Exploring Automated Distractor Generation for Math Multiple-choice Questions via Large Language Models

Anonymous ACL submission

Abstract

Multiple-choice questions (MCQs) are ubiquitous in almost all levels of education since they are easy to administer, grade, and are a reliable format in both assessments and practices. 004 005 An important aspect of MCQs is the distractors, i.e., incorrect options that are designed to target specific misconceptions or insufficient knowledge among students. To date, the task of crafting high-quality distractors largely remains a labor-intensive process for teachers 011 and learning content designers, which has limited scalability. In this work, we study the task of automated distractor generation in the domain of math MCQs and explore a wide va-015 riety of large language model (LLM)-based approaches, from in-context learning to fine-016 017 tuning. We conduct extensive experiments using a real-world math MCQ dataset and find 019 that although LLMs can generate some mathematically valid distractors, they are less adept at anticipating common errors or misconceptions among real students.

1 Introduction

027

Multiple-choice questions (MCQs) are widely used to evaluate students' knowledge because they enable quick and accurate administration and grading. MCQs are reliable because they are designed to measure specific learning objectives consistently (Nitko, 1996; Airasian, 2001; Kubiszyn and Borich, 2016). MCQs are constructed in a specific format; see Figure 1 for an example. The stem refers to the statement on the problem setup and context, followed by a question that needs to be answered. Among the options, the correct one can be referred to as the key, while incorrect ones can be referred to as distractors. As the name implies, distractors in MCQs are typically formulated to align with the common errors students would make or misconceptions students would exhibit. These distractors are chosen because students either i) lack the necessary

knowledge of the skills tested in the MCQ to accurately identify the key as the correct answer, or ii) hold misconceptions that result in selecting a specific distractor as the correct answer. While MCQs offer many advantages for students' knowledge evaluation, manually crafting high-quality MCQs is a demanding and labor-intensive process (Kelly et al., 2013). Specifically, high-quality distractors should be plausible enough to mislead students and not so evidently incorrect to be identified easily.

041

042

043

044

045

047

049

052

053

055

059

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

075

076

077

078

081

Prior work on automatic distractor generation primarily focuses on language learning and reading comprehension tasks, where distractors are used to assess students' comprehension of a given text or article. Early works use a ranking approach based on semantic similarity and word collocation information or a pre-defined ontology to produce distractors (Susanti et al., 2018; Stasaski and Hearst, 2017; Alsubait et al., 2014). More recent works use encoder-decoder models with attention mechanisms for distractor generation, resulting in longer and higher-quality distractors (Qiu et al., 2020; Shuai et al., 2023; Xie et al., 2021; Gao et al., 2019). Additionally, several recent works use pre-trained large language models (LLMs) such as BERT and T5 for distractor generation in the context of Swedish reading and Cloze test (Kalpakchi and Boye, 2021; Chiang et al., 2022; Rodriguez-Torrealba et al., 2022). Other works prompt LLMs such as ChatGPT and GPT-4 to generate distractors, either by providing detailed instructions or in-context examples in their prompts, for computer science course quiz questions and questions testing language mastery or factual knowledge (Tran et al., 2023; Bitew et al., 2023).

However, there is limited work on automatic distractor generation for math MCQs. This problem is more challenging than generating distractors for reading comprehension tasks because plausible distractors are not necessarily contained or can be inferred from the passage. A model for math

Stem

Write 35 as a fraction of 80. Answer in the simplest form.



Figure 1: Different parts of math MCQs and the terminology we use, illustrated with an example.

MCQ distractor generation should have some math problem-solving capability and more importantly, an understanding of the common errors or misconceptions among real students. Existing works either use constraint logic programming (Tomás and Leal, 2013) or manually constructed rules (Prakash et al., 2023) to generate distractors. However, these works only applies to math MCQs generated by templates. The work in (Dave et al., 2021) explores generating distractors using a neural network. However, their approach is training a math problem solver model and treating the incorrect outputs as distractors, which cannot capture common errors or misconceptions among real students.

1.1 Contributions

083

087

091

093

100

101

102

103

104

106

108

109

110

111

112

113

114

115

In this work, we investigate the task of automatically generating plausible distractors for math MCQs using LLMs. Our contributions include:

- We explore a variety of approaches to this task, including in-context learning, fine-tuning, and chain-of-thought prompting, together with rule- and sampling-based baselines.
- We conduct extensive quantitative and qualitative experiments on a real-world dataset of math MCQs. We find that the most effective approach is in-context learning, where we select a few example MCQs as input to the LLM, which can serve as a baseline for future work.
- We conduct a human evaluation and find that although the LLM-generated distractors are close to the human-authored ones in terms of mathematical validity, they do not necessarily reflect common errors or misconceptions among real students.

2 Task and Approaches

In this section, we first formally define relevant mathematical notation in MCQs and the automated distractor generation task. We then detail the LLMbased approaches and baselines that we explore.

2.1 Task Definition

We define an MCQ Q as a set of textual components, i.e., $Q = \{s, k, e_k, D, F\}$.¹ Each MCQ contains a stem s, a key k, an (optional) explanation of the key e_k , and a set of distractors D; each of which has an (optional) corresponding feedback message f_i which is shown to a student upon selecting a distractor $d_i \in D$. All of these components are sequences of words and math symbols (e.g., $s = \{w_1, \ldots, w_L\}$ where L is the length of the sequence s). Similar to (Qiu et al., 2020), we formulate the task of distractor generation as learning a function g^{dis} that outputs a set of distractors \hat{D} for an MCQ given the question stem and (optionally) key and its explanation, i.e.,

$$g^{\text{dis}}(s,k,e_k) \to \hat{D}.$$
 (1)

Our goal is to generate distractors that students with insufficient knowledge on skills required for the MCQ or specific misconceptions will select. This way, the MCQ can better distinguish between students that master all the required skills and those who do not. Below, we detail various LLM-based distractor generation approaches and several baselines that we explore.

We note that in this work, we study the problem of generating a set of distractors \hat{D} given a single 136

137

138

140

141

142

143

144

145

146

116

117

118

119

¹In this paper, we do not consider MCQs that contain diagrams or images; extending our work to multi-modal MCQ content is left for future work.



Figure 2: Overview of the kNN approach illustrated with a math MCQ on "compound percentage decrease".

question stem. This setting is different from a pos-147 sible alternative setting where we generate distrac-148 tors one-by-one, each corresponding to a common 149 error or misconception among real students. The 150 latter is applicable to the related problem of feed-151 152 back generation (Prihar et al., 2023), which investigates the task of generating a feedback message f_i for a distractor d_i . Providing feedback messages to students who select distractors can help them 155 identify their errors or misconceptions and guide 156 them towards the correct answer, which may expedite their learning process. In this work, we only 158 treat the feedback message as an additional rea-159 soning pathway to help LLMs generate plausible 160 distractors and do not study the quality of feedback 161 messages, which we leave for future work.

2.2 Approaches

The first approach is in-context learning or few-164 shot prompting, i.e., the LLM is expected to gen-165 erate desired outputs for a new task by learning 166 from the given examples (Brown et al., 2020). To 167 select examples, we select the k-nearest neighbor 168 (kNN) MCQs from a real-world math MCQ dataset, which we detail in Section 3.1, to the target MCQ. 170 After conducting tests with various values of k, we 171 find that this approach achieves the best distractor 172 generation performance when k = 3. To deter-173 mine similarity, we calculate the *cosine similarity* 174 between vectorized textual encodings of MCQs. Specifically, we use the pre-trained SBERT en-176 coder MPNet (Reimers and Gurevych, 2019) to calculate the textual encoding of the question stem 179 and (optionally) key and its explanation. Figure 2 provides a visual representation of this approach. The intuition for this approach is that MCQs with 181 similar question stems may have distractors that correspond to similar student errors or misconcep-183

tions that are feasible to the target MCQ, which may help the LLM to generate plausible distractors. Even though textual similarity may not be an appropriate representation for mathematical errors, these in-context examples should at least inform the LLM on distractor formatting (Chen et al., 2023; Lyu et al., 2023). We use ChatGPT in this approach for its proficiency in understanding tasks and delivering strong performance when provided with in-context examples. 184

185

186

187

188

189

190

191

192

193

194

195

196

197

198

199

200

202

203

204

205

206

207

208

209

210

211

212

213

214

215

216

217

218

219

220

The second approach is LLM fine-tuning (**FT**) to help pre-trained LLMs to adapt to the distractor generation task. We use the real-world math MCQ dataset to fine-tune the LLM in the format of Eq. 1, i.e., outputting all distractors given the question stem and (optionally) key and its explanation as input. We use ChatGPT (gpt-3.5-turbo-1106) (OpenAI, 2022), the largest base LLM that can be fine-tuned, in this approach.

The third approach is chain-of-thought prompting (**CoT**) (Wei et al., 2022). We provide the LLM with the question stem and (optionally) key and its explanation and detailed guidelines on distractor generation as input and ask it to first generate potential erroneous steps a student may take, followed by an incorrect answer, which we use as a distractor. This approach operates in a zero-shot manner and requires no access to any real MCQ data. Therefore, the performance depends solely on the LLM's ability in mathematical reasoning and anticipating students' errors or misconceptions. Given the demanding nature of this approach, we use a strong base LLM GPT-4 (OpenAI, 2023).

The fourth approach is a rule-based (**RB**) baseline, which can be used to generate different versions of the same MCQ with different numerical values. We emphasize that in many real-world ed-

ucational platforms, content creators do not use 221 rules to design distractors; in practice, not a lot of MCQs are created from templates and only differ by numerical values or named entities in their question stems. Therefore, we approximately follow the baseline approach in (Dave et al., 2021) and manually construct 444 distinct error explana-227 tions, such as "confuses factor and multiples" for question-distractor pairs that correspond to common errors or misconceptions among real students. 230 This process is extremely time-consuming and requires significant manual effort. We then provide the LLM with the question stem and (optionally) key and its explanation and a pool of error explana-234 tions that are feasible under the MCQ's topic (i.e., 235 fractions, rounding, etc.), and ask LLM to select 3 relevant ones and generate the corresponding distractors. We use GPT-4 in this approach for the same reason as CoT.

> The fifth approach is an improved version of the sampling-based (**SB**) baseline in (Dave et al., 2021). This approach fine-tunes a base LLM on MCQ answering, i.e., outputting the key given the question stem as input. Then, we randomly sample up to 20 output answers from the trained LLM given a question stem as input and choose 3 distinct incorrect ones as distractors. This approach implicitly assumes that LLMs make similar errors as real students. We use ChatGPT in this approach for the same reason as FT.

3 Experiments

240

241

242

243

244

245

246

247

248

251

255

263

265

269

In this section, we detail the specifics of our dataset, the evaluation metrics, the experimental setup, and report results from a series of quantitative, qualitative experiments, and human evaluation.

3.1 Dataset

Our dataset consists of 1.4K MCQs from a large digital learning platform, and all MCQs are written in English. We filter out questions with images/diagrams. Each question has 1 key and 3 distractors designed according to common student errors or misconceptions. The questions are sourced from the broad mathematical topic titled "Number" with subtopics including "Basic Arithmetic", "Fractions", and "Rounding and Estimating". The questions are primarily targeted towards students aged between 10 to 13. Each MCQ also has some additional metadata, e.g., the "topic" on 3 different granularity levels and the option selection distribution, i.e., the proportion of students who selected each option. The option selection distribution is computed on an average of 4000 student responses, with more than 900 student responses available in over 75% of the MCQs. We divided the dataset into two subsets, namely a training set and a test set, using an 80:20 ratio. We use the training set to select MCQs as in-context examples or fine-tune LLMs and the test set for evaluation.

3.2 Evaluation Metrics

Our main evaluation metric is a set of *alignmentbased* metrics, which quantifies the extent to which the LLM-generated distractors align with the human-authored ones. We denote the LLMgenerated distractors as \hat{D} where $|\hat{D}| = N$. We utilize 3 measures for this evaluation, two binary and one continuous. The binary metrics are **Exact** match h_e , i.e., whether all LLM-generated distractors match human-authored ones, and **Partial** match h_p , i.e., whether at least one LLM-generated distractor matches human-authored ones. These measures are formally defined as

$$h_e(D, \hat{D}) = \begin{cases} 1 & \forall \hat{d}_i \in \hat{D} : \hat{d}_i = d_i \\ 0 & \text{otherwise.} \end{cases}$$
292

and

1

$$h_p(D, \hat{D}) = \begin{cases} 1 & \exists \hat{d}_i \in \hat{D} : \hat{d}_i = d_i \\ 0 & \text{otherwise} \end{cases}$$
 294

We also use a continuous measure in the range [0, 1] that we call **Proportional** match h_n , i.e., the portion of LLM-generated distractors that match human-authored ones, defined as

$$h_n(D, \hat{D}) = \frac{\sum_{i=1} \mathbf{1}_{i:\hat{d}_i = d_i}}{N}$$
 299

where 1 denotes an indicator function. We report all metrics by averaging across all MCQs in the test set and scale the values of metrics by a factor of 100 into percentages.

Additionally, we experiment with a nonstandard, *distribution-based* metric, which tries to predict how often a distractor is selected by real students. This metric is motivated by the observation (See Section 3.5 for a detailed qualitative analysis) that human-authored distractors are sometimes not plausible or complete: for some MCQs, there may only be one highly common error or misconception among real students so teachers often have to throw 293

295

296

297

300

301

302

303

304

305

306

307

308

309

310

311

312

279

270

271

272

273

274

275

276

277

278

81

283

285

286

287

288

290

Approach	Exact	Partial	Proportional
kNN	9.89	72.44	37.10
CoT	4.24	63.96	29.92
RB	4.59	59.01	27.80
FT	2.83	57.60	25.32
SB	0.00	10.25	3.65

Table 1: Results on distractor generation on alignmentbased metrics, where in-context learning with kNN example selection outperforms other approaches.

in a few more that will be selected by almost no 313 314 one, while for other MCQs, there may be numerous plausible distractors that cannot all be included. 315 Therefore, our goal is to use the percentages of students who selected each option to train a model 317 that predicts how feasible a distractor is. Since we 318 cannot reach high predictive accuracy on the real 319 dataset we have, we relegate the details on this met-320 ric and experimental results to the Supplementary Material Section A. 322

3.3 Experimental Setup

323

325

326

327

329

332

333

339

For all approaches except SB, we use a uniform format to represent the target MCQ. This format comprises a concatenation of 3 elements: the question stem, key, and its explanation. We use this structure since it encapsulates the most comprehensive information about the target MCQ. Furthermore, based on CoT, we instruct the LLM to first generate feedback message and then the distractor, which intends to simulate a reasoning pathway, providing a scaffold that guides the subsequent generation of plausible distractors. We use greedy decoding and a maximum output length of 350 tokens for distractor generation. Additional hyperparameters and model details are in Supplementary Material Section B. We also provide our prompts for CoT, RB, and kNN in Tables 7, 8, and 9 respectively.

3.4 Results and Discussion

341Table 1 shows the results on distractor generation342for the 5 approaches we explore. Overall, kNN out-343performs the other approaches. This result is not344surprising since examples that are textually sim-345ilar to the target MCQ often contain distractors346that correspond to plausible errors or misconcep-347tions among real students for both MCQs. There-348fore, the LLM can generate distractors that match349the human-authored ones by simply replicating the350style of the in-context examples. This approach is

Approach	Exact	Partial	Proportional
kNN ^{all}	9.89	72.44	37.10
kNN ^{key}	10.95	69.26	36.75
kNN ^{none}	8.83	66.08	34.39
Random	2.12	54.77	23.44
Prompt ^{key}	8.13	65.72	33.33
Prompt ^{none}	2.83	36.04	16.96
$kNN^{all}_{\neg T}$	3.20	57.60	26.15
FT ^{gpt3.5}	2.83	57.60	25.32
FT ^{llama2}	0.35	40.99	16.02
RB ^{select}	4.59	59.01	27.80
RB ^{random}	1.06	52.65	23.20

Table 2: Results on ablation study on alignment-based metrics with different settings of kNN, FT, and RB.

especially effective for MCQs that have highly similar structures and differ in only numerical values. 351

352

353

354

356

357

358

360

361

362

363

364

365

366

367

368

369

370

371

372

373

374

375

376

377

378

381

383

The advantage of CoT over FT reflects the strong mathematical reasoning capability of GPT-4, which results in a performance gap that not even fine-tuning ChatGPT on human-authored distractors can make up. Instead of acting as an oracle, RB underperforms this expectation and does not even outperform CoT, despite requiring significant human expertise and effort. This result is likely due to the fact that despite extensive effort in labeling error explanations, we cannot come up with a comprehensive list of them; as a result, many target MCQs are not matched with error explanations for GPT-4 to select from. Overall, we observe that GPT-4 can often generate mathematically valid distractors but is unaware of what errors or misconceptions are common among real students. Therefore, CoT and RB do not perform as well as kNN. Among all approaches we explored, SB has by far the worst performance, which is not surprising, since when we train LLMs to answer math MCQs correctly, the incorrect answers they generate are either only marginally different than the key or completely unrelated to the question stem. Therefore, this approach generates distractors without coherent reasoning and do not resemble how real students make mistakes.

3.4.1 Ablation Study

In this ablation study, we investigate the impact of different configurations of kNN on its performance and summarize these results in the first part of Table 2. We explore how different ways of using different parts of the MCQ in the textual encoder

for nearest neighbor search could affect kNN's performance. We experiment with 3 different settings: using just the question stem (kNN^{none}); using the question stem and key (kNN^{key}); and using the question stem, key, and its explanation (kNN^{all}), which is the best performing setting. For comparison, we also experiment with a simple random 391 heuristic (Random) that chooses examples from the training set randomly without any specific criteria. We see that although using only the ques-394 tion stem captures the math skill covered by an MCQ and helps kNN find examples that have the 396 same format as the target MCQ, adding the key and explanation helps kNN find better examples that use similar problem-solving strategies to the target MCQ. We also explore how different prompt 400 formats could affect kNN's performance. We ex-401 periment with 3 different prompt formats. The best-402 performing setting (kNN^{all}) includes the question 403 stem, key, and explanation for both the target MCQ 404 and the in-context examples. The in-context exam-405 ples also contain feedback on the distractors, and 406 we ask the LLM to generate the feedback, followed 407 by the distractor. The other settings are to not in-408 clude feedback messages for the distractors and 409 the explanation for the key (**Prompt**^{key}), and not 410 including the key either (**Prompt**^{none}). We see that 411 including the key significantly improves kNN's per-412 formance and asking the LLM to generate feedback 413 followed by the distractor further improves perfor-414 mance. This result again reinforces the importance 415 of math problem-solving strategies and CoT reason-416 ing on the distractor generation performance. We 417 also explore the impact of not allowing MCQs with 418 the same topic to be selected as examples on kNN's 419 performance (**kNN**^{all}_{$\neg T$}). We see that doing so re-420 sults in a huge performance drop-off from kNN^{all}. 421 This result suggests that most errors or misconcep-422 tions behind distractors are topic-specific and do 423 not generalize across topics. 424

Next, we investigate the impact of different 425 base LLMs on FT's performance and summarize 426 these results in the second part of Table 2. We 427 compare ChatGPT against LLAMA2-7B (Tou-428 vron et al., 2023), which is one of the biggest 429 open-sourced generative LLMs (FT^{llama2}). We see 430 that ChatGPT outperforms LLAMA2-7B on all 431 432 3 alignment-based metrics. This result suggests that larger models that are better at mathematical 433 reasoning are more likely to generate plausible dis-434 tractors. We also investigate the impact of different 435 error selection approaches on RB's performance 436

Target

Quesiton stem: which multiplier can be used to find the value after an amount has decreased in value by 8% for 4 years?

Explanation: As its is a decrease, we need 100% - 8% which is 92% which is the same as 0.92. We then use the number of years as the power of 4.

Answer: $\times 0.92^4$

Example 1

Quesiton stem: which multiplier can be used to find the value after an amount has decreased in value by 5% for 5 years?

Explanation: As its is a decrease, we need 100% - 5% which is 95% which is the same as 0.95. We then use the number of years as the power of 5.

Answer: $\times 0.95^5$

Example 2

Quesiton stem: the value of a laptop that initially cost \$1100, declines in value by 15% a year. if you wanted to calculate the value of the tablet at the end of 6 years, what number would replace the square? $1100 \times \Box^6$

Explanation: As the value decreases by 15%, we have 100% - 15% = 85% = 0.85 as the multiplier.

Answer: 0.85

Example 3

Quesiton stem: a car depreciates in value by 15% each year. if a car was bought for \$3500, which of the following calculations would find the new value of the car after 3 years? Explanation: The multiplier is 1 - 0.15 = 0.85, and as we are using compound interest, we raise this to the power of 3. Answer: 3500×0.85^3

Table 3: Three in-context learning examples retrieved by kNN; we see that Example 1 is very similar to the target MCQ, except for different numerical values.

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

and summarize these results in the third part of Table 2. We experiment with a variant of RB that randomly selects error explanations under the same math topic (\mathbf{RB}^{random}) instead of asking GPT-4 to select 3 relevant ones (\mathbf{RB}^{select}). We see that asking the LLM to select error explanations outperforms selecting error explanations randomly, but not by a significant margin compared to other ablations. This result suggests that even though LLMs can generate many mathematically valid distractors, their ability to recognize which error explanations are popular among students is limited.

3.5 Qualitative Analysis

We now qualitatively investigate the distractors generated by the best approach, kNN, to extract some insights on the distractor generation task and how to improve performance. We group the 283 total

473

474

475

476

477

478

479

480

481

482

483 484

485

486

487

488

454

MCQs in the test set into 4 categories, according to the number of LLM-generated distractors that match the human-authored ones, from 0 to 3.

For the group where all LLM-generated distractors match the human-authored ones (3 out of 3), we find that, in all but 2 of the 28 such cases, there is an in-context example that is very similar to the target MCQ, with the only difference being different numerical values or named entities. See Table 3 for an example. However, this situation sometimes appears in other groups too, which is perhaps surprising since it implies that the presence of a near-identical in-context example alone is not sufficient for an LLM to generate plausible distractors. We investigate further into such cases and find that even for two MCQs with near-identical question stem, their sets of distractors and the errors or misconceptions underlying each distractor may differ even though both are plausible. This situation occurs when there are more than 3 plausible errors or misconceptions given a question stem.

Question Stem

Craig and Isaac share some fruit. Isaac gets threequarters of the fruit. In what ratio do they share the fruit? (Isaac's part second)

Key				
1:3				
LLM-	generat	ed Distra	ctors	
3:1	3:4	4:1		
Human-authored Distractors				
1:4	1:2	4:3		

Table 4: Example of LLM-generated distractors that are mathematically valid and plausible but do not match human-authored ones.

For the group where none of the LLM-generated distractors match the human-authored ones, we randomly select 20 of the 78 cases to analyze. We find that in 14 of the 20 cases (70%), the LLM-generated distractors are plausible and the human-authored ones are not superior to the LLMgenerated distractors. See Table 4 for an example. While this observation is entirely subjective, it highlights that alignment-based metrics may not be an appropriate metric to measure the quality of LLMgenerated distractors because human-authored ones may not be optimal. This observation is also part of our motivation in developing *distribution-based* metrics to predict how likely a LLM-generated distractor will be selected by real students with insuf-
ficient knowledge. Moreover, since many LLM-
generated distractors are valid and plausible even
if they are not the same with the human-authored
ones, there is promise in using automated distractor
generation for teacher support during the genera-
tion of MCQs.489
490
490
490
490
491

Question StemConvert 0.6 to a fraction in its simplest form.Key $\frac{3}{5}$ $\frac{3}{5}$ LLM-generated Distractors $\frac{6}{10}$ $\frac{5}{3}$ $\frac{6}{5}$ Human-authored Distractors $\frac{6}{10}$ $\frac{60}{100}$ $\frac{1}{6}$

Table 5: Example of LLM-generated distractors where the plausible one, $\frac{6}{10}$ matches the human-authored ones, while the rest of human-authored ons are placeholders. In this case, $\frac{6}{10}$ is selected by 28% of students while other distractors are rarely being selected.

495

496

497

498

499

500

501

502

503

504

505

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

Finally, for the group where 1 or 2 LLMgenerated distractors match the human-authored ones, we examine which human-authored distractor(s) are generated and which are not. We find that in many cases, the generated distractors that match to the human-authored ones seem to contain typical errors or misconceptions related to the question stem, while the other human-authored ones are not. See Table 5 for an example. This observation is further supported by selections made by real students, where the distractor that corresponds to the typical error or misconception is the one most often selected by students in 44 of 108 (40.7%) cases and 46 of 63 (73%) cases for 1 and 2 matches, respectively, while the rest are rarely being selected. This result suggests that many MCQs have only one highly plausible distractor while the others are placeholders. Again, using human-authored ones as the ground truth on alignment-based metrics is not ideal, which justifies our motivation in developing the distribution-based metric.

3.6 Human Evaluation

We conduct a human evaluation to assess the quality of LLM-generated distractors. This evaluation is motivated by observations from the qualitative analysis that the generated distractors are often mathematically valid even though they may differ

	QWK		Average Ratings	
	LLM	Human	LLM	Human
Validity Plausibility	0.34 0.54	0.23 0.54	3.28 2.68	3.99* 3.72*

Table 6: QWK and average ratings among human evaluators on LLM-generated and human-authored distractors for validity and plausibility. Under a Student's t-test, human evaluators prefer human-authored distractors with statistical significance ($p < 0.05^*$).

from human-authored ones.

523

524

525

527

529

530

531

533

535

536

537

538

541

542

544

547

548

549

550

551

552

553

555

557

558

561

3.6.1 Evaluation Design

We recruited 2 graduate students who have experience teaching math or related topics as human evaluators. They were presented with the same set of 20 MCQs that were randomly sampled from the test set, each accompanied by a mixture of 4 or 6 distractors. To ensure a balanced assessment, half of these were LLM-generated distractors, while the remaining were human-authored ones. To eliminate any potential ordering bias, the sequence of the distractors was randomized for each question. They were asked to rate the distractors on two aspects: mathematical validity (validity) and plausibility for middle school math students (plausibility). Validity measures whether a distractor is relevant to the question stem and can be tangibly reached by some incorrect reasoning. Plausibility measures how likely a distractor is to be selected by real students. Each aspect is scored on a scale from 1 to 5, with 1 being the lowest: a distractor that is irrelevant to the question stem or one that no student would select, while 5 being the highest: a distractor that is highly relevant to the question stem or one that is highly likely to trick students with insufficient math skills into selecting it. Additional evaluation setup details are in Supplementary Material Section C.

3.6.2 Evaluation Result and Discussion

Table 6 shows the inter-rater agreement, measured using quadratic weighted Kappa (QWK) (Brenner and Kliebsch, 1996) and the average rating, across 2 human evaluators for both LLM-generated and human-authored distractors. The QWK scores indicate a fair to moderate level of agreement between two human evaluators regarding both the validity and plausibility aspects of distractors. This observation suggests that measuring the quality of distractors based on their validity and plausibility is consistent at certain level and can be used in future assessments of distractors. We conduct a Student's t-test (Semenick, 1990) to compare the ratings for LLM-generated and human-authored distractors and find that on both aspects, there is a statistically significant difference (p < 0.05). This result shows that human evaluators generally prefer human-authored ones over LLM-generated distractors on both aspects. Furthermore, we observe that the gap between validity and plausibility is much higher for LLM-generated distractors than for human-authored ones. This observation indicates that LLMs exhibit a higher proficiency in generating mathematically valid distractors compared to understanding what errors or misconceptions are plausible among real students. This result is not surprising since LLMs, which have not been extensively trained on erroneous answers provided by real students, may struggle to understand students' misconceptions or the various ways in which students are prone to making errors. Therefore, there is still considerable room for improvement for LLMs in their capacity to understand errors or misconceptions among real students.

562

563

564

565

566

567

568

569

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586

587

589

590

591

592

593

594

595

596

597

598

599

600

601

602

603

604

605

606

607

608

609

610

4 Conclusions and Future Work

In this paper, we explore automated distractor generation for math multiple-choice questions via large language models. We conduct experiments on a real-world math MCQ dataset and find that the incontext learning-based approach kNN, achieves the best performance when compared to other approaches such as fine-tuning, chain-of-thought prompting, and various baselines. We also conduct human evaluation and observe that LLMs are capable of generating mathematically valid distractors but are not fully aware of common errors or misconceptions among real students. Our initial exploration of this task opens up many avenues for future work. For example, we need to further refine the distribution-based metrics that predict the percentage of students who select each distractor. We also need to develop modified text encoding approaches that are closely aligned with errors or misconceptions among real students for in-context example selection. Furthermore, we aim to explore the generation of distractors, each of which corresponds to a specific error or misconception, as well as the generation of high-quality feedback messages for each distractor.

Limitations 611

Being the attempt at the task of generating plausible distractors for math MCQs using LLMs, we 613 find several limitations in our current setup. First, 614 our best approach, kNN, is constrained to generalize beyond the range of topics represented within 616 the question pool. This is due to its reliance on 617 selecting in-context examples that have the closest 618 semantic meaning of the question stem, key, and 619 its explanation. Second, we find that some humanauthored distractors have low quality, and using 621 them as demonstrations may lead to the generation of distractors that do not contain common student 623 errors or misconceptions to effectively evaluate the students' knowledge. Third, the alignment-based 625 metrics may not accurately measure the quality of generated distractors because some MCQs may have more than 3 plausible distractors.

Ethical Considerations

The focus of our work is to automatically generate plausible distractors for math MCQs using 631 LLM. By automating part of the MCQ generation, 632 we aim to save educators and teachers from timeconsuming MCQ generation and allow them to dedicate more effort to teaching and student engagement. Based on our analysis on the generated distractors, we acknowledge that not every distractor generated by our work is plausible. Therefore, we strongly advise that our work should be adopted 639 as an auxiliary tool in the generation of MCOs. All automatically generated distractors should undergo a careful review by educators and teachers before being utilized in real tests for students.

References

647

649

651

652 653

654

658

- Peter Airasian. 2001. Classroom assessment: Concepts and applications. McGraw-Hill, Ohio, USA.
- Tahani Alsubait, Bijan Parsia, and Uli Sattler. 2014. Generating multiple choice questions from ontologies: Lessons learnt. In OWLED, pages 73-84. Citeseer.
- Semere Kiros Bitew, Johannes Deleu, Chris Develder, and Thomas Demeester. 2023. Distractor generation for multiple-choice questions with predictive prompting and large language models. arXiv preprint arXiv:2307.16338.
- Hermann Brenner and Ulrike Kliebsch. 1996. Dependence of weighted kappa coefficients on the number of categories. Epidemiology, pages 199-202.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie 659 Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind 660 Neelakantan, Pranav Shyam, Girish Sastry, Amanda 661 Askell, et al. 2020. Language models are few-shot 662 learners. Advances in neural information processing 663 systems, 33:1877–1901. 664 Yulin Chen, Ning Ding, Xiaobin Wang, Shengding Hu, 665 Hai-Tao Zheng, Zhiyuan Liu, and Pengjun Xie. 2023. 666 Exploring lottery prompts for pre-trained language 667 models. arXiv preprint arXiv:2305.19500. 668 Shang-Hsuan Chiang, Ssu-Cheng Wang, and Yao-669 Chung Fan. 2022. Cdgp: Automatic cloze distractor 670 generation based on pre-trained language model. In 671 Findings of the Association for Computational Lin-672 guistics: EMNLP 2022, pages 5835-5840. 673 Paul F Christiano, Jan Leike, Tom Brown, Miljan Mar-674 tic, Shane Legg, and Dario Amodei. 2017. Deep 675 reinforcement learning from human preferences. Ad-676 vances in neural information processing systems, 30. 677 Neisarg Dave, Riley Bakes, Barton Pursel, and C Lee 678 Giles. 2021. Math multiple choice question solv-679 ing and distractor generation with attentional gru 680 networks. International Educational Data Mining 681 682 Yifan Gao, Lidong Bing, Piji Li, Irwin King, and 683 Michael R Lyu. 2019. Generating distractors for 684 reading comprehension questions from real exami-685 nations. In Proceedings of the AAAI Conference on 686 Artificial Intelligence, volume 33, pages 6423-6430. 687 Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan 688 Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, 689 and Weizhu Chen. 2021. Lora: Low-rank adap-690 tation of large language models. arXiv preprint 691 arXiv:2106.09685. 692 Albert Q Jiang, Alexandre Sablayrolles, Arthur Men-693 sch, Chris Bamford, Devendra Singh Chaplot, Diego 694 de las Casas, Florian Bressand, Gianna Lengyel, Guil-695 laume Lample, Lucile Saulnier, et al. 2023. Mistral 696 7b. arXiv preprint arXiv:2310.06825. 697 Dmytro Kalpakchi and Johan Boye. 2021. Bert-based 698 distractor generation for swedish reading compre-699 hension questions using a small-scale dataset. arXiv 700 *preprint arXiv:2108.03973*. 701 Kim Kelly, Neil Heffernan, Sidney D'Mello, Namais 702 Jeffrey, and Amber C. Strain. 2013. Adding teacher-703 created motivational video to an its. In Proceedings 704 of 26th Florida Artificial Intelligence Research Soci-705 ety Conference, pages 503-508. 706 Tom Kubiszyn and Gary Borich. 2016. Educational 707 testing and measurement. John Wiley & Sons, New 708 Jersey, USA. 709 Xinxi Lyu, Sewon Min, Iz Beltagy, Luke Zettlemoyer, 710 and Hannaneh Hajishirzi. 2023. Z-ICL: Zero-shot 711 in-context learning with pseudo-demonstrations. In 712

Society.

713

715

- 736 737 740 741 742 743 744 745 747 750 751
- 753
- 755 756 757 758 759
- 761 762

- Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 2304–2317.
- Anthony J. Nitko. 1996. Educational assessment of students. Prentice-Hall, Iowa, USA.
- OpenAI. 2022. Introducing chatgpt.
 - OpenAI. 2023. Gpt-4 technical report.
 - Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, et al. 2011. Scikit-learn: Machine learning in python. the Journal of machine Learning research, 12:2825-2830.
 - Vijay Prakash, Kartikay Agrawal, and Syaamantak Das. 2023. Q-genius: A gpt based modified mcq generator for identifying learner deficiency. In International Conference on Artificial Intelligence in Education, pages 632-638. Springer.
 - Ethan Prihar, Morgan Lee, Mia Hopman, Adam Tauman Kalai, Sofia Vempala, Allison Wang, Gabriel Wickline, Aly Murray, and Neil Heffernan. 2023. Comparing different approaches to generating mathematics explanations using large language models. In International Conference on Artificial Intelligence in Education, pages 290-295. Springer.
 - Zhaopeng Qiu, Xian Wu, and Wei Fan. 2020. Automatic distractor generation for multiple choice questions in standard tests. arXiv preprint arXiv:2011.13100.
 - Nils Reimers and Iryna Gurevych. 2019. Sentence-bert: Sentence embeddings using siamese bert-networks. In Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing. Association for Computational Linguistics.
 - Ricardo Rodriguez-Torrealba, Eva Garcia-Lopez, and Antonio Garcia-Cabot. 2022. End-to-end generation of multiple-choice questions using text-to-text transfer transformer models. Expert Systems with Applications, 208:118258.
 - Doug Semenick. 1990. Tests and measurements: The ttest. Strength & Conditioning Journal, 12(1):36–37.
 - Pengju Shuai, Li Li, Sishun Liu, and Jun Shen. 2023. Qdg: A unified model for automatic questiondistractor pairs generation. Applied Intelligence, 53(7):8275-8285.
 - Katherine Stasaski and Marti A Hearst. 2017. Multiple choice question generation utilizing an ontology. In Proceedings of the 12th Workshop on Innovative Use of NLP for Building Educational Applications, pages 303-312.
 - Yuni Susanti, Takenobu Tokunaga, Hitoshi Nishikawa, and Hiroyuki Obari. 2018. Automatic distractor generation for multiple-choice english vocabulary questions. Research and practice in technology enhanced *learning*, 13:1–16.

Ana Paula Tomás and José Paulo Leal. 2013. Automatic generation and delivery of multiple-choice math quizzes. In Principles and Practice of Constraint Programming: 19th International Conference, CP 2013, Uppsala, Sweden, September 16-20, 2013. Proceedings 19, pages 848-863. Springer.

767

768

770

773

774

776

779

780

781

782

783

784

785

786

787

788

789

790

791

792

793

794

795

796

797

798

799

800

801

802

803

804

805

806

- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. arXiv preprint arXiv:2307.09288.
- Andrew Tran, Kenneth Angelikas, Egi Rama, Chiku Okechukwu, David H Smith IV, and Stephen Mac-Neil. 2023. Generating multiple choice questions for computing courses using large language models.
- Pauli Virtanen, Ralf Gommers, Travis E Oliphant, Matt Haberland, Tyler Reddy, David Cournapeau, Evgeni Burovski, Pearu Peterson, Warren Weckesser, Jonathan Bright, et al. 2020. Scipy 1.0: fundamental algorithms for scientific computing in python. Na*ture methods*, 17(3):261–272.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. 2022. Chain-of-thought prompting elicits reasoning in large language models. abs/2201.11903.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. 2019. Huggingface's transformers: State-of-the-art natural language processing. abs/1910.03771.
- Jiayuan Xie, Ningxin Peng, Yi Cai, Tao Wang, and Qingbao Huang. 2021. Diverse distractor generation for constructing high-quality multiple choice questions. IEEE/ACM Transactions on Audio, Speech, and Language Processing, 30:280–291.

808

809

812

813

814

815

816

818

821

823

824

825

830

832

833

837

838

839

841

845

846

847

849

853

854

857

Supplementary Material

A Distribution ranking metric

Since our qualitative analysis in Section 3.5 found that human-authored distractors are sometimes unplausible or incomplete, using them as the ground truth is not ideal. Therefore, we explore a distribution-based metric to evaluate the quality of LLM-generated distractors, based on one intuition: good distractors are ones that are likely going to be selected by many real students. Therefore, our goal is to train a model that can predict the portion of students that select each option in an MCQ. However, due to the highly noisy nature of this distribution, we opt to train a model that predicts the more often selected distractor among a pair, given a question stem, which is similar to the pairwise preference reward model in reinforcement learning from human feedback (RLHF) (Christiano et al., 2017). After training such a model, we can use it to compare generated distractors to human-authored ones in head-to-head matchups, giving us a proxy for how good an LLM is in terms of generating distractors that are likely to be selected by students.

Formally, we train an LLM-based model $\mathbf{r}_{\phi}(d_1, d_2, s, k, e_k) \rightarrow \{d_1, d_2\},$ where ϕ denotes the set of model parameters. We train this model by first constructing a dataset of all pairs of humanauthored distractors for each MCQ and include both orders of each pair to avoid ordering bias, resulting in $N \times \binom{3}{2} \times 2$ total pairs, where N denotes the number of MCQs. Each pair is associated with a binary-valued label indicating whether d_1 or d_2 is selected by more students, which we can calculate from the student response records in our dataset. We then use this dataset to fine-tune an LLM in a text generation task, where the LLM receives the question and distractor information in its prompt and outputs its preference. We show our prompt for this task in Table 10.

We use the same train/test split as the distractor generation experiments, and reserve 20% of the train split for validation after each epoch and early stopping. We fine-tune the mistralai/Mistral-7B-v0.1 (Jiang et al., 2023) model, which contains 7 billion parameters, from HuggingFace (Wolf et al., 2019) using LoRA (Hu et al., 2021) with adaptors on the q_proj, k_proj, v_proj, and o_proj matrices, set r = 32, $\alpha = 16$, dropout = 0.05, and use 8-bit quantization. We train the model using the AdamW optimizer for 10 epochs with a learning rate of 3e-5, a batch size of 16, accumulate gradients for 4 batches. The model converges on the validation set after 6 epochs. The GPU we use to train the model is NVIDIA RTX A6000. The training process is completed in 10 hours. When evaluated on the test set, the ranking model correctly identifies the preferred distractor 61.60% of the time (random guessing corresponds to 50% accuracy). This accuracy is low overall but high on subsets of distractor pairs whose student selection percentages differ by a large margin: on pairs with a larger than 20% margin, which accounts for 6% of pairs, the accuracy jumps to 74.47%. This result is not surprising since the selection percentage data is very noisy. 858

859

860

861

862

863

864

865

866

867

868

869

870

871

872

873

874

875

876

877

878

879

880

881

882

883

888

889

890

891

892

893

894

895

896

897

898

899

900

901

Using this trained model, we can evaluate the quality of LLM-generated distractors: we compare all possible head-to-head matchups between generated distractors and human-authored ones, and record the portion of times that the generated distractors are preferred by the ranking model. If the two distractors are the same then we record a tie. In cases where the generated distractors are invalid or repeated, we treat them as null and record a win for the human-authored ones. Formally, we define a preference score as

$$s = \frac{1}{18N} \sum_{i=1}^{N} \sum_{a=1}^{3} \sum_{b=1}^{3} r_{\phi}^{(i)}(\hat{d}_{a}^{(i)}, d_{b}^{(i)})$$
884

$$+ (1 - \mathbf{r}_{\phi}^{(i)}(d_b^{(i)}, \hat{d}_a^{(i)})), \qquad 885$$

$$\mathbf{r}_{\phi}^{(i)}(d_1, d_2) = \begin{cases} 0.5 & a_1 = a_2 \\ 1 & d_2 \text{ is null} \\ 0 & d_1 \text{ is null} \\ p & \text{otherwise} \end{cases},$$
886

$$p = \mathbf{1}_{\mathbf{r}_{\phi}(d_1, d_2, s^{(i)}, k^{(i)}, e_k^{(i)}) = d_1},$$
887

where d are generated distractors. This score has a range of [0, 1] where higher values indicate LLMgenerated distractors are likely to be selected by more students than the human-authored ones. We found that kNN scores 0.46 on the test set, which indicates that the distractors it generates are almost as plausible to students as human-authored ones.

We emphasize that this evaluation metric should only be considered exploratory due to several obvious limitations. First, student option selection percentages create noisy labels for the ranking model, limiting its accuracy. Second, using the overall selection percentages also ignores the individual learning context of each student since students with different knowledge levels may have different tendencies among MCQ options. Therefore, we leave
a more thorough treatment of the distribution-based
metric to future work.

906 907

937

938

939

944 945

946

947

948

B Hyperparameters and Implementation Details

We fine-tune the LLAMA2-7B model, which con-908 tains 7 billion parameters, from HuggingFace using 909 LoRA with adaptors on the q_proj and v_proj 910 matrices, set r = 8, $\alpha = 32$, dropout = 0.05, 911 and use 8-bit quantization. We use 20% of the 912 training set for validation. We train the model us-913 ing the AdamW optimizer for 15 epochs with a 914 learning rate of 3e-4, a batch size of 16, and ac-915 cumulated gradients for 16 batches. The model 916 converges on the validation set after 12 epochs. 917 The selection of the aforementioned hyperparam-918 eters is guided by exploratory evaluations and no 919 substantial hyper-parameter search is conducted. 920 The GPU we use to train the model is NVIDIA 921 RTX A6000. The training process is completed 922 in 3 hours and 43 minutes. We fine-tune the ChatGPT model using the first 200 data points 924 from the training set. We train the model using 925 the OpenAI's default fine-tuning settings, which we found providing the best performance, via Ope-927 928 nAI API. The training process is completed in 20 minutes. We use the scikit-learn (Pedregosa et al., 2011) implementation to calculate QWK, and 930 use the scipy (Virtanen et al., 2020) implemen-931 tation to calculate Student's t-test. For prompting 932 GPT-4 and ChatGPT using OpenAI API, we use temperature = 0, max_tokens = 350, top_p = 1, 934 frequency_penalty = 0.0, presence_penalty = 0.0 as our setup for greedy decoding.

> All our experiments are implemented in Python or Pytorch code, and We note that all software employed in this work is open-source, or the license is unspecified.

C Human Evaluation Details

In this work, we obtained approval from the ethics review board for human evaluation. We show the evaluation instructions to human evaluators in Table 11. We do not provide any compensation for human evaluators because their participation is entirely voluntary and we appreciate their contribution to this work.

D Prompt Format

We provide the prompts for CoT, RB, and kNN in the work below. We use <> to indicate that a variable is filled in dynamically.

Prompt	You are given the following math question along with the correct answer and explanation. Please use the following template to give 3 alternative incorrect answers to be used as multiple-choice options in a multiple-choice exam. Prior to the incorrect answer, provide feedback to be displayed to the student as an explanation of why that is not the correct answer. [Template] Distractor1 Feedback: Distractor2 Feedback: Distractor2: Distractor3 Feedback: Distractor3: Question: <question> Explanation: <explanation> Answer: <answer></answer></explanation></question>
	Table 7: CoT prompt format
Prompt	You are given the following math question along with the correct answer, explanation, and a list of errors. Please follow the template to first select 3 most likely errors for this question and use the se- lected errors to generate 3 alternative incorrect answers to be used as multiple-choice options in a multiple-choice exam. Prior to the incorrect answer, provide feedback to be displayed to the student as an explanation of why that is not the correct answer. If the list of errors is not given, generate 3 errors instead and do not contain any explanation in the 3 incorrect answer. [Template] Error1: Error2: Error3: Distractor1 Feedback: Distractor1 Feedback: Distractor2 Feedback: Distractor2 Feedback: Distractor3 Fe

Table 8: RB prompt format

Prompt Q E A D D D D D D D D D D D D D D D D D D	Question: <in-context question=""> Explanation: <in-context explanation=""> Answer: <in-context answer=""> Distractor1 Feedback: <in-context distractor1="" feedback=""> Distractor1:<in-context distractor1=""> Distractor2 Feedback: <in-context distractor2="" feedback=""> Distractor2:<in-context distractor2=""> Distractor3 Feedback:<in-context distractor3="" feedback=""> Distractor3:<in-context distractor3=""> stop] Question: <target question=""> Explanation: <target explanation=""> Answer: <target answer=""></target></target></target></in-context></in-context></in-context></in-context></in-context></in-context></in-context></in-context></in-context>
--	--

Table 9: kNN prompt format, in practice, we use 3 in-context examples

E Ranking Metric Examples

Prompt	A teacher assigns the following math multiple choice question to a class of middle school students.
	Question: $\frac{3}{5}$ of $50 = \frac{6}{10}$ of \square Correct Answer: 50 Solution: 3/5 and 6/10 are equivalent, so 3/5 of 50 is the same as 6/10 of 50.
	Here are 2 incorrect options that some students choose: Option A: 30 Option B: 18 Which incorrect option are the students more likely to pick?
Output	Preferred Answer: A

Table 10: Example prompt and output for the ranking model used in the distribution ranking metric.

F Instruction

You are given a csv file. Each row corresponds to a question stem and a distractor.

Your job is to rate the distractor on two aspects: mathematical validity and plausibility for middle school math students.

Mathematical validity measures whether a distractor is relevant to the question stem and can be tangibly reached by some incorrect reasoning. Mathematical validity is scored on a scale from 1 to 5, where 1 indicates a distractor that is irrelevant to the question stem, and 5 indicates a distractor that is highly relevant to the question stem.

Plausibility measures how likely a distractor is to be selected by middle school students learning math. Plausibility is scored on a scale from 1 to 5, where 1 indicates that no student would select it and 5 indicates that the distractor is highly likely to trick students with insufficient math skills into selecting it.

please use numbers on mac to rate distractors and give 1 and 1 for both metric if the distractor is the correct answer.

Your ratings will be used to quantitatively measures and analyzes the quality of distractors on validity and plausibility.

Table 11: Instruction for Human Evaluation