TREECUT: A Synthetic Unanswerable Math Word Problem Dataset for LLM Hallucination Evaluation

Anonymous ACL submission

Abstract

Large language models (LLMs) now achieve near-human performance on standard math word problem benchmarks (e.g., GSM8K), yet their true reasoning ability remains disputed. A key concern is that models often produce confident, yet unfounded, answers to unanswerable problems. We introduce TREECUT, a synthetic dataset that systematically generates in*finite* unanswerable math word problems and their answerable counterparts, by representing each question as a tree and removing chosen necessary conditions. Experiments show TREE-CUT effectively induce hallucinations in large language models, including GPT-40 and o3mini, with rates of 61% and 42% in their respective worst-case scenarios. Further analysis highlights that deeper or more complex trees, composite item names, and removing necessary condition near the middle of a path all increase the likelihood of hallucinations, underscoring the persistent challenges LLMs face in identifying unanswerable math problems.

1 Introduction

011

013

017

019

021

037

041

Mathematical reasoning is a crucial part of human intelligence. Recent years have witnessed remarkable advancements in the mathematical reasoning capabilities of large language models (LLMs). By leveraging techniques such as chain-of-thought prompting (Wei et al., 2022), state-of-the-art LLMs (e.g., Achiam et al. (2023); Team et al. (2024); Dubey et al. (2024)) achieved human-level performance on benchmarks like GSM8K (Cobbe et al., 2021). However, it remains controversial whether this performance implies reasoning capability beyond pattern matching.

A substantial body of research highlights the capability of Large Language Models in mathematical reasoning. Achiam et al. (2023); Team et al. (2024); Dubey et al. (2024); Yang et al. (2024), among others, achieved over 90% accuracy on GSM8K (Cobbe et al., 2021), a dataset consists of 8K grade school math word problems. Yang et al. (2024); Zhou et al. (2023), among others, achieved over 80% accuracy on the more difficult MATH dataset (Hendrycks et al., 2021), which consists of 12.5K high school math competition problems.

042

043

044

047

048

053

054

056

058

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

076

077

078

079

081

Meanwhile, there is a line of research questioning the reasoning ability of LLMs by showing their vulnerability under superficial changes of the input that do not alter the underlying logic. Works like Shi et al. (2023); Jiang et al. (2024) find that LLMs are easily distracted by irrelevant context or token level perturbation that does not change the underlying logic of the reasoning task. Mirzadeh et al. (2024) further demonstrate that the performance of LLMs declines when numerical values are altered in the questions from the GSM8K dataset.

There is yet another line of research that challenges the ability of LLMs to refrain from answering unanswerable problems. Ma et al. (2024); Li et al. (2024); Sun et al. (2024); Zhou et al. (2024a); Saadat et al. (2024) introduce minor modifications to existing math word problems to create unanswerable variants, and find that LLMs often generate hallucinatory answers for these unanswerable questions, even when they perform well on the original answerable datasets. However, these efforts rely on pre-existing math word problem sources, making them susceptible to training data contamination, limited in scope, and lacking rich structures for extended research.

To address these shortcomings, we propose TREECUT, a synthetic dataset capable of systematically generating an infinite number of unanswerable math word problems and their answerable counterparts. Our unanswerable dataset proves to be challenging even for GPT-40 and 03-mini. In addition, TreeCut allows precise control over the structural components of each problem, enabling detailed investigations into when and why LLMs produce hallucinations. We will release the dataset generation code upon publication.



Figure 1: The left and middle panels depict the tree structures corresponding to the answerable and unanswerable questions, respectively. In the right panel, the strike-through sentence represents the formula removed by the *cut*. The variable mappings to items are as follows: x_1 represents a burger, x_2 represents a scrambled egg, x_3 represents a BLT sandwich, and x_4 represents a pie.

2 Related Work

Math Word Problem Benchmark Numerous math word problem datasets of different difficulty have been proposed in previous research, most notable examples including GSM8K (Cobbe et al., 2021) and MATH (Hendrycks et al., 2021).

Many benchmarks have been developed to measure the robustness of mathematical reasoning. (Li et al., 2024; Zhou et al., 2024b; Yu et al., 2023; Shi et al., 2023) perturb or rewrite math word problems to measure the robustness of mathematical reasoning. Mirzadeh et al. (2024) developed GSM-Symbolic, a new benchmark derived from the GSM8K dataset by modifying numerical values, entity names, and question complexity.

Recently, Opedal et al. (2024) introduced Math-GAP, a framework for evaluating LLMs using synthetic math word problems with controllable proof tree characteristics. In contrast to their approach, the tree structure in our problem-generation procedure is fundamentally different. In our work, each node represents a variable, and the questioned variable appears as a leaf. In their work, however, each node represents a logical statement, with the answer represented by the root. More importantly, we focus on unanswerable math word problems, an aspect that their study did not address.

110Unanswerable Math ProblemsYin et al. (2023)111introduced SelfAware, consisting of unanswerable112questions from five diverse categories. It includes113less than 300 unanswerable mathematical problems.114Li et al. (2024) and Zhou et al. (2024a) generate115unanswerable questions by prompting GPT4 to116eliminate a necessary condition from the original117problem, an then the modified questions are fur-

ther checked or refined by human annotators. Sun et al. (2024) task human annotators to modify original questions in existing MWP datasets to make them unanswerable, creating a dataset composed of 2,600 answerable questions and 2,600 unanswerable questions. Ma et al. (2024) prompt GPT4 to modify problems from GSM8K, generating the Unreasonable Math Problem benchmark. 118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

151

3 TREECUT: a Synthetic (Un)answerable Math Word Problem Dataset

For the purpose of our investigation, we aim to have full control over the various aspects that determine the underlying *structure* of a math word problem: the name of the entities, the numeric values, and the complexity of the problem. Furthermore, we seek to reliably generate unanswerable problems by precisely removing specific necessary conditions of our choosing.

To this end, we start with a special kind of answerable math word problem that can be represented as a tree, as illustrated in Figure 1. Within such a tree, each non-root node represents a variable, while the root is a uniquely reserved node. An edge from root gives value to a variable, while an edge between two variables represents a linear formula of the two neighboring nodes. Given such a tree, any variable can be calculated following the unique path from the root to the node that represents the variable. Such a solving procedure does not require solving a linear equation system, as the solution only consists of carrying out basic arithmetic operations along the path. To guarantee that the arithmetic operations are well within the capacity of current frontier LLMs, we further restrict the

ansDepth	Llama-8B	Llama-70B	Qwen-7B	Qwen-72B	GPT-40	o3-mini
2	79%	22%	85%	61%	14%	42%
4	84%	35%	89%	86%	19%	23%
6	84%	63%	95%	89%	45%	18%
8	84%	63%	93%	84%	61%	23%

Table 1: Percentage of hallucination of various LLMs at different ansDepth values for unanswerable problems

ansDepth	Llama-8B	Llama-70B	Qwen-7B	Qwen-72B	GPT-40	o3-mini
2	68% (14%)	95% (1%)	87% (2%)	95% (1%)	99% (1%)	100% (0%)
4	28% (12%)	82% (6%)	31% (6%)	86% (6%)	94% (0%)	100% (0%)
6	17% (16%)	83% (3%)	12% (9%)	80% (7%)	85% (3%)	100% (0%)
8	5% (12%)	76% (7%)	7% (10%)	68% (8%)	84% (2%)	100% (0%)

Table 2: Accuracy of various LLMs at different ansDepth levels for answerable problems. The percentage in parentheses represents the proportion of answerable questions incorrectly identified as unanswerable.

unit price of each food item to be an integer between 5 and 15, and the coefficients of each linear equation taking non-zero integer values between -3 and 3. Finally, variables are randomly mapped to items, and then the formulas are translated to natural language using templates.

From an answerable math word problem described above, we generate an unanswerable problem by removing an edge along the path from the root to the questioned variable. In Figure 1, x_3 is the questioned variable. Along the path to the root, we remove the edge between x_1 and x_2 (denoted by a *cut*), rendering x_2 and x_3 undetermined, thus making the question unanswerable, as all we know about x_2 and x_3 is one single linear equation. A key benefit of such a generation procedure is that the distance from the questioned variable to the *cut* is also fully controlled, as we will see that this factor plays an important role in triggering LLM hallucination.

In summary, we can control the *structure* of problems via the following parameters:

- numVars: total number of variables,
- ansDepth: distance from the root to the questioned variable,
- compositeName: boolean, whether the items in the question have composite names (e.g. "a burger at Bistro Nice" versus "a burger"),
- cutDepth: distance from the questioned variable to the *cut*, if an unanswerable problem is to be generated.

Appendix A contains the detailed problem generation algorithm.

4 Experiments

We evaluate several state-of-the-art LLMs using TREECUT. Additionally, we analyze the halluci-

nation rate of GPT-40 on unanswerable problems generated under different parameter configurations of TREECUT.

190

191

192

193

194

195

196

197

199

200

201

202

203

204

205

206

207

208

209

210

211

212

213

214

215

216

217

218

219

220

4.1 Experimental Setup

For each set of generation parameters, we randomly generate 100 problems. During evaluation, we employ a zero-shot prompting template that explicitly directs the model to indicate when a question is unanswerable due to insufficient conditions. A chain-of-thought system message is incorporated for all models except o3-mini¹.

4.2 Evaluating LLMs

In the first set of experiments, we generate unanswerable math word problems of varying difficulty to evaluate the following LLMs: Llama 3.1 Instruct with 8B and 70B parameters(Dubey et al., 2024), Qwen2.5 Instruct with 7B and 72B parameters(Yang et al., 2024), GPT-4o(Achiam et al., 2023), and o3-mini(OpenAI, 2025).

Table 1 summarizes the results. None of the LLMs gives satisfactory results. Llama 3.1 8B, Qwen2.5 7B and 72B barely have any success identifying unanswerable problems. Llama 3.1 70B and GPT-40 struggle with more complex problems (ansDepth = 6, 8). o3-mini has the lowest hallucination for ansDepth = 6, 8. However, for the easiest case where ansDepth = 2 (in this setting, only 4 variables are mentioned in each problem), o3-mini displays a bias of making hallucinatory assumptions (see Appendix C.2 for examples).

To investigate whether the unsatisfactory accuracy of identifying unanswerable problems comes from the incapability of the necessary mathematical

176

177

178

179

181

182

183

184

187

152

¹Following OpenAI's guidelines of reasoning models.

operations, we evaluate the LLMs on the *answer-able* counterparts of the unanswerable questions
 using the same prompting template. We observe
 that almost every model displays a significant gap
 between its ability of solving answerable problems
 and identifying unanswerable problems. For in stance, GPT-40 correctly solves 84% of answerable
 problems for ansDepth = 8, but only correctly rec ognizes 39% of unanswerable problems.

4.3 Unanswerable Problem Structure and Hallucination

For a more fine-grained investigation of LLM's hallucination behavior under different *structures* of unanswerable problems, we analyze GPT-4o's hallucination rate on unanswerable problems generated under different parameter choices of numVars, ansDepth, compositeName and cutDepth.



Figure 2: Hallucination percentage under different configurations of unanswerable problems, plotted against varying ansDepth.

Tree Structure and Item Names To investigate the effect of (i) a deeper tree structure, (ii) a more complex tree structure, and (iii) composite item names, we consider the following parameter configurations:

- ansDepth ∈ {4, 5, 6, 7, 8}, which controls the depth of the questioned variable,
- cutDepth = $\lfloor ansDepth/2 \rfloor$
- numVars = ansDepth + 2 (generates a more complex tree structure) or numVars = ansDepth (the tree structure degenerates into a single path),
- compositeName: *true* or *false*.

There are $5 \times 2 \times 2 = 20$ configurations in total. We randomly generate 100 unanswerable problems for each configuration, and summarize GPT-4o's hallucination rate in Figure 2. In the figure,

* Orange line represents complex tree structure,

★ blue line represents simple tree structure,

256

257

258

261

262

263

264

265

266

267

268

269

270

272

273

274

275

276

277

278

279

280

281

282

286

287

290

291

293

295

- Solid line stands for composite item names,
- Dashed line stands for simple item names.

Examining each line individually, we observe that the hallucination rate increases as the depth of the questioned variable grows. Comparing solid and dashed lines of the same color, a more complex tree structure consistently results in a higher likelihood of hallucination across different ansDepth values. Comparing orange and blue lines of the same linestyle, composite item names consistently lead to a higher likelihood of hallucination compared to simple item names.



Figure 3: Hallucination percentage versus cutDepth. Left panel has ansDepth = 7. Right panel has ansDepth = 8.

Location of the *Cut* For each unanswerable problem, the *cut* always happens along the path from the root to the questioned variable. Does the location of the *cut* change hallucination ratio? We vary cutDepth from 1 to 7 while keeping ansDepth = 8 and other parameters fixed. In the right panel of Figure 3, we see that cutDepth = 3, 4, 5, 6 all trigger over 60% hallucination for GPT-40 (with cutDepth = 5 triggering over 70%), but a cutDepth = 1, 2, 7 only triggers less than 50% of hallucination, which means that GPT-40 is more confused when the *cut* happens around the middle point along the path, comparing to that happens near the root or the questioned variable.

4.4 Conclusion of Experiments

Our findings indicate that the unanswerable math word problems generated by TREECUT effectively induce hallucinations in large language models, including GPT-40 and 03-mini, with rates of 61% and 42% in their respective worst-case scenarios. Focusing on GPT-40, we further observe that hallucinations are more likely to occur when the problem exhibits (i) a deeper tree structure, (ii) a more complex tree structure, (iii) composite item names, or (iv) a *cut* positioned around the middle of the path. These results underscore the challenges LLMs face in handling unanswerable math problems.

255

238

234

350 351 352 353 354 355 357 360 361 362 363 364 365 366 367 368 369 370 371 372 373 374 375 376 378 379 381 383 384 386 389 390 391 392 393 394 397

398

399

400

401

402

403

348

349

5 Limitations

296

304

305

310

312

313

314

315

316

317

319

320

321

327

329

331

332

333

334

336

337

339

340

341

342

343

345

346

347

Our synthetic dataset is specifically designed for math word problems, representing only a small subset of the broader field of mathematics. Additionally, our evaluations are based solely on zero-shot chain-of-thought prompting. We do not explore alternative prompting techniques commonly used in LLM-based mathematical reasoning studies, which may impact performance comparisons.

References

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. 2021. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.
- Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. 2021. Measuring mathematical problem solving with the math dataset. *arXiv preprint arXiv:2103.03874*.
- Bowen Jiang, Yangxinyu Xie, Zhuoqun Hao, Xiaomeng Wang, Tanwi Mallick, Weijie J Su, Camillo J Taylor, and Dan Roth. 2024. A peek into token bias: Large language models are not yet genuine reasoners. *arXiv preprint arXiv:2406.11050*.
- Qintong Li, Leyang Cui, Xueliang Zhao, Lingpeng Kong, and Wei Bi. 2024. Gsm-plus: A comprehensive benchmark for evaluating the robustness of llms as mathematical problem solvers. *arXiv preprint arXiv:2402.19255*.
- Jingyuan Ma, Damai Dai, Lei Sha, and Zhifang Sui. 2024. Large language models are unconscious of unreasonability in math problems. *arXiv preprint arXiv:2403.19346*.
- Iman Mirzadeh, Keivan Alizadeh, Hooman Shahrokhi, Oncel Tuzel, Samy Bengio, and Mehrdad Farajtabar. 2024. Gsm-symbolic: Understanding the limitations of mathematical reasoning in large language models. *arXiv preprint arXiv:2410.05229*.
- Andreas Opedal, Haruki Shirakami, Bernhard Schölkopf, Abulhair Saparov, and Mrinmaya Sachan. 2024. Mathgap: Out-of-distribution evaluation on

problems with arbitrarily complex proofs. *arXiv* preprint arXiv:2410.13502.

- OpenAI. 2025. Openai o3 mini. Accessed: Feb. 5, 2025.
- Asir Saadat, Tasmia Binte Sogir, Md Taukir Azam Chowdhury, and Syem Aziz. 2024. When not to answer: Evaluating prompts on gpt models for effective abstention in unanswerable math word problems. *arXiv preprint arXiv:2410.13029*.
- Freda Shi, Xinyun Chen, Kanishka Misra, Nathan Scales, David Dohan, Ed Chi, Nathanael Schärli, and Denny Zhou. 2023. Large language models can be easily distracted by irrelevant context. In *Proceedings of the 40th International Conference on Machine Learning*, pages 31210–31227.
- YuHong Sun, Zhangyue Yin, Qipeng Guo, Jiawen Wu, Xipeng Qiu, and Hui Zhao. 2024. Benchmarking hallucination in large language models based on unanswerable math word problem. In *Proceedings of the* 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024), pages 2178–2188.
- Gemini Team, Petko Georgiev, Ving Ian Lei, Ryan Burnell, Libin Bai, Anmol Gulati, Garrett Tanzer, Damien Vincent, Zhufeng Pan, Shibo Wang, et al. 2024. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv preprint arXiv:2403.05530*.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837.
- An Yang, Beichen Zhang, Binyuan Hui, Bofei Gao, Bowen Yu, Chengpeng Li, Dayiheng Liu, Jianhong Tu, Jingren Zhou, Junyang Lin, et al. 2024. Qwen2.5math technical report: Toward mathematical expert model via self-improvement. *arXiv preprint arXiv:2409.12122.*
- Zhangyue Yin, Qiushi Sun, Qipeng Guo, Jiawen Wu, Xipeng Qiu, and Xuan-Jing Huang. 2023. Do large language models know what they don't know? In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 8653–8665.
- Longhui Yu, Weisen Jiang, Han Shi, Jincheng Yu, Zhengying Liu, Yu Zhang, James T Kwok, Zhenguo Li, Adrian Weller, and Weiyang Liu. 2023. Metamath: Bootstrap your own mathematical questions for large language models. *arXiv preprint arXiv:2309.12284*.
- Aojun Zhou, Ke Wang, Zimu Lu, Weikang Shi, Sichun Luo, Zipeng Qin, Shaoqing Lu, Anya Jia, Linqi Song, Mingjie Zhan, et al. 2023. Solving challenging math word problems using gpt-4 code interpreter with code-based self-verification. *arXiv preprint arXiv:2308.07921*.

404	Zihao Zhou, Shudong Liu, Maizhen Ning, Wei Liu,
405	Jindong Wang, Derek F Wong, Xiaowei Huang, Qi-
406	ufeng Wang, and Kaizhu Huang. 2024a. Is your
407	model really a good math reasoner? evaluating math-
408	ematical reasoning with checklist. arXiv preprint
409	arXiv:2407.08733.

Zihao Zhou, Qiufeng Wang, Mingyu Jin, Jie Yao, Jianan
Ye, Wei Liu, Wei Wang, Xiaowei Huang, and Kaizhu
Huang. 2024b. Mathattack: Attacking large language
models towards math solving ability. In *Proceedings*of the AAAI Conference on Artificial Intelligence,
pages 19750–19758.

A Generation of the Math Word Problems

```
Algorithm 1 Generating Math Word Problem using Random Tree
Require: numVars \geq ansDepth \geq 2
Require: unanswerable \in \{true, false\}, order \in \{"forward", "backward", "random"\}
Require: cutDepth: int
 1: if unanswerable = true then
Require: cutDepth: int, satisfying 1 \le \text{cutDepth} < \text{ansDepth}
 2: end if
                                                                 \triangleright (i) Sample a dictionary of variable values
 3: varDict \leftarrow {}
 4: for i \leftarrow 1 to numVars do
        Sample an integer v \in [5, 15]
 5:
        varDict[x_i] \leftarrow v
 6:
 7: end for
                                                                                  \triangleright (ii) Build the random tree
 8: Assign root as the parent of x_1
 9: for i \leftarrow 2 to ansDepth do
        Assign x_{i-1} as the parent of x_i
10:
11: end for
                                         ▷ Finish building the path from the root to the questioned variable
                                                                                ▷ Assign the remaining nodes
12: for i \leftarrow ansDepth + 1 to numVars do
        Randomly select a node x_p in the tree
13:
14:
        Assign x_p as the parent of x_i
15: end for
                                                  ▷ (iii) Get the list of all edges via a breadth-first traversal
16: edgeList \leftarrow the list of edges collected by a breadth-first traversal (see Algorithm 2)
                                                            \triangleright (iv) For unanswerable problems, create the cut
17: if unanswerable = true then
18:
        Remove (x_{ansDepth-cutDepth-1}, x_{ansDepth-cutDepth}) from edgeList
19: end if
                                        \triangleright (v) Generate a formula for each edge, and store in forumlaList
20: formulaList \leftarrow []
21: for edge (x_i, x_j) in edgeList do
        Sample a, b \in \{-3, -2, -1, 1, 2, 3\}
22:
        Define formula \leftarrow a \cdot x_i + b \cdot x_j = a \cdot \text{varDict}[x_i] + b \cdot \text{varDict}[x_j]
23:
24:
        Append formula to formulaList
                                                   ▷ So that formulaList has the same order as edgeList
25: end for
                                            ▷ (vi) Adjust the ordering of formulaList according to order
26: if order = "backward" then
        Reverse formulaList
27:
28: end if
29: if order = "random" then
        Random Shuffle formulaList
30:
31: end if
          return formulaList
                                                           ▷ Formulas serving as conditions of the problem.
```

Algorithm 1 generates formulaList, which contains the formulas that will serve as the conditions of the problem. To translate that into natural language, item names will be sampled according to the *compositeName* option. Then, formulaList can be translated to natural language using pre-defined templates. The question sentence will simply be "what is the price of {item name of the questioned

417

418

419

variable}".

We want to point out that although all the variables are assigned a value in varDict, this is purely for the sake of (i) subsequently generating the random formulas (ii) guaranteeing that all calculable variables will have values between 5 and 15. When unanswerable = true, the *cut* will guarantee that the problem is unanswerable.

In the following, we also detail the simple breadth-first traversal algorithm for getting all the edges from the tree, which enables us to control the order of the conditions in the problem.

Require: root: the root of a tree

▷ Get the list of all edges via a breadth-first traversal

```
1: edgeList \leftarrow [], q \leftarrow a queue containing root
```

```
2: while q is not empty do
```

- 3: node \leftarrow q.dequeue()
- 4: **for** child \in node.children **do**
- 5: Add (node, child) to edgeList
- 6: Add child to q
- 7: end for
- 8: end while

return edgeList

427

428

429

430

431

432

421 422

423

424 425

426

B Details of Experiments

B.1 Prompt Template

Below is the prompt template used for evaluation, which is a 0-shot prompting template with chain-ofthought system message. For reasoning models like o3-mini, we do not use the chain-of-thought system message.

Evaluation Prompt Template

```
// system message (for non-reasoning models)
As an expert problem solver, solve step by step the following mathematical
questions.
// user message
Please solve the following math question, and then answer in the form 'Answer:
x'. If the known conditions are not sufficient to answer the question, please
answer in the form 'Answer: unknown.'.
Question: {{question}}
```

Your solution:

B.2 Evaluation

For all the models other than o3-mini, we set {max_tokens=4000, temperature=0}. o3-mini does not allow setting temperature, and we set {max_completion_tokens=32000, reasoning_effort="high"}. GPT-40 and o3-mini are evaluated via calling OpenAI API.

When evaluating model's answer to an unanswerable problem, we change the response to all lower case, and if "unknown" is mentioned after the last occurrence of the trigger word "answer", it is counted as a correct answer.

433

434 435

436

437

438

439

C Examples of Hallucinatory Answers	441
C.1 GPT-40	442
Example 1: A BLT sandwich at Urban Plate costs 13 dollars. 3 burgers at Taste Good Cuisine cost 6 dollars more than 3 fruit tarts at Taste Good Cuisine. A BLT sandwich at Taste Good Cuisine costs 5	443 444
dollars A fruit tart at Taste Good Cuisine costs 9 dollars less than 3 burgers at Urban Plate 2 burgers at	
Taste Good Cuisine cost 8 dollars more than a fruit tart at Urban Plate. Question: how much does a burger	446
at Taste Good Cuisine cost?	447
Model answer:	448
To solve this problem, we need to set up equations based on the information given and solve for the unknowns.	449
Latio define the verichles.	450
- Let \(B T \) be the cost of a burger at Taste Good Cuisine.	451
- Let \(F_T \) be the cost of a fruit tart at Taste Good Cuisine.	453
- Let (B_U) be the cost of a burger at Urban Plate.	454
- Let ((F_U)) be the cost of a fruit tart at urban Plate.	455
We are given the following information:	457
	458
1. A BLT sandwich at Urban Plate costs 13 dollars.	459 460
	461
$3B_T = 3F_T + 6$	462
<pre>\]</pre>	463
4. A fruit tart at Taste Good Cuisine costs 9 dollars less than 3 burgers at Urban Plate.	464
ν[466
F_T = 3B_U - 9	467
\] 5 2 burgers at Taste Good Cuisine cost 8 dollars more than a fruit tart at Urban Plate	468
	470
$2B_T = F_U + 8$	471
	472
We need to find the cost of a burger at Taste Good Cuisine, (B_T) .	474
	475
From equation (2), we can express \(F_I \) in terms of \(B_I \):	476
$F_T = B_T - 2$	478
\]	479
Substitute this expression for (E,T,V) into equation (A):	480
	482
$B_T - 2 = 3B_U - 9$	483
	484
$B_T = 3B_U - 7$	485
	487
Now late substitute $\langle (E,T, -B,T, -2) \rangle$ into equation (E).	488
$\langle [$	409
$2B_T = (B_T - 2) + 8$	491
	492
2B T = B T + 6	493
	495
	496
\]	497 498
	499
Thus, the cost of a burger at Taste Good Cuisine is 6 dollars.	500 501
Answer: 6	502
Mistake: When the model did "Now, let's substitute $F_T = B_T - 2$ into equation (5)", it	503
actually confused F_T with F_U in equation (5).	504

Example 2: 3 Greek salads at Texas BBQ cost 6 dollars less than 3 pieces of cheese cake at Texas BBQ.
3 burgers at Texas BBQ and 2 pieces of cheese cake at Texas BBQ cost 59 dollars. 3 burgers at Taste
Good Cuisine and 2 Greek salads at Taste Good Cuisine cost 42 dollars. 2 burgers at Texas BBQ and a
piece of cheese cake at Taste Good Cuisine cost 32 dollars. 3 burgers at Texas BBQ cost 15 dollars more
than 3 fruit tarts at Taste Good Cuisine. A burger at Taste Good Cuisine costs 10 dollars. A fruit tart at
Texas BBQ costs 2 dollars more than a Greek salad at Texas BBQ. Question: how much does a Greek
salad at Texas BBQ cost?

512 Model answer:

513

514

515 516

517

518 519

520

521

522

524 525

526

527

530 531

534

536

537

538

539

540

541

542

543

545

547

548

549

550

553 554

556

557

558

559

560 561

564 565

567

568

569 570

571

To solve this problem, we need to set up equations based on the information given and solve for the cost of a Greek salad at Texas BBQ. Let's define the variables:

```
- \( G \): Cost of a Greek salad at Texas BBQ
- \( C \): Cost of a piece of cheesecake at Texas BBQ
- \( B \): Cost of a burger at Texas BBQ
- \( F \): Cost of a fruit tart at Taste Good Cuisine
- \( B_T \): Cost of a burger at Taste Good Cuisine (given as $10)
- \( F_T \): Cost of a fruit tart at Texas BBQ
Now, let's translate the information into equations:
1. "3 Greek salads at Texas BBQ cost 6 dollars less than 3 pieces of cheesecake at Texas BBQ."
   ١Ľ
   3G = 3C - 6
   \]
2. "3 burgers at Texas BBQ and 2 pieces of cheesecake at Texas BBQ cost 59 dollars."
   ١Ľ
   3B + 2C = 59
   1
3. "3 burgers at Taste Good Cuisine and 2 Greek salads at Taste Good Cuisine cost 42 dollars."
   ١L
   3B_T + 2G = 42
   1
   Since (B_T = 10), substitute:
   ١L
   3(10) + 2G = 42 \implies 30 + 2G = 42 \implies 2G = 12 \implies G = 6
   \backslash ]
4. "2 burgers at Texas BBQ and a piece of cheesecake at Taste Good Cuisine cost 32 dollars."
   ١L
   2B + C_T = 32
   \1
  (Note: We don't have enough information about (C_T ), the cost of cheesecake at Taste Good Cuisine,
   so we can't use this equation directly.)
5. "3 burgers at Texas BBQ cost 15 dollars more than 3 fruit tarts at Taste Good Cuisine."
   ١Ľ
   3B = 3F + 15
   \٦
6. "A fruit tart at Texas BBQ costs 2 dollars more than a Greek salad at Texas BBQ."
   ١Ľ
   F_T = G + 2
   \]
Now, let's solve the equations:
From equation 3, we found \langle G = 6 \rangle.
Let's verify if this value is consistent with other equations:
- From equation 1:
  ١L
  3G = 3C - 6 \implies 3(6) = 3C - 6 \implies 18 = 3C - 6 \implies 3C = 24 \implies C = 8
  \backslash ]
```

```
- From equation 2:
                                                                                                                        572
                                                                                                                        573
  ١Ľ
  3B + 2C = 59 \implies 3B + 2(8) = 59 \implies 3B + 16 = 59 \implies 3B = 43 \implies B = \frac{43}{3}
                                                                                                                        574
  \]
                                                                                                                        575
  (This value seems unusual, but let's proceed.)
                                                                                                                        576
                                                                                                                        577
- From equation 5:
                                                                                                                        578
  ١L
                                                                                                                        579
 3B = 3F + 15 \implies \frac{43}{3} = F + 5 \implies F = \frac{43}{3} - 5 = \frac{43}{3} - \frac{15}{3} = \frac{28}{3}680
  \]
                                                                                                                        582
                                                                                                                        583
- From equation 6:
                                                                                                                        584
  ١Ľ
  F_T = G + 2 \setminus implies F_T = 6 + 2 = 8
                                                                                                                        585
                                                                                                                        586
                                                                                                                        587
```

The calculations are consistent with the given conditions. Therefore, the cost of a Greek salad at Texas BBQ is:

Answer: 6

Mistake: At step 3, model calculated the value of "G", but from the start of the solution, "G" denotes the "cost of a Greek salad at Texas BBQ", not the cost of a Greek salad at Taste Good Cuisine.

C.2 o3-mini

Example 1: A pizza at Taste Good Cuisine and 3 lasagnas at Taste Good Cuisine cost 48 dollars. A lasagna at Urban Plate costs 6 dollars. 3 lasagnas at Urban Plate and 3 pizzas at Urban Plate cost 63 dollars. Question: how much does a lasagna at Taste Good Cuisine cost?

Model answer:

Let p be the price of a pizza at Taste Good Cuisine and 1 be the price of a lasagna at Taste Good Cuisine. We are given:

p + 31 = 48. (Taste Good Cuisine)
 A lasagna at Urban Plate costs 6 dollars.
 3 lasagnas at Urban Plate and 3 pizzas at Urban Plate cost 63 dollars.

From (3), the total cost of 3 lasagnas at Urban Plate is $3 \times 6 = 18$ dollars. This means that the 3 pizzas at Urban Plate cost 63 - 18 = 45 dollars, so one pizza at Urban Plate costs 45/3 = 15 dollars.

A common (though unstated) assumption in these puzzles is that dishes with the same name cost the same across restaurants. Thus, we assume that a pizza at Taste Good Cuisine also costs 15 dollars. Substitute p = 15 into equation (1):

15 + 31 = 48 31 = 48 - 15 = 331 = 33/3 = 11

Answer: 11

Mistake: The model made an assumption: "A common (though unstated) assumption in these puzzles is that dishes with the same name cost the same across restaurants." However, this cannot be a mathematically valid assumption because following this assumption, it arrives at the conclusion that the price for lasagna at Urban Plate is not the same as that at Taste Good Cuisine (6 dollars versus 11 dollars).

We randomly selected 20 problems from the ansDepth = 2 set that triggered hallucinatory answer from o3-mini and manually checked the model answer, all of them shares this same mistake. Below is a similar example:

Example 2: A scrambled egg at Bistro Nice and a piece of cheese cake at Mike's Place cost 21 dollars. 2 scrambled eggs at Bistro Nice cost 6 dollars more than 2 pieces of cheese cake at Bistro Nice. A scrambled egg at Mike's Place costs 6 dollars. Question: how much does a scrambled egg at Bistro Nice cost? **Model answer:**

Let S_B = price of a scrambled egg at Bistro Nice, C_B = price of a piece of cheesecake at Bistro Nice, C_M = price of a piece of cheesecake at Mike's Place.

The problem gives:
1. S_B + C_M = 21 dollars.
2. Two scrambled eggs at Bistro Nice cost 6 dollars more than two cheesecakes at Bistro Nice: 2S_B = 2C_B + 6 → S_B = C_B + 3.
3. A scrambled egg at Mike's Place costs 6 dollars (this is extra information).

Without further assumptions the system involves three unknowns (S_B, C_B, C_M) and only two equations. However, it is common in these problems to assume that when the same item is mentioned (here "cheesecake") it has a consistent price from one restaurant to the other. That is, we assume $C_M = C_B$.

Under that assumption the first equation becomes: $S_B + C_B = 21$. Then using $S_B = C_B + 3$ we substitute: $(C_B + 3) + C_B = 21 \rightarrow 2C_B + 3 = 21 \rightarrow 2C_B = 18 \rightarrow C_B = 9$. Thus, $S_B = 9 + 3 = 12$.

Answer: 12

631

633

634

641

642

644

645

647

651

657

659

660 661

664

669

673 674

675

677

678

679 680

681

Mistake: Same as above.

Example 3: 3 BLT sandwiches at Bistro Nice cost 6 dollars more than 3 pieces of cheese cake at Urban Plate. 3 burgers at Urban Plate cost 18 dollars less than 3 pieces of cheese cake at Bistro Nice. 2 BLT sandwiches at Urban Plate and 2 scrambled eggs at Urban Plate cost 34 dollars. A BLT sandwich at Urban Plate costs 8 dollars. 3 pieces of cheese cake at Bistro Nice cost 12 dollars more than 3 burgers at Bistro Nice. 2 burgers at Urban Plate cost 10 dollars less than 2 scrambled eggs at Bistro Nice. 3 pieces of cheese cake at Bistro Nice. 3 pieces of cheese cake at Bistro Nice and a piece of cheese cake at Urban Plate cost 54 dollars. Question: how much does a piece of cheese cake at Urban Plate cost?

Model answer:

Let CC_U be the cost of one piece of cheese cake at Urban Plate.

Step 1. From "3 BLT sandwiches at Bistro Nice cost 6 dollars more than 3 pieces of cheese cake at Urban Plate": $3 \cdot BLT_B = 3 \cdot CC_U + 6 \rightarrow BLT_B = CC_U + 2$ (1)

Step 2. From "3 burgers at Urban Plate cost 18 dollars less than 3 pieces of cheese cake at Bistro Nice": $3 \cdot BUR_U = 3 \cdot CC_B - 18 \rightarrow With BUR_U$ determined later, we use this after finding BUR_U.

Step 3. "2 BLT sandwiches at Urban Plate and 2 scrambled eggs at Urban Plate cost 34 dollars": $2 \cdot BLT_U + 2 \cdot Eggs_U = 34$ and we are given in Step 4 that $BLT_U = 8$. Thus, $2 \cdot 8 + 2 \cdot Eggs_U = 34 \rightarrow 16 + 2 \cdot Eggs_U = 34 \rightarrow 2 \cdot Eggs_U = 18 \rightarrow Eggs_U = 9$.

Step 4. "A BLT sandwich at Urban Plate costs 8 dollars": BLT_U = 8.

Step 5. "3 pieces of cheese cake at Bistro Nice cost 12 dollars more than 3 burgers at Bistro Nice": $3 \cdot CC_B = 3 \cdot BUR_B + 12 \rightarrow CC_B = BUR_B + 4$ (2)

Step 6. "2 burgers at Urban Plate cost 10 dollars less than 2 scrambled eggs at Urban Plate": $2 \cdot BUR_U = 2 \cdot Eggs_U - 10 \rightarrow 2 \cdot BUR_U = 18 - 10 \rightarrow 2 \cdot BUR_U = 8 \rightarrow BUR_U = 4.$

Step 7. Now return to the burgers at Urban Plate in Step 2: $3 \cdot BUR_U = 3 \cdot CC_B - 18$ With $BUR_U = 4$, $3 \cdot 4 = 3 \cdot CC_B - 18 \rightarrow 12 = 3 \cdot CC_B - 18 \rightarrow 3 \cdot CC_B = 30 \rightarrow CC_B = 10$.

Step 8. "3 pieces of cheese cake at Bistro Nice and a piece of cheese cake at Urban Plate cost 54 dollars": $3 \cdot CC_B + CC_U = 54 \rightarrow 3 \cdot 10 + CC_U = 54 \rightarrow 30 + CC_U = 54 \rightarrow CC_U = 24$.

The question asks: "how much does a piece of cheese cake at Urban Plate cost?"

Answer: 24

Mistake: Step 6. "2 burgers at Urban Plate cost 10 dollars less than 2 scrambled eggs at Urban Plate". That is a hallucinatory condition. The actual condition is "2 burgers at Urban Plate cost 10 dollars less than 2 scrambled eggs at Bistro Nice".