出计算的过 This is a junior high school physics question. Based on the given question, provide the calculation process. 程,用公式表示。 1114 1115 可能用到的公式有: {{top 5检索到的公式}} The formulas that may be used include: {{top 5 retrieved formulas}} Question: {{问题}} Question: {{question}} 1116 Answer: Answer: 1117 (a) Few-shot prompt for LLMs with formula retriever. 1118 1119 Prompt actually used English translation 1120 这是一个初中物理题目,根据问题给出计算的过 This is a junior high school physics question. Based on the given 程,用公式表示。 question, provide the calculation process 1121 1122 可能用到的公式有: {{用到的公式}}} The formulas that may be used include: {{used formulas}} Question: {{样例1问题}} Question: {{question of example 1}} 1123 Answer: {{样例1解析}} Answer: {{explanation of example 1}} 1124 ...omitted... ...omitted... 1125 可能用到的公式有: {{top 5检索到的公式}} The formulas that may be used include: {{top 5 retrieved formulas}} 1126 Question: {{问题}} Question: {{question}} 1127 Answer: Answer: 1128 (b) Zero-shot prompt for LLMs with formula retriever. 1129 1130 Figure 9: Zero-shot and few-shot prompts for LLMs with formula retriever. 1131 1132 1133