
FormulaReasoning: A Dataset for Formula-Based Numerical Reasoning

Xiao Li Bolin Zhu Sichen Liu Yin Zhu Yiwei Liu Gong Cheng*
State Key Laboratory for Novel Software Technology, Nanjing University, Nanjing, China
{xiaoli.nju, bolinzhu, sichenliu, yinzhu, ywliu}@smail.nju.edu.cn
gcheng@nju.edu.cn

Abstract

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

14 1 Introduction

15 Numerical reasoning constitutes one of the significant forms within natural language reason-
16 ing (Frieder et al., 2023). The study of numerical reasoning has seen substantial progress in recent
17 years, largely driven by the development of LLMs (OpenAI, 2023; Touvron et al., 2023; Li et al.,
18 2023c) and specialized datasets (Wang et al., 2017; Dua et al., 2019; Amini et al., 2019; Cobbe
19 et al., 2021a). Current datasets for numerical reasoning typically include simple, commonsense
20 numerical questions that do not reflect the complexity of real-world problems. These datasets have
21 not fully addressed the interpretability issue in numerical reasoning, as they often rely on implicit
22 commonsense knowledge without explicit guidance knowledge during the reasoning process. This
23 issue becomes particularly evident when LLMs meet hallucination (Frieder et al., 2023; Bang et al.,
24 2023). Consequently, one might naturally ask “*What knowledge could I use to guide numerical
25 reasoning process?*”. Formulas exactly represent such knowledge that has been largely overlooked in
26 research but is frequently utilized in real-life applications.

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

*Corresponding author

Question

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 $C_{\text{water}} = 4.2 \times 10^3 \text{ J}/(\text{kg} \cdot ^\circ\text{C})$.
Q: If the total electrical energy consumed during the heating process is $1 \times 10^7 \text{ J}$, what is the thermal efficiency of the water heater?

Explanation (Reasoning Steps)

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

32 Recently, Liu et al., 2023 constructed two formula-based datasets, Math23K-F and MAWPS-F. How-
33 ever, the formulas in these datasets primarily consist of commonsense formulas (such as total_amount
34 $= \text{unit_amount} \times \text{total_number}$), and only 33.5% and 38.4% of the questions in these datasets,
35 respectively, require the use of formulas.
36 To address this gap, we constructed a dataset for numerical reasoning that requires the use of formulas
37 called FormulaReasoning. We annotated formulas for each question in FormulaReasoning. An exam-
38 ple of FormulaReasoning is shown in Figure 1.² The formula-based feature makes FormulaReasoning
39 a more challenging dataset for developing systems that can tackle real-world numerical reasoning
40 problems. Indeed, in fields such as mathematics and physics, formulas serve as an important vessel for
41 representing domain knowledge. However, existing datasets scarcely consider explicit incorporation
42 of formulas into numerical reasoning.

Dataset	Math23K-F	MAWPS-F	GSM8K	FormulaReasoning
# questions	23,162	2,373	8,792	5,420
# formulas (and variants)	51 (131)	18 (46)	0 (0)	272 (824)
# questions requiring formula (proportion)	7,750 (33.46%)	911 (38.39%)	N/A	5,420 (100%)
Avg. # reasoning steps	1.16	1.01	3.59	2.37

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

43 We collected questions requiring formula-based numerical reasoning from Chinese junior high
44 school physics examinations. With the *combined efforts of manual annotation and assistance from*
45 *LLMs*, we annotated each question with an explanation text, a final answer, and a set of relevant
46 formulas (including formula structures, parameter names, symbols, numerical values, and units) and
47 built a *formula database*. The formula database functions as an external knowledge base, which can
48 be used to evaluate retrieval-based/augmented systems. In Table 1, we compare FormulaReasoning
49 with two existing formula-based datasets and the well-known GSM8K. In comparison to Math23K-F
50 and MAWPS-F, FormulaReasoning contains a *larger number of formulas* (272), whereas the other
51 two datasets contain 51 and 18 formulas. Additionally, all questions in FormulaReasoning require

²Please note that FormulaReasoning is in Chinese. For the convenience of understanding, we translated Chinese into English in all the examples presented in this paper.

52 the use of formulas. The *higher average number of reasoning steps* (2.37 vs. 1.16/1.01) implies
53 that FormulaReasoning is more challenging and better suited for evaluating existing models as a
54 multi-step formula-based reasoning task.

55 We used FormulaReasoning to evaluate LLMs ranging from 7B to >100B parameters, as well as
56 fine-tuned models such as Qwen-1.8B (Bai et al., 2023) and ChatGLM-6B (Zeng et al., 2022) with
57 a proposed Chain-of-Thought supervised fine-tuned method and a data augmentation method. We
58 also trained an encoder for formula retrieval and experimented with retrieval-augmented generative
59 models. Our empirical findings show that the best existing models only achieve an accuracy of around
60 74%, lagging behind an accuracy 92% of humans, indicating that there is still significant room for
61 exploration in formula-based numerical reasoning.

62 Our contributions are summarized as follows:

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

69 The dataset is available on <https://zenodo.org/doi/10.5281/zenodo.11408109> under
70 the CC BY 4.0 License and our code is available on [https://github.com/nju-websoft/
71 FormulaReasoning](https://github.com/nju-websoft/FormulaReasoning) under the Apache License 2.0.

72 **2 Related Work**

73 **2.1 Numerical Reasoning Datasets**

74 Numerical reasoning is one of the fundamental capabilities of natural language reasoning. The
75 study of numerical reasoning in natural language has existed for several years. Numerous datasets,
76 such as DROP (Dua et al., 2019), GSM8K (Cobbe et al., 2021b), TSQA (Li et al., 2021) and
77 MATH (Hendrycks et al., 2021), have introduced natural language numerical reasoning. Another line
78 of research focusing on numerical reasoning in natural language is math word problem (MWP). MWP
79 tasks typically provide a short passage (i.e., a question) and require the generation of an arithmetic
80 expression that can compute an answer. Representative datasets include MAWPS (Koncel-Kedziorski
81 et al., 2016), Math23K (Wang et al., 2017), MathQA (Amini et al., 2019), etc.

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

89 **2.2 Numerical Reasoning Methods**

90 The methods for solving numerical reasoning have evolved from statistical approaches (Hosseini
91 et al., 2014; Kushman et al., 2014) to those based on rules and templates (Shi et al., 2015; Wang et al.,
92 2019) and further to methods based on deep learning models (Gupta et al., 2019; Chen et al., 2022;
93 Kim et al., 2022; Li et al., 2023a). In the past two years, with the rapid development of LLMs, LLMs
94 have demonstrated strong capabilities in resolving numerical reasoning questions. Consequently,
95 several methods aimed at enhancing the reasoning abilities of LLMs have been proposed, including
96 the notable Chain of Thoughts (CoTs) method (Wei et al., 2022), along with many subsequent variant
97 approaches (Kojima et al., 2022; Wang et al., 2022; Zhou et al., 2022; Li et al., 2023b).

98 We established representative existing methods as baselines for FormulaReasoning, including
 99 zero/few-shot CoTs prompting methods to LLMs ranging from 7B to over 100B parameters. We
 100 trained a specialized formula retriever for retrieving formulas and explored retrieval-enhanced numer-
 101 ical reasoning. We also divided the reasoning process into formula generation, parameter extraction,
 102 and calculation, and used data augmentation to enhance fine-tuned models with fewer than 7B
 103 parameters.

104 3 Dataset Construction

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

114 3.1 Preprocessing

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

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

Original explanation.

The change in water temperature is $60 - 20 = 40$ °C. Therefore, the heat absorbed by the water is $Q_{\text{absorbed}} = 50 \text{ kg} \times 4.2 \times 10^3 \text{ J/(kg}\cdot\text{°C)} \times 40 \text{ °C} = 8.4 \times 10^6 \text{ J}$. 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: $\eta = Q_{\text{absorbed}} / W_{\text{total}} \times 100\% = (8.4 \times 10^6 \text{ J}) / (1.0 \times 10^7 \text{ J}) \times 100\% = 84\%$. 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: $[\text{Degree of water temperature increase}] = [\text{Final temperature}] - [\text{Initial temperature}] = 60 \text{ °C} - 20 \text{ °C} = 40 \text{ °C}$. The degree of water temperature increase = 40 °C.
 2. Calculating the heat absorbed by water: $[\text{Heat absorbed by water}] = [\text{Mass of water}] \times [\text{Specific heat capacity of water}] \times [\text{Degree of water temperature increase}] = 50 \text{ kg} \times 4.2 \times 10^3 \text{ J/(kg}\cdot\text{°C)} \times 40 \text{ °C} = 8400000 \text{ J}$. The heat absorbed by water = 8400000 J.
 3. The thermal efficiency of the water heater can be obtained from: $[\text{Thermal efficiency of the water heater}] = [\text{Heat absorbed by water}] / [\text{Total electrical energy consumed}] \times 100\% = 8400000 \text{ J} / (1 \times 10^7 \text{ J}) * 100\% = 84\%$. The thermal efficiency of the water heater = 84%.
- Answer = 84%

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

125 3.2 Formula Normalization

126 We found that the reasoning steps (i.e. the explanation) in the obtained raw dataset lacked a normalized
 127 format and were expressed quite casually. Some formulas mixed parameter names (e.g., “mass of

128 water”) and symbols (e.g., “ m_{water} ”), while others simply provided calculations in numerical form
 129 without parameter names or symbols. In order to ensure that all explanations adopted a normalized
 130 form of formulas, we normalized the formula annotations in the explanations. An example can
 131 be found in Table 2. In this process, we need to *identify the formulas used within the original*
 132 *explanations* and to *correct any formatting issues*. Manually undertaking such tasks would require
 133 significant effort. However, since the process is not open-ended, but rather structured and verifiable,
 134 we could automatically, e.g., *using a LLM*, extract formulas from the explanations, calculate each step,
 135 and compare the result with the given answer to ensure the accuracy of this normalization process.

136 Specifically, to enhance the efficiency of the annotation, we adopted a coarse-to-fine annotation
 137 approach with the help of a LLM³. We first prompted the LLM in a few-shot manner to generate
 138 accurate explanations of the reasoning process. Then, we used few-shot prompts to guide the LLM in
 139 correcting minor errors within the normalized explanations, including formatting errors in formula
 140 annotations and inaccuracies in the parameters used during computations. Both prompts can be found
 141 in Appendix C.1.1. Next, we will provide a detailed description of this process.

142 Initially, we introduced the question along with its original explanation and the corresponding answer
 143 to guide the LLM through few-shot prompting to revise the original explanation. We observed that
 144 the ability of the LLM to revise explanations towards normalized explanations remained satisfactory.
 145 To assess the correctness of the revised explanations, we extracted formulas from these explanations
 146 and then computed the answer using the numbat tool⁴. In addition to providing explanations, we also
 147 required the LLM to present the values, symbols, and units of each parameter in the formulas in the
 148 form of a table. An example is shown in Figure 1.

149 At this stage, we checked the correctness of the formula format in the explanations by automatic rules,
 150 including whether there were omissions in parameter names, parameter symbols, or corresponding
 151 units, and these issues were all correctable. Therefore, if our program detected that the LLM had not
 152 successfully generated an accurate normalized explanation, we used few-shot prompting to identify
 153 and correct these specific errors. More details can be found in Appendix C.1.1. We observed that
 154 the questions which remained incorrect despite multiple attempts by the LLM were of notably poor
 155 quality, including missing important reasoning steps, unclear question formulation, and so on. Some
 156 examples of these questions can be found in Appendix C.1.2. These questions were removed from
 157 our dataset. Following this step, our dataset contains a remaining total of 5,420 questions.

158 3.3 Formula Database Construction

159 Our next step was to *construct a unified formula*
 160 *database for the entire dataset*. Given that pa-
 161 rameters in the same formula can be expressed
 162 differently across various problem contexts, for
 163 instance, the two formulas “[weight of water]
 164 = [mass of water] * [gravitational acceleration]”
 165 and “[weight] = [mass] * [gravitational acceler-
 166 ation]” both calculate the weight of an object,
 167 we need to merge these formulas into a single
 168 representation.

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

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

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

³During dataset construction, we accessed Qwen-max via API (<https://help.aliyun.com/zh/dashscope/developer-reference/quick-start>). Qwen-max is a LLM with over 100B parameters and a strong capability in Chinese.

⁴<https://numbat.dev>. Numbat is designed for scientific computations with support for physical units.

173 **Symbolic rules based merging.** In this step, we merged formulas through symbolic rules. Specif-
174 ically, this was achieved by *comparing the structure of the formulas and the symbols*. Take the
175 following as an example of judging whether two formulas have the same structure: the formulas
176 “ $f_1 : a_1=(b_1+c_1)/d_1$ ”, “ $f_2 : a_2=(b_2+c_2)/d_2$ ” and “ $f_3 : b_1=a_1*d_1-c_1$ ” have the same structure because
177 f_2 can be derived from f_1 by renaming parameters, and f_3 can be obtained from f_1 by transformation.
178 Moreover, in physics, certain physical quantities are conventionally represented by specific symbols.
179 For example, the mass of an object is often denoted by “ m ” and the density of an object is frequently
180 represented by the symbol “ ρ ”. Subscripts are then used to distinguish which specific object a
181 physical quantity refers to, such as “ ρ_{water} ” for the density of water. For any two formulas, we first
182 computed all the transformations of each formula to obtain a set of all its variants. Then, we compared
183 the formula structures in the two sets to determine if two formulas were structurally equivalent. If
184 they shared the same structure, we then compared whether their symbols, with subscripts removed,
185 were identical. If they were, we considered these two formulas to be mergeable. When merging, we
186 retained the parameter with the shorter length from the two. After merging based on symbolic rules,
187 we reduced the number of formulas in the formula database from 12,906 to 1,163.

188 **Semantic-based merging.** In the symbolic rules based merging process, the semantic information
189 of the parameter names was neglected. This led us to *perform merges grounded on the semantics*
190 *of the parameter names*. For instance, two formulas that were not merged during the symbolic
191 fusion stage, “[density] = [mass] / [volume]” and “[density of water] = [mass of water] / [volume

192 of water]”, can actually be merged. We would carry out the merging of these two formulas based
193 on the semantic information of the parameter names (for example, “density” and “density of water”
194 are semantically similar). Specifically, for formulas with identical structures, we tokenized each
195 pair of corresponding parameters to create two sets of words⁵. When the two sets overlapped, the
196 parameters were considered to have semantic connection, and the formulas became candidates for
197 merging. Utilizing this approach, we identified a set of pairs of potentially mergeable formulas
198 and then consulted the LLM for a thorough evaluation of each pair. The prompts can be found in
199 Appendix C.1.3. After this step, the number of formulas in the formula database was reduced to 439.

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

205 4 Experiments Setup

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

213 4.1 Dataset Split

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

⁵We used jieba: <https://github.com/fxsjy/jieba>.

219 **4.2 Evaluation**

220 **4.2.1 Human Performance**

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

225 **4.2.2 LLMs**

226 Following Kojima et al., 2022, we incorporated the phrase “Let’s think step by step” into the zero-shot
227 prompt to guide LLMs in generating the reasoning steps. For the few-shot setting, we randomly
228 sampled five questions from the training set to serve as examples for in-context learning. Each
229 example includes the question text and the reasoning steps (i.e., the explanation). Examples of the
230 prompts can be found in Appendix C.4.1.

231 We conducted experiments on GPT-4-turbo, GPT-3.5-turbo, GLM4, and Qwen-max, with each of
232 these models having over 100 billion parameters. We also evaluated on Llama2-7B (Touvron et al.,
233 2023), Llama3-8B (Meta, 2024), Qwen-7B/14B (Bai et al., 2023), InternLM2-7B/20B (Team, 2023),
234 ChatGLM3-6B (Zeng et al., 2022), including the base and chat versions of these models. We followed
235 the common practice that few-shot experiments were performed on the base versions, while zero-shot
236 experiments were conducted on the chat or instruct versions.

237 **4.2.3 Formula Retriever**

238 We trained a formula retriever on the training set. Specifically, we encoded each question using the
239 Chinese-BERT-wwm-base (Devlin et al., 2019; Cui et al., 2021) model to obtain the CLS vector of
240 the question. Each formula in the formula database was represented by a randomly initialized vector.
241 During training, we calculated the cosine score between the question vector and the formula vector.
242 The retriever was then trained with in-batch negatives and contrastive learning loss (Gao et al., 2021).
243 Subsequently, for each question in the id test, we retrieved the top five formulas with the highest
244 scores and included them in the prompt to observe the change in the performance of the LLM when
245 provided with relevant formulas. More details can be found in Appendix C.4.2.

246 **4.2.4 Supervised Fine-tuned Models**

247 We found that directly prompting models possessing fewer than 7B parameters failed to produce
248 satisfactory outcomes (for example, ChatGLM3-6B attained merely 8.99 points in a zero-shot setting).
249 Therefore, we conducted supervised fine-tuning of models with fewer than 7B parameters, yet
250 discerned that, dissimilar to larger models (such as GPT-4-turbo), smaller models did not exhibit
251 proficient performance in numerical extraction and calculation. In order to augment the reasoning
252 capabilities of smaller models, we explored two approaches for improvement.

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

259 **Data Augmentation (DA)** We augmented the training dataset with the assistance of larger models.
260 Firstly, we utilized a few-shot approach to prompt the LLM (Qwen-max) to generate new question-
261 answer pairs. The correctness of the computation process generated by the LLM was meticulously
262 verified using a calculator. Subsequently, the formulas generated by the model were extracted and
263 normalized. More details could be found in Appendix C.3.1.

264 **4.3 Metric**

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

268 **5 Experiments Results**

269 In this section, we presented the experimental results and analysis. Due to space constraints, the error
 270 analysis can be found in Appendix C.2 and the implementation details can be found in Appendix C.4.

271 **5.1 Human Performance**

272 In FormulaReasoning, humans achieved impressive performance, with a score of 93.49 on the id test,
 273 90.47 on the ood test, and an average score of 92.03.

274 **5.2 Results of LLMs**

Model	Size	zero-shot CoT			few-shot CoT		
		id test	ood test	Avg.	id test	ood test	Avg.
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
GLM4	>100B	65.32	65.22	65.27	62.47	65.98	64.16
Qwen-max	>100B	58.67	57.80	58.25	58.91	63.94	61.33
InternLM*	20B	5.70	4.60	5.17	18.29	11.25	14.90
Qwen*	14B	32.07	37.60	34.73	44.89	36.83	41.01
Llama3*	8B	26.66	17.98	20.41	12.81	8.87	10.91
Llama2*	7B	0.00	0.26	0.13	1.43	0.26	0.87
Qwen*	7B	7.36	8.70	8.01	21.14	18.16	19.71
InternLM*	7B	7.84	7.67	7.76	9.50	8.18	8.86
ChatGLM3*	6B	9.36	8.62	8.99	23.89	19.95	21.92
Human	-	93.49	90.47	92.03	93.49	90.47	92.03

Table 4: Results of LLMs with zero-shot and few-shot prompting. * indicates that the chat or instruct version of the model was used in the zero-shot setting, while the base version of the model was used in the few-shot setting.

275 The evaluation results on LLMs are shown in Table 4. *GPT-4-turbo exhibited the best performance*
 276 *in both zero-shot and few-shot settings*, surpassing the second-ranked GLM4 by an average of 6.16
 277 points in zero-shot setting and 10.22 in few-shot setting. Among models with size not exceeding
 278 20B, Qwen-14B demonstrated commendable performance in both zero-shot and few-shot settings.
 279 The subpar performance of Llama2 might be due to its pre-training data being primarily in English.
 280 We also conducted few-shot testing on the chat version of LLMs with size not exceeding 20B,
 281 and the results can be found in Appendix C.4.3. After incorporating few-shot examples, GPT-4-
 282 turbo, GPT-3.5-turbo and Qwen-max demonstrated performance improvements, ranging from 0.24
 283 to 6.14. However, similar performance changes were not observed on GLM4, possibly due to its
 284 supervised fine-tuning and alignment with human preferences which enhanced GLM4’s understanding
 285 of instructions but probably also compromised its in-context learning ability.

286 Human performance surpassed the performance of few-shot GTP-4-turbo on the id and ood tests by
 287 margins of 21.99 and 13.25 points, respectively. Such results demonstrated that there remained a
 288 substantial gap between the current capabilities of state-of-the-art LLMs and human performance.
 289 This was even more pronounced when considering smaller-scale models. These findings underscored
 290 *the challenging nature of FormulaReasoning as an unresolved dataset*, and that there was significant
 291 room for improvement in LLMs as they struggled to match human levels of reasoning.

292 **5.3 Results of LLMs with Formula Retriever**

293 The results of LLMs utilizing the formula retriever are shown
 294 in Table 5. We found that the impact on performance varied
 295 among different LLMs when incorporating retrieved formulas
 296 into prompts. We observed a positive enhancement on GLM4,
 297 with score increments of 4.99 and 3.33 with zero-shot and
 298 few-shot, respectively, on the id test. However, we observed
 299 a performance decline with GPT-4-turbo. Specifically, we
 300 found that the top 5 retrieved formulas often included irrele-
 301 vant ones, as the number of formulas required varies for each
 302 problem. The presence of these extraneous formulas affected the model’s performance, indicating
 303 that there is considerable room for further research in utilizing a formula database.

Model	zero-shot	few-shot
GLM4	65.32	62.47
+ formula retriever	70.31	65.80
GPT-4-turbo	70.07	71.50
+ formula retriever	68.17	67.00

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

304 5.4 Results of Supervised Fine-tuned Models

305 Table 6 shows the results for the supervised fine-tuned
 306 models, with and without CoT-SFT and DA, which were
 307 detailed in Section 4.2.4. In most settings, both models
 308 achieved higher scores on the id test than the ood test, yet
 309 they still exhibited considerable performance on the ood
 310 test. This indicates that 1) *the ood formulas indeed im-*
 311 *pacted model performance* and 2) *the models still demon-*
 312 *strate generalizability*. We hope that the division of id test
 313 and ood test will be helpful for assessing the generalization
 314 ability of fine-tuned models in future works.

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

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

315 It was noteworthy that with CoT-SFT, Qwen-1.8B and
 316 ChatGLM3-6B, with a mere parameter count of 1.8B and 6B, respectively, achieved performance
 317 comparable to GPT-4-turbo (though such a comparison may not be entirely fair). This indicated that
 318 the incorporation of CoT-SFT and the use of calculators could significantly enhance the reasoning
 319 capabilities of small models. Our findings revealed that *focusing on reasoning with CoT while*
 320 *delegating numerical calculation to a calculator could enhance the performance of small models,*
 321 *given their limited calculating capability. The assistance of LLMs for data augmentation could also*
 322 *enhance smaller models’ reasoning capability.* This discovery provides valuable insights for future
 323 deployment of numerical reasoning systems on small models.

324 6 Conclusion and Limitations

325 We introduced FormulaReasoning, a dataset for formula-based numerical reasoning. We annotated
 326 the reasoning steps with formulas for each question with both manual and LLM-assisted efforts.
 327 Furthermore, we constructed a formula database after merging formulas with similar meanings,
 328 serving as an external knowledge base for subsequent retrieval-based/augmented approaches. We
 329 evaluated FormulaReasoning across various sizes of LLMs, supervised fine-tuned models, and
 330 retrieval-augmented LLMs, demonstrating its challenging nature as an unresolved task. Our findings
 331 indicate substantial room for improvement of existing models on formula-based numerical reasoning,
 332 thus motivating future research efforts.

333 We have also translated the dataset into English unitizing LLMs. However, we have not yet accurately
 334 assessed the quality of the translated dataset. At present, we have not released the English version
 335 of the dataset, but we will do so later after ensuring the quality of the English dataset. Additionally,
 336 our dataset is limited to the domain of physics. Although junior high school physics is not overly
 337 complex and can be understood by most people, it is still possible to explore formula-based question
 338 answering data in other domains.

339 Acknowledgments and Disclosure of Funding

340 This work was supported by the CIPSC-SMP-Zhipu.AI Large Model Cross-Disciplinary Fund. We
341 are grateful to Chao Li for his suggestions and efforts in the annotation.

342 References

- 343 Aida Amini, Saadia Gabriel, Shanchuan Lin, Rik Koncel-Kedziorski, Yejin Choi, and Hannaneh
344 Hajishirzi. MathQA: Towards interpretable math word problem solving with operation-based
345 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)*, pages 2357–2367, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1245. URL <https://aclanthology.org/N19-1245>.
- 349 Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenheng Ge,
350 Yu Han, Fei Huang, Binyuan Hui, Luo Ji, Mei Li, Junyang Lin, Runji Lin, Dayiheng Liu, Gao
351 Liu, Chengqiang Lu, K. Lu, Jianxin Ma, Rui Men, Xingzhang Ren, Xuancheng Ren, Chuanqi
352 Tan, Sinan Tan, Jianhong Tu, Peng Wang, Shijie Wang, Wei Wang, Shengguang Wu, Benfeng Xu,
353 Jin Xu, An Yang, Hao Yang, Jian Yang, Jian Yang, Shusheng Yang, Shusheng Yang, Bowen Yu,
354 Yu Bowen, Hongyi Yuan, Zheng Yuan, Jianwei Zhang, Xing Zhang, Yichang Zhang, Zhenru Zhang,
355 Chang Zhou, Jingren Zhou, Xiaohuan Zhou, and Tianhang Zhu. Qwen technical report. *ArXiv*,
356 abs/2309.16609, 2023. URL <https://api.semanticscholar.org/CorpusID:263134555>.
- 357 Yejin Bang, Samuel Cahyawijaya, Nayeon Lee, Wenliang Dai, Dan Su, Bryan Wilie, Holy Lovenia,
358 Ziwei Ji, Tiezheng Yu, Willy Chung, Quyet V. Do, Yan Xu, and Pascale Fung. A multitask,
359 multilingual, multimodal evaluation of chatgpt on reasoning, hallucination, and interactivity. *ArXiv*,
360 abs/2302.04023, 2023. URL <https://api.semanticscholar.org/CorpusID:256662612>.
- 361 Jiayi Chen, Xiao-Yu Guo, Yuan-Fang Li, and Gholamreza Haffari. Teaching neural module networks
362 to do arithmetic. In *Proceedings of the 29th International Conference on Computational Linguistics*,
363 pages 1502–1510, Gyeongju, Republic of Korea, October 2022. International Committee on
364 Computational Linguistics. URL <https://aclanthology.org/2022.coling-1.129>.
- 365 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser,
366 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John
367 Schulman. Training verifiers to solve math word problems. *ArXiv*, abs/2110.14168, 2021a. URL
368 <https://api.semanticscholar.org/CorpusID:239998651>.
- 369 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser,
370 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to solve
371 math word problems. *arXiv preprint arXiv:2110.14168*, 2021b.
- 372 Yiming Cui, Wanxiang Che, Ting Liu, Bing Qin, and Ziqing Yang. Pre-training with whole word
373 masking for chinese bert. *IEEE Transactions on Audio, Speech and Language Processing*, 2021. doi:
374 10.1109/TASLP.2021.3124365. URL <https://ieeexplore.ieee.org/document/9599397>.
- 375 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep
376 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)*, pages 4171–4186, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1423. URL <https://aclanthology.org/N19-1423>.
- 381 Bosheng Ding, Chengwei Qin, Ruochen Zhao, Tianze Luo, Xinze Li, Guizhen Chen, Wenhan Xia,
382 Junjie Hu, Anh Tuan Luu, and Shafiq Joty. Data Augmentation using LLMs: Data Perspectives,
383 Learning Paradigms and Challenges. *arXiv e-prints*, art. arXiv:2403.02990, March 2024. doi:
384 10.48550/arXiv.2403.02990.

- 385 Dheeru Dua, Yizhong Wang, Pradeep Dasigi, Gabriel Stanovsky, Sameer Singh, and Matt Gardner.
386 DROP: A reading comprehension benchmark requiring discrete reasoning over paragraphs. In
387 *Proceedings of the 2019 Conference of the North American Chapter of the Association for Compu-*
388 *tational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages
389 2368–2378, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi:
390 10.18653/v1/N19-1246. URL <https://aclanthology.org/N19-1246>.
- 391 Simon Frieder, Luca Pinchetti, Ryan-Rhys Griffiths, Tommaso Salvatori, Thomas Lukasiewicz,
392 Philipp Christian Petersen, Alexis Chevalier, and J J Berner. Mathematical capabilities of chat-
393 gpt. *ArXiv*, abs/2301.13867, 2023. URL [https://api.semanticscholar.org/CorpusID:](https://api.semanticscholar.org/CorpusID:256415984)
394 256415984.
- 395 Tianyu Gao, Xingcheng Yao, and Danqi Chen. SimCSE: Simple contrastive learning of sentence
396 embeddings. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language*
397 *Processing*, pages 6894–6910, Online and Punta Cana, Dominican Republic, November 2021.
398 Association for Computational Linguistics. doi: 10.18653/v1/2021.emnlp-main.552. URL [https:](https://aclanthology.org/2021.emnlp-main.552)
399 [//aclanthology.org/2021.emnlp-main.552](https://aclanthology.org/2021.emnlp-main.552).
- 400 Nitish Gupta, Kevin Lin, Dan Roth, Sameer Singh, and Matt Gardner. Neural module networks for
401 reasoning over text. In *International Conference on Learning Representations*, 2019.
- 402 Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn
403 Song, and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset. In
404 *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track*
405 *(Round 2)*, 2021.
- 406 Mohammad Javad Hosseini, Hannaneh Hajishirzi, Oren Etzioni, and Nate Kushman. Learning to
407 solve arithmetic word problems with verb categorization. In *Proceedings of the 2014 Conference*
408 *on Empirical Methods in Natural Language Processing (EMNLP)*, pages 523–533, Doha, Qatar,
409 October 2014. Association for Computational Linguistics. doi: 10.3115/v1/D14-1058. URL
410 <https://aclanthology.org/D14-1058>.
- 411 Edward J Hu, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen,
412 et al. Lora: Low-rank adaptation of large language models. In *International Conference on*
413 *Learning Representations*, 2021.
- 414 Jeonghwan Kim, Junmo Kang, Kyung-min Kim, Giwon Hong, and Sung-Hyon Myaeng. Exploiting
415 numerical-contextual knowledge to improve numerical reasoning in question answering. In
416 *Findings of the Association for Computational Linguistics: NAACL 2022*, pages 1811–1821,
417 Seattle, United States, July 2022. Association for Computational Linguistics. doi: 10.18653/v1/
418 2022.findings-naacl.138. URL <https://aclanthology.org/2022.findings-naacl.138>.
- 419 Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large
420 language models are zero-shot reasoners. *Advances in neural information processing systems*, 35:
421 22199–22213, 2022.
- 422 Rik Koncel-Kedziorski, Subhro Roy, Aida Amini, Nate Kushman, and Hannaneh Hajishirzi. MAWPS:
423 A math word problem repository. In *Proceedings of the 2016 Conference of the North American*
424 *Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages
425 1152–1157, San Diego, California, June 2016. Association for Computational Linguistics. doi:
426 10.18653/v1/N16-1136. URL <https://aclanthology.org/N16-1136>.
- 427 Nate Kushman, Yoav Artzi, Luke Zettlemoyer, and Regina Barzilay. Learning to automatically
428 solve algebra word problems. In *Proceedings of the 52nd Annual Meeting of the Association*
429 *for Computational Linguistics (Volume 1: Long Papers)*, pages 271–281, Baltimore, Maryland,
430 June 2014. Association for Computational Linguistics. doi: 10.3115/v1/P14-1026. URL [https:](https://aclanthology.org/P14-1026)
431 [//aclanthology.org/P14-1026](https://aclanthology.org/P14-1026).

- 432 Xiao Li, Yawei Sun, and Gong Cheng. Tsqa: tabular scenario based question answering. In
433 *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 13297–13305,
434 2021.
- 435 Xiao Li, Yin Zhu, Sichen Liu, Jiangzhou Ju, Yuzhong Qu, and Gong Cheng. Dyrren: A dynamic
436 retriever-reranker-generator model for numerical reasoning over tabular and textual data. In
437 *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pages 13139–13147,
438 2023a.
- 439 Yifei Li, Zeqi Lin, Shizhuo Zhang, Qiang Fu, Bei Chen, Jian-Guang Lou, and Weizhu Chen.
440 Making language models better reasoners with step-aware verifier. In *Proceedings of the 61st*
441 *Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*,
442 pages 5315–5333, Toronto, Canada, July 2023b. Association for Computational Linguistics. doi:
443 10.18653/v1/2023.acl-long.291. URL <https://aclanthology.org/2023.acl-long.291>.
- 444 Yuan-Fang Li, Sébastien Bubeck, Ronen Eldan, Allison Del Giorno, Suriya Gunasekar, and Yin Tat
445 Lee. Textbooks are all you need ii: phi-1.5 technical report. *ArXiv*, abs/2309.05463, 2023c. URL
446 <https://api.semanticscholar.org/CorpusID:261696657>.
- 447 Jia-Yin Liu, Zhenya Huang, Zhiyuan Ma, Qi Liu, Enhong Chen, Tianhuang Su, and Haifeng Liu.
448 Guiding mathematical reasoning via mastering commonsense formula knowledge. *Proceedings*
449 *of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, 2023. URL
450 <https://api.semanticscholar.org/CorpusID:260499357>.
- 451 Meta. Meta llama 3, 2024. URL <https://llama.meta.com/llama3/>.
- 452 OpenAI. Gpt-4 technical report. *ArXiv*, 2023. URL [https://api.semanticscholar.org/](https://api.semanticscholar.org/CorpusID:257532815)
453 [CorpusID:257532815](https://api.semanticscholar.org/CorpusID:257532815).
- 454 Shuming Shi, Yuehui Wang, Chin-Yew Lin, Xiaojiang Liu, and Yong Rui. Automatically solving
455 number word problems by semantic parsing and reasoning. In *Proceedings of the 2015 Conference*
456 *on Empirical Methods in Natural Language Processing*, pages 1132–1142, Lisbon, Portugal,
457 September 2015. Association for Computational Linguistics. doi: 10.18653/v1/D15-1135. URL
458 <https://aclanthology.org/D15-1135>.
- 459 Kashun Shum, Shizhe Diao, and Tong Zhang. Automatic prompt augmentation and selection with
460 chain-of-thought from labeled data. In Houda Bouamor, Juan Pino, and Kalika Bali, editors,
461 *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 12113–12139,
462 Singapore, December 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.
463 findings-emnlp.811. URL <https://aclanthology.org/2023.findings-emnlp.811>.
- 464 InternLM Team. Internlm: A multilingual language model with progressively enhanced capabilities.
465 <https://github.com/InternLM/InternLM>, 2023.
- 466 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay
467 Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation
468 and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- 469 Lei Wang, Dongxiang Zhang, Jipeng Zhang, Xing Xu, Lianli Gao, Bing Tian Dai, and Heng Tao
470 Shen. Template-based math word problem solvers with recursive neural networks. In *Proceedings*
471 *of the AAAI Conference on Artificial Intelligence*, volume 33, pages 7144–7151, 2019.
- 472 Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V Le, Ed H Chi, Sharan Narang, Aakanksha
473 Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language
474 models. In *The Eleventh International Conference on Learning Representations*, 2022.
- 475 Yan Wang, Xiaojiang Liu, and Shuming Shi. Deep neural solver for math word problems. In
476 *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages
477 845–854, Copenhagen, Denmark, September 2017. Association for Computational Linguistics.
478 doi: 10.18653/v1/D17-1088. URL <https://aclanthology.org/D17-1088>.

- 479 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny
480 Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in*
481 *Neural Information Processing Systems*, 35:24824–24837, 2022.
- 482 Chenxi Whitehouse, Monojit Choudhury, and Alham Fikri Aji. LLM-powered data augmentation for
483 enhanced cross-lingual performance. In *The 2023 Conference on Empirical Methods in Natural*
484 *Language Processing*, 2023. URL <https://openreview.net/forum?id=wFWwyXE1N>.
- 485 Aohan Zeng, Xiao Liu, Zhengxiao Du, Zihan Wang, Hanyu Lai, Ming Ding, Zhuoyi Yang, Yifan Xu,
486 Wendi Zheng, Xiao Xia, et al. Glm-130b: An open bilingual pre-trained model. In *The Eleventh*
487 *International Conference on Learning Representations*, 2022.
- 488 Chujie Zheng, Sahand Sabour, Jiaxin Wen, Zheng Zhang, and Minlie Huang. AugESC: Dialogue
489 augmentation with large language models for emotional support conversation. In *Findings of the*
490 *Association for Computational Linguistics: ACL 2023*, pages 1552–1568, Toronto, Canada, July
491 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.findings-acl.99. URL
492 <https://aclanthology.org/2023.findings-acl.99>.
- 493 Denny Zhou, Nathanael Schärli, Le Hou, Jason Wei, Nathan Scales, Xuezhi Wang, Dale Schuurmans,
494 Claire Cui, Olivier Bousquet, Quoc V Le, et al. Least-to-most prompting enables complex
495 reasoning in large language models. In *The Eleventh International Conference on Learning*
496 *Representations*, 2022.

497 **Checklist**

- 498 1. For all authors...
- 499 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
- 500 contributions and scope? [Yes]
- 501 (b) Did you describe the limitations of your work? [Yes] Section 6.
- 502 (c) Did you discuss any potential negative societal impacts of your work? [Yes] Section 6.
- 503 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
- 504 them? [Yes]
- 505 2. If you are including theoretical results...
- 506 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 507 (b) Did you include complete proofs of all theoretical results? [N/A]
- 508 3. If you ran experiments (e.g. for benchmarks)...
- 509 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
- 510 mental results (either in the supplemental material or as a URL)? [Yes] Section 1.
- 511 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
- 512 were chosen)? [Yes] Appendix C.4.
- 513 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
- 514 ments multiple times)? [Yes] Appendix C.4.
- 515 (d) Did you include the total amount of compute and the type of resources used (e.g., type
- 516 of GPUs, internal cluster, or cloud provider)? [Yes] Appendix C.4.
- 517 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 518 (a) If your work uses existing assets, did you cite the creators? [Yes] Section 4.
- 519 (b) Did you mention the license of the assets? [Yes] Section 1.
- 520 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
- 521 Section A.
- 522 (d) Did you discuss whether and how consent was obtained from people whose data you’re
- 523 using/curating? [Yes] Section A.
- 524 (e) Did you discuss whether the data you are using/curating contains personally identifiable
- 525 information or offensive content? [Yes] Section A.
- 526 5. If you used crowdsourcing or conducted research with human subjects...
- 527 (a) Did you include the full text of instructions given to participants and screenshots, if
- 528 applicable? [Yes] Section 3.
- 529 (b) Did you describe any potential participant risks, with links to Institutional Review
- 530 Board (IRB) approvals, if applicable? [Yes] Section A.
- 531 (c) Did you include the estimated hourly wage paid to participants and the total amount
- 532 spent on participant compensation? [Yes] Section 4.

533 A Dataset Card

534 A.1 Motivation

535 **1. For what purpose was the dataset created? Was there a specific task in mind? Was there a**
536 **specific gap that needed to be filled? Please provide a description.**

537 The motivation behind constructing FormulaReasoning comes from the need to address the limitations
538 of existing numerical reasoning datasets. While numerical reasoning has seen significant advance-
539 ments with the rise of LLMs and specialized datasets, current datasets often lack knowledge-guided
540 reasoning process. They typically rely on implicit commonsense knowledge rather than explicit
541 formulas, which becomes problematic when LLMs encounter hallucinations.

542 To overcome these limitations, FormulaReasoning was created to emphasize the use of specific
543 formulas in numerical reasoning. Unlike previous datasets that primarily rely on implicit knowledge,
544 FormulaReasoning requires explicit formula-based reasoning. This shift introduces a higher level of
545 challenge and reflects real-world numerical problem-solving scenarios better.

546 **2. Who created the dataset (e.g., which team, research group) and on behalf of which entity**
547 **(e.g., company, institution, organization)?**

548 FormulaReasoning is created by Xiao Li, Bolin Zhu, Sichen Liu, Yin Zhu, Yiwei Liu and Gong
549 Cheng from the State Key Laboratory for Novel Software Technology, Nanjing University.

550 **3. Who funded the creation of the dataset? If there is an associated grant, please provide the**
551 **name of the grantor and the grant name and number.**

552 This work was supported by the CIPSC-SMP-Zhipu.AI Large Model Cross-Disciplinary Fund.

553 A.2 Composition

554 **1. What do the instances that comprise the dataset represent (e.g., documents, photos, people,**
555 **countries)? Are there multiple types of instances (e.g., movies, users, and ratings; people and**
556 **interactions between them; nodes and edges)? Please provide a description.**

557 The data within the dataset exclusively comprises elementary physics questions based on daily
558 life scenarios, all organized in text format, without photos, specific people information or specific
559 countries.

560 **2. How many instances are there in total (of each type, if appropriate)?**

561 We divided FormulaReasoning into training, *id* (in-distribution) test, and *ood* (out-of-distribution)
562 test, comprising 4,608, 421 and 391 questions, respectively.

563 **3. Does the dataset contain all possible instances or is it a sample (not necessarily random)**
564 **of instances from a larger set? If the dataset is a sample, then what is the larger set? Is the**
565 **sample representative of the larger set (e.g., geographic coverage)? If so, please describe how**
566 **this representativeness was validated/verified. If it is not representative of the larger set, please**
567 **describe why not (e.g., to cover a more diverse range of instances, because instances were**
568 **withheld or unavailable).**

569 FormulaReasoning is not from a larger set.

570 **4. What data does each instance consist of? “Raw” data (e.g., unprocessed text or images) or**
571 **features? In either case, please provide a description.**

572 Each instance consists of a question, the formulas, the parameters within these formulas and
573 their corresponding numerical values, textual explanations, and the final numerical answer. See
574 <https://github.com/nju-websoft/FormulaReasoning> for more details.

575 **5. Is there a label or target associated with each instance? If so, please provide a description.**

576 Yes, each instance contains textual explanations, and the final numerical answer.

577 **6. Is any information missing from individual instances? If so, please provide a description,**
578 **explaining why this information is missing (e.g., because it was unavailable). This does not**
579 **include intentionally removed information, but might include, e.g., redacted text.**

580 No.

581 **7. Are relationships between individual instances made explicit (e.g., users' movie ratings, social**
582 **network links)? If so, please describe how these relationships are made explicit.**

583 N/A.

584 **8. Are there recommended data splits (e.g., training, development/validation, testing)? If so,**
585 **please provide a description of these splits, explaining the rationale behind them.**

586 Yes. We divided FormulaReasoning into training, *id* (in-distribution) test, and *ood* (out-of-distribution)
587 test, comprising 4,608, 421 and 391 questions, respectively. We required that all formulas in the *id*
588 test must appear in the training set, whereas in the *ood* test, each question involves at least one formula
589 that has not been seen in the training set. This division is designed to evaluate the generalization
590 capabilities of fine-tuned models on formulas that they have not previously encountered.

591 **9. Are there any errors, sources of noise, or redundancies in the dataset? If so, please provide a**
592 **description.**

593 Currently, there are no known errors, noise, or redundancies. We have addressed these occurrences
594 during the annotation process.

595 **10. Is the dataset self-contained, or does it link to or otherwise rely on external resources (e.g.,**
596 **websites, tweets, other datasets)? If it links to or relies on external resources, a) are there**
597 **guarantees that they will exist, and remain constant, over time; b) are there official archival**
598 **versions of the complete dataset (i.e., including the external resources as they existed at the time**
599 **the dataset was created); c) are there any restrictions (e.g., licenses, fees) associated with any of**
600 **the external resources that might apply to a dataset consumer? Please provide descriptions of**
601 **all external resources and any restrictions associated with them, as well as links or other access**
602 **points, as appropriate.**

603 Yes, FormulaReasoning is self-contained, and it doesn't rely on any external resources.

604 **11. Does the dataset contain data that might be considered confidential (e.g., data that is**
605 **protected by legal privilege or by doctor-patient confidentiality, data that includes the content**
606 **of individuals' non-public communications)? If so, please provide a description.**

607 No.

608 **12. Does the dataset contain data that, if viewed directly, might be offensive, insulting, threaten-**
609 **ing, or might otherwise cause anxiety? If so, please describe why.**

610 No. Firstly, it is unlikely for harmful information to appear in the questions designed for middle
611 school education. Secondly, we have not identified such information within the dataset.

612 **13. Does the dataset relate to people? If not, you may skip the remaining questions in this**
613 **section.**

614 No.

615 **A.3 Collection Process**

616 **1. How was the data associated with each instance acquired?**

617 See Section 3.

618 **2. What mechanisms or procedures were used to collect the data (e.g., hardware apparatuses or**
619 **sensors, manual human curation, software programs, software APIs)?**

620 See Section 3.

621 **3. If the dataset is a sample from a larger set, what was the sampling strategy (e.g., deterministic,**
622 **probabilistic with specific sampling probabilities)?**

623 Our FormulaReasoning is not sampled from a larger set.

624 **4. Who was involved in the data collection process (e.g., students, crowdworkers, contractors)**
625 **and how were they compensated (e.g., how much were crowdworkers paid)?**

626 A total of 5 graduate students participated in the annotation work, and 108 high school students were
627 involved in the human performance tasks. For more details, see Section 3 and Section 4.

628 **5. Over what timeframe was the data collected?**

629 The questions in FormulaReasoning were derived from junior high school physics examinations in
630 China over the past 14 years (2010 – 2024).

631 **6. Were any ethical review processes conducted (e.g., by an institutional review board)?**

632 The ethical review board of our department has approved our experiment.

633 **A.4 Preprocessing/cleaning/labeling**

634 **1. Was any preprocessing/cleaning/labeling of the data done (e.g., discretization or bucketing,**
635 **tokenization, part-of-speech tagging, SIFT feature extraction, removal of instances, processing**
636 **of missing values)?**

637 Yes. For more details, see Section 3.

638 **2. Was the “raw” data saved in addition to the preprocessed/cleaned/labeled data (e.g., to**
639 **support unanticipated future uses)?**

640 Yes, the raw data has been included in the released dataset.

641 **3. Is the software that was used to preprocess/clean/label the data available?**

642 Yes, they are included in our GitHub repository.

643 **A.5 Uses**

644 **1. Has the dataset been used for any tasks already? If so, please provide a description.**

645 Yes, in this paper, we utilized the dataset to evaluate the reasoning ability of language models.

646 **2. Is there a repository that links to any or all papers or systems that use the dataset? If so,**
647 **please provide a link or other access point.**

648 N/A. Currently, there have been no external works that have utilized FormulaReasoning.

649 **3. What (other) tasks could the dataset be used for?**

650 FormulaReasoning can be utilized for evaluating the reasoning ability of language models, particularly
651 in scenarios requiring knowledge (formulas). Additionally, the formula database we constructed can
652 be employed for evaluating retrieval-augmented generation models. Furthermore, we partitioned the
653 test set into id and ood tests for assessing the generalization ability of language models.

654 **4. Is there anything about the composition of the dataset or the way it was collected and**
655 **preprocessed/cleaned/labeled that might impact future uses? For example, is there anything**
656 **that a dataset consumer might need to know to avoid uses that could result in unfair treatment**
657 **of individuals or groups (e.g., stereotyping, quality of service issues) or other risks or harms**
658 **(e.g., legal risks, financial harms)? If so, please provide a description. Is there anything a**
659 **dataset consumer could do to mitigate these risks or harms?**

660 No. Our data originates from elementary physics questions based on everyday life scenarios, exclud-
661 ing any potentially harmful information.

662 **5. Are there tasks for which the dataset should not be used? If so, please provide a description.**

663 No.

664 **A.6 Distribution**

665 **1. Will the dataset be distributed to third parties outside of the entity (e.g., company, institution,**
666 **organization) on behalf of which the dataset was created? If so, please provide a description.**

667 No. We only open source the datasets through public channels: [https://github.com/nju-](https://github.com/nju-websoft/FormulaReasoning)
668 [websoft/FormulaReasoning](https://github.com/nju-websoft/FormulaReasoning).

669 **2. How will the dataset will be distributed (e.g., tarball on website, API, GitHub)? Does the**
670 **dataset have a digital object identifier (DOI)?**

671 Our code is available at <https://github.com/nju-websoft/FormulaReasoning> under the
672 Apache 2.0 License.

673 Our data is available at <https://zenodo.org/doi/10.5281/zenodo.11408109> under the Cre-
674 ative Commons Attribution 4.0 International (CC BY 4.0) license.

675 DOI: 10.5281/zenodo.11408109.

676 Croissant metadata: [https://huggingface.co/api/datasets/xli/FormulaReasoning/](https://huggingface.co/api/datasets/xli/FormulaReasoning/croissant)
677 [croissant](https://huggingface.co/api/datasets/xli/FormulaReasoning/croissant).

678 **3. When will the dataset be distributed?**

679 We have distributed FormulaReasoning.

680 **4. Will the dataset be distributed under a copyright or other intellectual property (IP) license,**
681 **and/or under applicable terms of use (ToU)? If so, please describe this license and/or ToU, and**
682 **provide a link or other access point to, or otherwise reproduce, any relevant licensing terms or**
683 **ToU, as well as any fees associated with these restrictions.**

684 Our code is distributed under the Apache License, Version 2.0. Our data is distributed under the
685 Creative Commons Attribution 4.0 International (CC BY 4.0) license.

686 **5. Have any third parties imposed IP-based or other restrictions on the data associated with the**
687 **instances? If so, please describe these restrictions, and provide a link or other access point to,**
688 **or otherwise reproduce, any relevant licensing terms, as well as any fees associated with these**
689 **restrictions.**

690 No.

691 **6. Do any export controls or other regulatory restrictions apply to the dataset or to individual**
692 **instances? If so, please describe these restrictions, and provide a link or other access point to,**
693 **or otherwise reproduce, any supporting documentation.**

694 No.

695 **A.7 Maintenance**

696 **1. Who will be supporting/hosting/maintaining the dataset?**

697 The Authors.

698 **2. How can the owner/curator/manager of the dataset be contacted (e.g., email address)?**

699 Contact authors via emails listed under the title or through GitHub issues.

700 **3. Is there an erratum? If so, please provide a link or other access point.**

701 No.

702 **4. Will the dataset be updated (e.g., to correct labeling errors, add new instances, delete**
703 **instances)? If so, please describe how often, by whom, and how updates will be communicated**
704 **to dataset consumers (e.g., mailing list, GitHub)?**

705 Updates, if any, will be provided on GitHub by the authors.

706 **5. If the dataset relates to people, are there applicable limits on the retention of the data**
707 **associated with the instances (e.g., were the individuals in question told that their data would**
708 **be retained for a fixed period of time and then deleted)? If so, please describe these limits and**
709 **explain how they will be enforced.**

710 No, FormulaReasoning doesn't relate to people.

711 **6. Will older versions of the dataset continue to be supported/hosted/maintained? If so, please**
712 **describe how. If not, please describe how its obsolescence will be communicated to dataset**
713 **consumers.**

714 N/A.

715 **7. If others want to extend/augment/build on/contribute to the dataset, is there a mechanism for**
716 **them to do so? If so, please provide a description. Will these contributions be validated/verified?**
717 **If so, please describe how. If not, why not? Is there a process for communicating/distributing**
718 **these contributions to dataset consumers? If so, please provide a description.**

719 Others can do anything subject to the license of our dataset.

720 **B The Machine Learning Reproducibility Checklist**

721 1. For all models and algorithms presented, check if you include:

722 (a) A clear description of the mathematical setting, algorithm, and/or model. [\[Yes\]](#) See
723 Section 4.

724 (b) A clear explanation of any assumptions. [\[N/A\]](#)

725 (c) An analysis of the complexity (time, space, sample size) of any algorithm. [\[Yes\]](#) See
726 Appendix C.4.

727 2. For any theoretical claim, check if you include:

728 (a) A clear statement of the claim. [\[N/A\]](#)

729 (b) A complete proof of the claim. [\[N/A\]](#)

730 3. For all datasets used, check if you include:

731 (a) The relevant statistics, such as number of examples. [\[Yes\]](#) See Section 4.

732 (b) The details of train / validation / test splits. [\[Yes\]](#) See Section 4.

733 (c) An explanation of any data that were excluded, and all pre-processing step. [\[Yes\]](#) See
734 Section 3 and Section 4.

735 (d) A link to a downloadable version of the dataset or simulation environment. [\[Yes\]](#) See
736 Appendix A.

737 (e) For new data collected, a complete description of the data collection process, such as
738 instructions to annotators and methods for quality control. [\[Yes\]](#) See Section 3.

739 4. For all shared code related to this work, check if you include:

740 (a) Specification of dependencies. [\[Yes\]](#)

741 (b) Training code. [\[Yes\]](#)

742 (c) Evaluation code. [\[Yes\]](#)

743 (d) (Pre-)trained model(s). [\[Yes\]](#)

744 (e) README file includes table of results accompanied by precise command to run to
745 produce those results. [\[Yes\]](#)

746 5. For all reported experimental results, check if you include:

- 747 (a) The range of hyper-parameters considered, method to select the best hyper-parameter
748 configuration, and specification of all hyper-parameters used to generate results. [Yes]
749 See Appendix C.4.
- 750 (b) The exact number of training and evaluation runs. [Yes] See Appendix C.4.
- 751 (c) A clear definition of the specific measure or statistics used to report results. [Yes] See
752 Section 4.
- 753 (d) A description of results with central tendency (e.g. mean) & variation (e.g. error bars).
754 [N/A]
- 755 (e) The average runtime for each result, or estimated energy cost. [Yes] See Appendix C.4.
- 756 (f) A description of the computing infrastructure used. [Yes] See Appendix C.4.

757 C Appendix

758 C.1 Dataset Construction

759 C.1.1 Prompts in Formula Normalization

760 The process of formula normalization is delineated into three distinct stages: the generation of natural
761 language explanations, the extraction of the associated parameters from the explanations, and the
762 subsequent error correction phase. The initial two stages are illustrated in Figures 3 and 4. The third
763 stage is further splited into three specific error categories, each addressed by a dedicated prompt: input
764 errors, where the parameters mentioned in the explanation are absent from the question; calculation
765 errors, which occur when the calculator reports an error during the computation process; and output
766 errors, where the final computed answer is incorrect. We provide an example here focusing on
767 prompts for correcting calculation errors, while prompts for the other two error types can be found in
768 our code submission. The prompts designed to correct calculation errors are depicted in Figure 5.
769 The entire normalization procedure employs a 6-shot prompting, an instance of which is provided
770 herein for illustrative purposes.

771 C.1.2 Examples of Deleted Questions

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

775 C.1.3 Semantic-based Merging for Formula Database Construction

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

781 C.2 Case Study and Error Analysis

782 We sampled 50 error cases from the id test (few-shot setting) of GPT-3.5-turbo and manually
783 categorized the types and proportions of errors. We divided the error types into two main categories:
784 *formula errors* and *calculation errors*. Formula errors encompass inappropriate formulas and omitted
785 formulas, while calculation errors primarily involve inaccuracies in numerical calculation and unit
786 errors. We found that 38% of errors were caused by incorrect formulas, while the remaining 62%
787 were attributable to calculation errors. We provide one example for each of the two types of errors
788 listed in Figure 2. It could be observed that FormulaReasoning poses challenges to existing models in
789 terms of formula application and numerical calculation (including unit calculation and arithmetic
790 calculation).

791 C.3 Experiments

792 C.3.1 Data Augmentation (DA) for FormulaReasoning

793 There have been several studies utilizing large language models (LLMs) for data augmentation (Ding
794 et al., 2024). The data generated in these related works (Zheng et al., 2023; Whitehouse et al.,
795 2023) primarily focus on daily conversations or sentiment analysis and do not require rigorous
796 numerical calculations. Some research on data augmentation involving numerical calculations (Shum
797 et al., 2023) employs LLMs to generate solutions to questions to aid in training, rather than creating
798 complete questions. In contrast to these approaches, our work generates complete questions that
799 involve numerical calculations (particularly formula calculations), along with automatic improvement
800 and selection to ensure data quality.

Question

The maximum power of a tank engine is $7.2 \times 10^5 \text{ W}$. 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. ($q_{\text{diesel}} = 4.0 \times 10^7 \text{ J/kg}$) Q: The efficiency of the tank's engine.

Output of GPT-3.5-turbo (wrong formula highlighted in red)

Solution: 1. Calculate the time it takes for the tank to travel: ... omitted ...
3. Calculate the engine efficiency: $[\text{Engine efficiency}] = [\text{Engine output power}] / [\text{Heat energy consumed by fuel}]$... omitted ...

Gold Explanation (correct formula highlighted in green)

Solution: 1. Calculate the heat released by the complete combustion of diesel fuel: ... omitted ...
3. Calculate the efficiency of the tank's engine: $[\text{Tank engine efficiency}] = [\text{Work done by the engine}] / [\text{Heat released by complete combustion of diesel fuel}] \times 100\%$ Equation ... omitted ...

(a) An error case caused by wrong formulas.

Question

Given that the calorific value of natural gas is $3.8 \times 10^7 \text{ J/m}^3$ and the specific heat capacity of water is $4.2 \times 10^3 \text{ J/(kg} \cdot \text{°C)}$. 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 ?

Output of GPT-3.5-turbo (wrong calculation highlighted in red. Note that $1 \text{ L} = 10^{-3} \text{ m}^3$)

Solution: 1. Calculate the heat content of natural gas: $[\text{Heat content of natural gas}] = [\text{Calorific value of natural gas}] \times [\text{Volume of natural gas}]$ Calculation = $3.8 \times 10^7 \text{ J/m}^3 \times 10 \text{ L} = 3.8 \times 10^8 \text{ J}$ Heat content of natural gas = $3.8 \times 10^8 \text{ J}$... omitted ...

Gold Explanation (correct calculation highlighted in green)

Solution: 1. Calculate the heat released from the complete burning of natural gas:
 $[\text{Heat released from complete combustion of natural gas}] = [\text{Volume of natural gas}] \times [\text{Calorific value of natural gas}]$
Calculation = $10 \text{ L} \times 3.8 \times 10^7 \text{ J/m}^3 = 3.8 \times 10^5 \text{ J}$... omitted ...

(b) An error case caused by wrong calculation.

Figure 2: Error cases.

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

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

810 We found that mixing the newly generated data into the original training set did not always bring
811 positive improvement, perhaps because the newly generated data has not undergone manual re-
812 view. We found that randomly selecting a small portion of the newly generated data can enable
813 the model to have performance improvement. We set several different mixing ratios selected from
814 $\{5\%, 10\%, 15\%, 20\%, 2\%, 30\%, 35\%, 40\%\}$. We fine-tuned the ChatGLM-6B-base using the aug-
815 mented data set. After training for a fixed number of steps (150k and 200k), we selected the
816 checkpoints with the smallest loss among models of different mixing ratios.

817 C.4 Implementation Details

818 We accessed to GPT-4-turbo, GPT-3.5-turbo⁶, GLM4⁷, and Qwen-max⁸ through API calls with the
819 default hyper-parameters. For other LLMs, we conducted experiments on NVIDIA V100-32G GPUs
820 for 7B models, and on NVIDIA A100-80G GPUs for 14B/20B models. These LLMs generated using
821 nucleus sampling with $\text{top}_p=0.8$. Models that require fine-tuning were experimented on NVIDIA
822 V100 GPUs with Huggingface Transformers and Pytorch 2.0. For mT5-base and mT5-large, we set
823 a learning rate of $5e-5$ and a batch size of 32, testing the model after training for 50 epochs. For
824 Qwen-1.8B, we used a learning rate of $1e-5$ and a batch size of 32, and tested the model after training

⁶<https://platform.openai.com/docs>

⁷<https://open.bigmodel.cn/>

⁸<https://help.aliyun.com/zh/dashscope/developer-reference/quick-start>

825 for 10 epochs. For ChatGLM3-6B, we fine-tuned with LoRA Hu et al. (2021) with r=8, alpha=32
 826 and learning rate of 5e-5, batch size of 1. The max input length and output length are both set to
 827 512. We utilized nucleus sampling with top_p=0.8 for generation. In the case of CoT-SFT, which
 828 directly outputted formulas along with corresponding parameter values and units, if the generation
 829 output contained formatting errors, we allowed the small model to retry up to 5 times until a correctly
 830 formatted output was generated. Training mT5-base, mT5-large, Qwen-1.8B, ChatGLM-6B models
 831 requires 6, 12, 12 and 24 hours respectively.

832 C.4.1 Zero-shot and Few-shot Prompts

833 Zero-shot and few-shot prompts are shown in Figure 8.

834 C.4.2 Formula Retriever

835 Let the number of formulas in the formula database be N . During training, we randomly initialized
 836 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
 837 representation of the i -th formula in formula database. We denoted a batch of questions with a batch
 838 size of B as $Q = \{q_1, q_2, \dots, q_B\}$. The indices of the gold-standard formulas corresponding to these
 839 B questions were denoted as $L = \{l_1, l_2, \dots, l_B\}$ (i.e. the label of q_i is l_i , where $1 \leq i \leq B$).

840 BERT was utilized to encode each question,

$$\mathbf{h}_{cls}^i, \mathbf{h}_1^i, \dots = \text{BERT}(q_i), 1 \leq i \leq B. \quad (1)$$

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

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

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

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

847 Models were evaluated with top-5 retrieved formulas. Prompts can be found in Appendix C.4.4. We
 848 utilized zero-shot CoTs.

849 C.4.3 Few-shot Experiments on the LLMs of Chat Versions

850 In this experiment, we compared the performance of the same version of the model under zero-shot
 851 and few-shot settings. Results are shown in Table 7. For the chat version of the LLMs, we could
 852 observe that few-shot can effectively improve model performance, with performance improvements
 853 ranging from 1.27 to 9.18 on average across id test and ood test. Comparing the performance of the
 854 base version and chat version of the same model under few-shot settings, except for minimal changes
 855 on InternLM-chat-7B and Llama2-chat-7B, the performance of the other models showed a decrease
 856 from base to chat versions.

857 C.4.4 Prompts for LLMs with Formula Retriever

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

Model	Size	id test	ood test	Avg.
zero-shot CoTs with LLMs of chat/instruct versions				
InternLM-chat	20B	5.70	4.60	5.17
Qwen-chat	14B	32.07	37.60	34.73
Llama3-instruct	8B	22.66	17.98	20.41
Llama2-chat	7B	0.00	0.26	0.13
Qwen-chat	7B	7.36	8.70	8.01
InternLM-chat	7B	7.84	7.67	7.76
few-shot CoTs with LLMs of base versions				
InternLM-base	20B	18.29	11.25	14.90
Qwen-base	14B	44.89	36.83	41.01
Llama3-base	8B	12.81	8.87	10.91
Llama2-base	7B	1.43	0.26	0.87
Qwen-base	7B	21.14	18.16	19.71
InternLM-base	7B	9.50	8.18	8.86
few-shot CoTs with LLMs of chat/instruct versions				
InternLM-chat	20B	11.58	10.10	10.87
Qwen-chat	14B	41.38	37.93	39.72
Llama3-instruct	8B	6.90	6.16	6.54
Llama2-chat	7B	1.97	1.00	1.50
Qwen-chat	7B	19.21	15.02	17.19
InternLM-chat	7B	10.10	7.88	9.03

Table 7: Results of different versions of the LLMs with zero-shot and few-shot on FormulaReasoning.

Prompt actually used	English translation
<p>我需要你修改问题原有的解析，给出规范格式的新解析，要求如下：</p> <ol style="list-style-type: none"> 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成基本公式进行求解 公式中的计算符号，如“+”、“-”、“×”、“/”和“^”不能省略 公式需要同时给出符号和有具体含义的两种形式，然后代入数值计算得出答案 涉及到单位换算的部分需要展示出来具体过程 使用“[]”标注公式中的变量，其中科学计数法形式的数字“$a \times 10^b$”以及复杂的单位，需要使用“()”标注 如果有latex格式的公式，比如“$\frac{Q_{吸}}{Q_{放}}$”需要改成正常算式的形式：“$Q_{吸}/Q_{放}$” 解析开始如果给出“(解析可能存在错误)”表明给出的解析中可能存在错误 <p>下面是一些例子：</p> <p>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倍，在0.5h内行驶了28km的路程，求：行驶的平均速度为多少？</p> <p>解析：解：车辆在0.5h内行驶了28km的路程，所以行驶的平均速度为：$V=s/t=28km/0.5h=56km/h$。</p> <p>output: 解： 原始解析比较清晰，下面将其转换为规范格式： 第一步，先计算平均速度： $v=s/t$ [平均速度]=[路程]/[时间] 算式=$28km/0.5h=56km/h$ 平均速度=$56km/h$ 答案=$56km/h$...omitted...</p> <p>下面是一个新的问题： 问题：{{question}} 解析：{{explanation}}</p>	<p>I need you to modify the original explanation of the question and provide a new explanation with the following requirements:</p> <ol style="list-style-type: none"> Please think step by step. If there has formula combination, you need to decompose the combination into basic formulas step by step. Calculation symbols such as "+", "-", "x", "/" and "^" in formulas cannot be omitted. The formula needs to be given in both symbolic and concrete forms. After that, you need to substitute into the numerical calculation to obtain the answer. The part related to unit conversion needs to show the specific process. Use "[]" to label variables in formula, "()" for numbers like "$a \times 10^b$" in scientific notation and for complex units. If you have a latex formula such as "$\frac{Q_{in}}{Q_{out}}$", change it to a normal formula: "Q_{in}/Q_{out}". Start with "(Explanation may be wrong)" indicates that there may be an error in the given explanation. <p>Here are some examples:</p> <p>input: Question: In a test, the resistance of the car at a constant speed is 0.02 times the weight of the car, and the car traveled a distance of 28km in 0.5h. What is the average speed? Explanation: The car has traveled a distance of 28km in 0.5h. The average speed of the vehicle is $V=s/t=28km/0.5h=56km/h$.</p> <p>output: Solution: The raw parsing is clear, so let's convert it to a canonical format: $v=s/t$ [Average speed]=[distance]/[time] expression=$28km/0.5h=56km/h$ Average speed =$56km/h$ Answer=$56km/h$...omitted...</p> <p>There's a new question: Question: {{question}} Explanation: {{explanation}}</p>

Figure 3: Prompt of the formula normalization stage 1.

Prompt actually used

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公式，以及公式对应的参数表：

- 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.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
- 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,而且变量需要使用有意义的文字标识,尽量避免直接使用数值
- 4.如果有latex格式的公式,比如 $\frac{Q_{吸}}{Q_{放}}$ 需要改成正常算式的形式: "[Q_吸]/[Q_放]", 算式中的单位换算部分不属于公式,不需要被提取
- 5.参数表中的参数是公式中使用到的参数(参数名称要与公式中的参数一致),表格包括:概念、符号、数值、单位,使用"|"分割单元格
- 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式(没有进行单位换算)
- 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,比如"3盏灯","转了8圈","4个人"等等,那么单位写入"null"
- 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息

下面是一些例子:

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倍,在0.5h内行驶了28km的路程,求:行驶的平均速度为多少?

解析: 解:

原始解析比较清晰,下面将其转换为规范格式:

第一步,先计算平均速度:

$v=s/t$

[平均速度]=[路程]/[时间]

算式=28km/0.5h=56km/h

平均速度=56km/h

答案=56km/h

output:

提取分析:

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"[]"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

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)

提取结果:

1.先计算平均速度,根据路程和时间的关系:

[平均速度]=[路程]/[时间]

参数表:

概念|符号|数值|单位

路程|s|28|km

时间|t|0.5|h

平均速度|v|56|km/h

...omitted...

下面是一个新的问题:

问题: {{question}}

解析: {{explanation}}

English translation

I need you to help me extract the formula and parameters table for the formula from given question and explanation:

1. Think step by step, analyze the explanation first, and then generate the extracted results. If there is a combination of formulas, the combination needs to be split into basic formulas step by step.
2. All calculation symbols such as "+", "-", "×", "/" and "^" in the formula cannot be omitted.
3. Each variable in the formula needs to be labeled with "[]", and the variable needs to be identified with meaningful text instead of numbers.
4. If a latex formula such as $\frac{Q_{in}}{Q_{out}}$ needs to be changed to a normal formula: "[Q_in]/[Q_out]". The unit conversion does not need to be extracted.
5. The parameters table come from the parameters in formula (the parameter name should be consistent with the parameters in the formula), the table include: concept, symbol, numeric, unit, using cell division "|".
6. The numeric and unit in the parameter table come from the problem itself and the intermediate results of analytical calculation. If the parameters are converted into different units, the parameter table should give the original parameter form (without unit conversion).
7. If the parameter in the parameters table has no units or represents the amount of an object, such as "3 lights", "8 revolutions", "4 people", etc., then the units are written as "null".
8. Apart from formula and parameter table, no additional information is required in the extraction results.

Here are some examples:

input:

Question: In a test, the resistance of the car at a constant speed is 0.02 times the weight of the car, and the car traveled a distance of 28km in 0.5h. What is the average speed?

Explanation:

1.we compute the average speed:

$v=s/t$

[Average speed]=[distance]/[time]

expression=28km/0.5h=56km/h

Average speed =56km/h

Answer=56km/h

output:

Extraction analysis:

1. The formula format in the analysis meets the requirements. The calculation symbols are not omitted. Variables are labeled with "[]", and variables are expressed with meaningful text.

2. No unit conversion was performed during the computation, and parameters can be directly extracted.

Below is the extraction result (the unit conversion part does not belong to the formula and does not need to be extracted; no additional explanatory information is required in the parameter table).

Extraction result:

1. First calculate the average speed, based on the relationship between distance and time:

[average speed]=[distance]/[time]

Parameter table:

Concept | Symbol | Numeric | Unit

distance | s | 28 | km

time | t | 0.5 | h

average speed | v | 56 | km/h

...omitted...

There's a new question:

Question: {{question}}

Explanation: {{explanation}}

Figure 4: Prompt of the formula normalization stage 2.

Prompt actually used	English translation
<p>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，下面是错误纠正的要求：</p> <ol style="list-style-type: none"> 1.你需要先进行错误分析，分析如何修改来纠正错误，然后给出错误纠正部分，纠正解析中的错误 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 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 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，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 5.公式的格式为“[待求解参数]=[参数1](+ - × ÷)[参数2]...”；参数的格式为：“概念 符号 数值 单位”，比如“水的沸点是100°C”，表示为“水的沸点 t_沸 100 °C” <p>下面是一些例子：</p> <p>input: 问题：假设13.0吨烟煤在煤炉中完全燃烧，放出的热量部分被水吸收，可以使$4 \times 10^5 \text{kg}$的水从20°C升高到100°C，求水吸收的热量是多少J [c_水=$4.2 \times 10^3 \text{J} / (\text{kg} \cdot \text{°C})$] 错误信息： 算式错误：1.计算水升高的温度差： 公式：[水升高的温度差]=[末温]-[初温] 算式=[末温]-[初温] 问题缺失参数：水升高的温度差=80 °C； output: 错误分析： 1.根据错误信息：算式存在错误，而且算式中存在“[参数]”的情况：“[末温]”、“[初温]”，需要对参数表增加缺失的参数根据错误信息，“[末温]-[初温]”，从题目中可以找到相关文本“从20°C升高到100°C”，按照要求的参数格式表示为： 初温 t_0 20 °C 末温 t 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公式可以得到： 算式=$(100 \text{°C}) - (20 \text{°C}) = 80 \text{°C}$ 水升高的温度差=80 °C 2.根据错误信息，问题缺失参数，由于分析1中纠正算式后计算得到了“水升高的温度差=80 °C”，所以问题不再缺失参数，不需要进行修改</p> <p>错误纠正： 参数表：null->初温 t_0 20 °C 参数表：null->末温 t 100 °C ...omitted... 下面是一个新的问题： 问题：{{question}} 错误：{{error}}</p>	<p>I need your help to correct the error in the explanation. I will provide the question and error information. The following are the requirements for error correction:</p> <ol style="list-style-type: none"> 1. You need to first conduct error analysis, analyze how to modify to correct the error, and then provide the error correction to correct the error in the explanation. 2. The error correction section does not require any additional explanatory information. The format of the error correction section is: "Content: Pre modified Content ->Modified Content". When adding content, "Pre modified Content" is null, and when deleting content, "Modified Content" is null. 3. Missing parameters in the question: If there are no missing parameters in the question, add the missing parameters to the question; If the parameters in the question have the same meaning as the missing parameters but different formats, modify the parameters in the question to be the same as the missing parameters. 4. Expression error: The formula and incorrect parameters need to be modified. If there is "[parameter]" or "null" in the expression, the missing parameters need to be filled in; If there are no issues with the parameters, it may be necessary to modify the formula. 5. The format of the formula is "[parameter to be solved]=[parameter 1] (+ - × ÷) [Parameter 2]..."; The format of the parameter table is: "concept symbol numeric unit", for example, "The boiling point of water is 100 °C", which is represented as "the boiling point of water t_boiling 100 °C" <p>Here are some examples:</p> <p>input: question: Assuming that 13.0 tons of bituminous coal is completely burned in a coal furnace, the heat released is partially absorbed by water, which can make $4 \times 10^5 \text{kg}$ of water is raised from 20 °C to 100 °C, how much heat is absorbed by the water ? [c_water=$4.2 \times 10^3 \text{J} / (\text{kg} \cdot \text{°C})$] Error message: Expression error: 1. Calculate the temperature difference of water rise: Formula: [temperature difference of water rise]=[final temperature] - [initial temperature] expression=[final temperature] - [initial temperature] Question missing parameter: Temperature difference of water rise=80 °C; output: 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 temperature]" and "[initial temperature]". It is necessary to add them to the parameter table. According to the error message, "[Last Temperature] - [Initial Temperature]", the relevant text "Increase from 20 °C to 100 °C" can be found in the title. According to the required parameter format, it is represented as: Initial temperature t_0 20 °C Final temperature t 100 °C After adding missing parameters to the parameter table, substitute them into the formula for calculating the temperature difference can be obtained as follows: expression=$(100 \text{°C}) - (20 \text{°C}) = 80 \text{°C}$ 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 be modified. Error correction: Parameter table: null ->Initial temperature t_0 20 °C Parameter table: null ->final temperature t 100 °C ...omitted... There's a new question: Question: {{question}} Error: {{error}}</p>

Figure 5: Prompt of the formula normalization stage 3: error correction for “calculation error”.

Question:
As shown in the figure, the Xuelong 2 scientific research icebreaker designed in China.
...*omitted*... When traveling at a constant speed of 3.6km/h in thick ice covered waters, the resistance experienced by the icebreaker is approximately $2 \times 10^7 \text{N}$. Calculate the propulsion power of the icebreaker at this time.
Reference answer: $2 \times 10^7 \text{ W}$

Formula:
[thrust]=[resistance]
[propulsion power]=[thrust] \times [constant speed]

Parameter table:

Parameter	symbol	value	unit
resistance	f	2×10^7	N
ship speed	v	1	m/s

Explanation:
1. Calculate thrust:
thrust=resistance= $2 \times 10^7 \text{N}$
2. Calculate propulsion power:
propulsion power=thrust \times constant speed= $2 \times 10^7 \text{N} \times$ constant speed(*cannot find value*)

Error:
1. The parameter "resistance" in the question is in the incorrect format.
2. "constant speed" could not be located in the parameter table.

Figure 6: An example of deleted questions.

Prompt actually used	English translation
<p>下面我会给出两个公式，每个公式由参数和运算符构成，[]中的表示参数。 你需要判断我给出的两个公式中对应参数表达含义是否相同，是否是同一个公式： 如果含义不相同，不是同一个公式，只需要回答不是； 如果各个参数含义相同，是同一个公式，则需要给出最终的公式，并且给出一个三行的表格来表示参数的对应关系，每个单元格内容是一个参数，前两行填写两个公式的参数，第三行填写统一后的公式参数。 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，是否是同一个公式：</p>	<p>I will give two formulas below. Each formula consists of parameters and operation symbols. The text in [] represent parameter. You need to judge whether the corresponding parameters in the two formulas I gave have the same meaning and whether they are the same formula: If the meaning is different, and they are not the same formula, just answer no; If each pair of parameters have the same meaning, and they are the same formula, the final formula needs to be given, and a three-row table needs to be given to indicate the corresponding relationship between the parameters. The content of each cell is a parameter, and the first two rows are filled with two formulas. Parameters, fill in the unified formula parameters in the third row. Here is formula 1: {formula 1} Here is formula 2: {formula 2} Judge whether they are the same formula by their meanings:</p>

Figure 7: Prompt for semantic-based merging.

Prompt actually used	English translation
这是一个初中物理题目，根据问题给出计算的过程，让我们一步一步地思考，在最后用“###”作为开始给出最终答案（一个数字）和答案的单位。 Question: {{问题}} Answer:	This is a junior high school physics question. Based on the given question, provide the calculation process and let's think step by step. Finally, use "###" to start giving the final answer (a number) and the unit of the answer. Question: {{question}} Answer:

(a) Zero-shot prompt for LLMs.

Prompt actually used	English translation
这是一个初中物理题目，根据问题给出计算的过程，用公式表示。 Question: {{样例1问题}} Answer: {{样例1解析}} ...omitted... Question: {{问题}} Answer:	This is a junior high school physics question. Based on the given question, provide the calculation process. Question: {{question of example 1}} Answer: {{explanation of example 1}} ...omitted... Question: {{question}} Answer:

(b) Few-shot prompt for LLMs.

Figure 8: Zero-shot and few-shot prompts for LLMs.

Prompt actually used	English translation
这是一个初中物理题目，根据问题给出计算的过程，用公式表示。 可能用到的公式有: {{top 5检索到的公式}} Question: {{问题}} Answer:	This is a junior high school physics question. Based on the given question, provide the calculation process. The formulas that may be used include: {{top 5 retrieved formulas}} Question: {{question}} Answer:

(a) Few-shot prompt for LLMs with formula retriever.

Prompt actually used	English translation
这是一个初中物理题目，根据问题给出计算的过程，用公式表示。 可能用到的公式有: {{用到的公式}} Question: {{样例1问题}} Answer: {{样例1解析}} ...omitted... 可能用到的公式有: {{top 5检索到的公式}} Question: {{问题}} Answer:	This is a junior high school physics question. Based on the given question, provide the calculation process. The formulas that may be used include: {{used formulas}} Question: {{question of example 1}} Answer: {{explanation of example 1}} ...omitted... The formulas that may be used include: {{top 5 retrieved formulas}} Question: {{question}} Answer:

(b) Zero-shot prompt for LLMs with formula retriever.

Figure 9: Zero-shot and few-shot prompts for LLMs with formula retriever.