
CausalBench: A Comprehensive Benchmark for Evaluating Causal Reasoning Capabilities of Large Language Models

Abstract

1 Causal reasoning, a core aspect of human cognition, is essential for advancing
2 large language models (LLMs) towards artificial general intelligence (AGI) and
3 reducing their propensity for generating hallucinations. However, existing datasets
4 for evaluating causal reasoning in LLMs are limited by narrow domain coverage
5 and a focus on cause-to-effect reasoning through textual problems, which does
6 not comprehensively assess whether LLMs truly grasp causal relationships or
7 merely guess correct answers. To address these shortcomings, we introduce a
8 novel benchmark that spans textual, mathematical, and coding problem domains.
9 Each problem is crafted to probe causal understanding from four perspectives:
10 cause-to-effect, effect-to-cause, cause-to-effect with intervention, and effect-to-
11 cause with intervention. This multi-dimensional evaluation method ensures that
12 LLMs must exhibit a genuine understanding of causal structures by correctly
13 answering questions across all four dimensions, mitigating the possibility of cor-
14 rect responses by chance. Furthermore, our benchmark explores the relationship
15 between an LLM’s causal reasoning performance and its tendency to produce
16 hallucinations. We present evaluations of state-of-the-art LLMs using our bench-
17 mark, providing valuable insights into their current causal reasoning capabili-
18 ties across diverse domains. The dataset is publicly available for download at
19 <https://huggingface.co/datasets/CCLV/CausalBench>.

20 1 Introduction

21 Causal reasoning, the ability to understand and infer causal relationships between variables, is a
22 fundamental aspect of human cognition and plays a crucial role in decision-making, problem-solving,
23 and learning [1]. For large language models (LLMs), causal reasoning refers to the ability to
24 accurately identify, represent, and reason about causal relationships described in text, mathematical
25 equations, or code snippets [1]. Developing strong causal reasoning abilities in LLMs is essential
26 for progress toward artificial general intelligence (AGI), as it enables models to understand not just
27 correlations but the underlying mechanisms driving outcomes [3]. This understanding is crucial for
28 making accurate predictions, generating insightful explanations, and adapting to new situations, as
29 core components of AGI.

30 However, existing causal reasoning benchmarks have several limitations that hinder their ability to
31 comprehensively evaluate the causal reasoning capabilities of LLMs. First, current benchmarks often
32 focus on a single perspective of causal reasoning, such as cause-to-effect, lacking a multifaceted
33 assessment that considers effect-to-cause reasoning and the impact of interventions. This narrow
34 focus allows models to correctly answer causal questions by chance without truly understanding the
35 underlying causal relationships [5]. Second, current benchmarks are primarily text-based, lacking
36 diversity in problem types, such as mathematical and coding problems that can encapsulate causal
37 dependencies. Incorporating these diverse problem formats would enable a more robust evaluation

38 of LLMs’ capacity to reason about causality across various modalities. Third, the limited scale of
39 existing benchmarks may not provide a sufficiently comprehensive assessment of LLMs’ causal
40 reasoning abilities due to the limited scale of the benchmark dataset.

41 To address these limitations, we propose CausalBench, a comprehensive benchmark for evaluating the
42 causal reasoning capabilities of LLMs. CausalBench comprises four perspectives of causal reasoning
43 for each scenario: cause-to-effect, effect-to-cause, cause-to-effect with intervention, and effect-to-
44 cause with intervention. This multi-perspective approach mitigates the potential for correct answers
45 by chance and provides a more accurate evaluation of LLMs’ understanding of causal relationships.
46 Moreover, CausalBench includes a diverse set of problem types spanning textual, mathematical, and
47 coding domains, enabling a comprehensive assessment of causal reasoning abilities across different
48 modalities. The benchmark consists of more than 60,000 problems and employs six evaluation
49 metrics to measure LLMs’ causal reasoning performance.

50 The major contributions of CausalBench are three-fold: (1) evaluating four causal reasoning perspec-
51 tives per scenario to robustly assess causal understanding, (2) incorporating a diverse problem set
52 spanning math, code, and natural language for cross-modal evaluation, and (3) implementing strict
53 quality control measures, including a causal inference engine check and human expert review, to
54 ensure the benchmark’s validity and reliability. By addressing the limitations of existing benchmarks,
55 CausalBench aims to provide a more comprehensive and accurate evaluation of the causal reasoning
56 capabilities of LLMs, facilitating progress towards AGI.

57 **2 Related Works**

58 Existing datasets and benchmarks for evaluating causal reasoning primarily focus on commonsense
59 causality [9, 31, 32], which assesses the alignment between commonsense knowledge about causal
60 relationships in humans and language models. These datasets, such as WikiWhy [9], CausalWorld
61 [31], and UCLM [32], provide valuable insights into how well language models capture and reason
62 about everyday causal relationships. However, they do not explicitly evaluate the models’ ability to
63 perform formal causal reasoning based on well-defined rules and principles from the field of causal
64 inference. Some recent works have started to explore more formal aspects of causal reasoning in
65 language models. For example, CRASS [28] focuses specifically on counterfactual reasoning, which
66 involves reasoning about alternative outcomes based on hypothetical changes to past events. While
67 counterfactual reasoning is an important aspect of causal inference, CRASS does not cover the full
68 spectrum of causal inference tasks, such as interventional and observational reasoning. Another
69 concurrent work by Kiciman et al. [16] evaluates language models on various causality-related tasks,
70 including causal sufficiency analysis, causal discovery, and counterfactual reasoning. However, their
71 evaluation primarily relies on the conceptual knowledge accrued from the training data rather than
72 formal causal inference, except for their causal sufficiency analysis. This means that the models’
73 performance may be influenced by spurious correlations or memorization from the training data
74 rather than a genuine understanding of causal principles.

75 In contrast, our proposed dataset, CausalBench, is grounded in the principles of causal inference
76 [11, 25, 26]. CausalBench provides a comprehensive and principled framework for assessing the
77 causal reasoning capabilities of language models, ensuring that the models are evaluated on their
78 ability to perform formal causal inference rather than relying on spurious correlations or memorization
79 from training data. By encompassing a diverse set of causal scenarios (text, code, and math),
80 four causal perspectives (cause to effect, effect to cause, cause to effect with intervention, and
81 effect to cause with intervention), and explanations associated with ground truth for each test case,
82 CausalBench offers a rigorous and systematic approach to benchmarking causal reasoning in LLMs.
83 It is designed to test the models’ ability to reason about causal relationships in a variety of domains,
84 including natural language, programming code, and mathematical equations. In summary, while
85 existing datasets and benchmarks have made contributions to the study of causal reasoning in language
86 models, CausalBench offers a more comprehensive, principled, and rigorous approach to evaluating
87 formal causal inference capabilities across multiple domains. By grounding the evaluation in the
88 principles of causal inference and providing a diverse set of test cases with associated explanations,
89 CausalBench aims to set a new standard for benchmarking causal reasoning in LLMs.

90 3 Dataset Construction Process and Method

91 The construction of CausalBench involves three key steps: manual generation of initial test cases,
92 scaling up using LLM such as GPT-4 Turbo, and quality control through causal inference engines
93 together with human verification. Initially, we manually create a set of test cases covering four aspects
94 of causal inference: (a) cause to effect, (b) effect to cause, (c) cause to effect with intervention, and (d)
95 effect to cause with intervention to ensure a comprehensive evaluation of causal reasoning capabilities
96 from different perspective. To expand the dataset, we then use GPT-4 Turbo with few-shot prompting,
97 leveraging the model’s ability to generate additional test cases that adhere to the desired format and
98 cover the four causal inference aspects. The few-shot prompts are designed to guide GPT-4 Turbo
99 in producing a diverse and extensive set of problems that maintain consistency with the manually
100 generated cases. Afterward, we implement a quality control process involving validation through
101 causal inference engines and review by human experts. The causal inference engines verify the
102 logical consistency and correctness of the generated test cases, while human experts review and refine
103 the dataset to maintain high standards of quality and relevance.

104 3.1 Workflow Overview

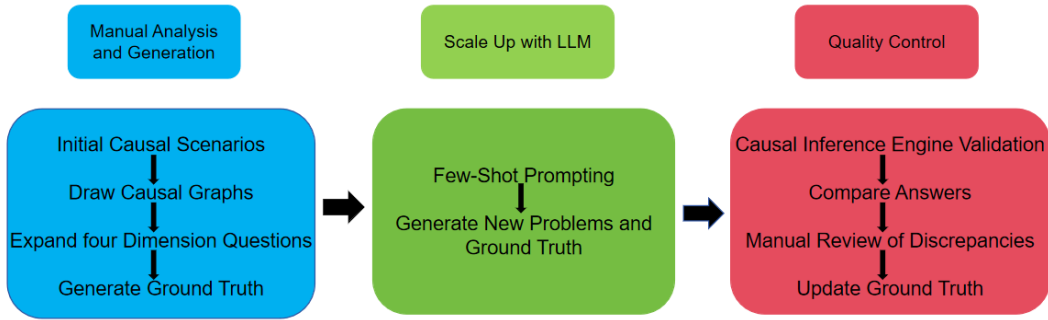


Figure 1: Workflow overview of the CausalBench dataset construction process.

105 3.2 Manual Analysis and Generation

106 For the text problems of our Benchmark, we randomly selected 100 questions from the CLADDER
107 dataset [10] and manually analyzed them to determine their category within (1) inference from
108 cause to effect, (2) effect to cause, (3) cause to effect with intervention, or (4) effect to cause with
109 intervention. These perspectives represent different dimensions of causal reasoning: (1) Cause to the
110 effect: Given the cause, what is the likelihood of the effect? (2) Effect to cause: Given the effect,
111 what is the likelihood of the cause? (3) Cause to effect with intervention: If an intervention is added
112 to the causal relationship, given the cause, what is the likelihood of the effect? and (4) Effect to cause
113 with intervention: If an intervention is added to the causal relationship, given the effect, what is the
114 likelihood of the cause?

115 After categorizing the selected cases from the CLADDER dataset, we expanded them by creating
116 additional questions for the other three perspectives. For example, if a case was classified as “cause to
117 effect”, we generated corresponding questions for “effect to cause”, “cause to effect with intervention”,
118 and “effect to cause with intervention” manually.

119 To correctly expand other perspective questions and their ground truths, we visualized the relationships
120 between variables using causal diagrams and analyzed these relationships by calculating conditional
121 probabilities. Causal diagrams represent variables as nodes and causal relationships as directed edges.
122 For example, consider the following hypothetical scenario:

123 *Imagine a self-contained, hypothetical world with only the following conditions, and without any*
124 *unmentioned factors or causal relationships: Parents’ intelligence has a direct positive effect on*
125 *parents’ social status and child’s intelligence. Other unobserved factors has a positive direct effect*
126 *on parents’ social status and child’s intelligence. If a child is intelligent, would it be more likely that*
127 *this child had intelligent parents?*

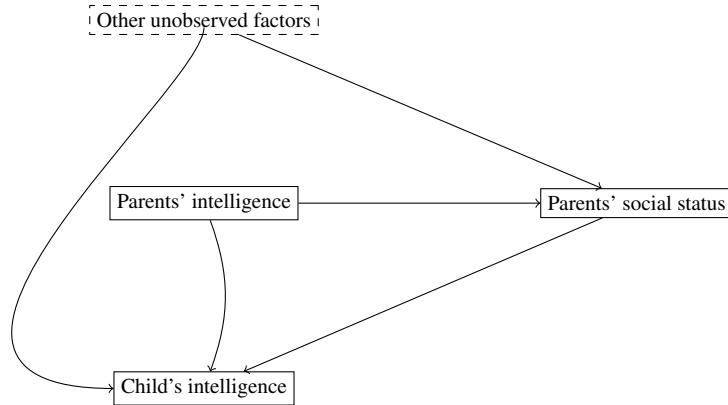


Figure 2: Causal Graph Example

128 In this scenario, the causal diagram would have four nodes: Parents' intelligence, Parents' social
 129 status, Child's intelligence, and Other unobserved factors. There would be directed edges from
 130 Parents' intelligence to Parents' social status and Child's intelligence, from Other unobserved factors
 131 to Parents' social status and Child's intelligence, and from Parents' social status to Child's intelligence.
 132 Conditional probabilities can be estimated based on the causal graph.

133 Using the causal graph and conditional probabilities, we can categorized the original questions as
 134 effect-to-cause. The probability of the child being intelligent given that the parents are intelligent is
 135 higher than the probability of the child being intelligent given that the parents are unintelligent, so the
 136 ground truth is yes. Then extend the questions to cover four perspectives by adjusting the questioning
 137 logic and incorporating interventions into the causal path diagram, and calculate ground truth for
 138 each questions.(examples are provided in the Appendix)

139 Finally, we obtained 100 causal scenarios, with 400 causal questions. They serve as the foundation
 140 for our few-shot prompting approach, providing examples for GPT-4 Turbo on how to identify the
 141 type of the initial question and generate additional questions for the remaining perspectives. By
 142 using these examples in a few-shot prompting setting, we guide the model to generate additional
 143 perspective questions with answers for all other causal scenarios in the CLADDER dataset.

144 For coding and mathematical problems, we manually created 100 code scenarios and 100 math
 145 scenarios, each containing causal relationships, and designed four perspective questions for each
 146 scenario. These questions addressed causal issues based on the relationships described in the scenarios
 147 (examples are provided in the Appendix). We then used causal graphs and conditional probabilities to
 148 manually generate the ground truths and employed few-shot prompts with GPT-4 Turbo to generate
 149 additional code, math scenarios and questions with corresponding answers.

150 In summary, the manual analysis and generation process involved visualizing causal relationships
 151 using causal diagrams and calculating conditional probabilities for each scenario. We modified the
 152 questioning approach and added interventions to expand each problem into four forms, covering cause-
 153 to-effect, effect-to-cause, cause-to-effect with intervention, and effect-to-cause with intervention, and
 154 generated ground truths for each question. By the end of this section, we had created 100 sets of 400
 155 text-based questions with ground truths, 100 sets of 400 coding questions with ground truths, and
 156 100 sets of 400 math questions with ground truths. These manually generated samples serve as the
 157 foundation for our few-shot prompting approach, which utilizes GPT-4 Turbo to generate additional
 158 test cases.

159 3.3 Scaling Up with LLMs

160 After manually generating and verifying an initial set of questions, we employed GPT-4 Turbo to scale
 161 up the dataset. The scale-up process was divided into three parts: text problems, coding problems,
 162 and mathematical problems.

163 For the text problems, we provided GPT-4 Turbo with original CLADDER dataset[10] questions with
 164 manually expanded questions along with their ground truths. By learning from these samples, GPT-4

165 Turbo was tasked with reading the remaining CLADDER scenarios (around 10,000 problems) and
166 their corresponding questions, determining the question perspective, expanding the scenario into the
167 other three perspectives, and generating the associated ground truths. This process ensures every text
168 causal scenario has four dimension questions and corresponding ground truths.

169 In the case of coding problems, we supplied GPT-4 Turbo with the 100 manually created code
170 examples containing causal relationships. Using these examples as a foundation, GPT-4 Turbo
171 generated an additional 2,000 code snippets, each incorporating causal relationships. For each
172 newly generated code snippet, GPT-4 Turbo created four perspectives of questions and provided the
173 corresponding ground truths, ensuring a comprehensive evaluation of causal reasoning in the context
174 of programming.

175 Similarly, for mathematical problems, GPT-4 Turbo was employed to generate 2,000 new mathemati-
176 cal scenarios across various domains, such as probability theory, mathematical statistics, differential
177 equations, and complex analysis. For each mathematical scenario, GPT-4 Turbo generated four types
178 of questions and their associated ground truths, assessing the model’s ability to reason about causal
179 relationships in mathematical contexts.

180 By leveraging the capabilities of GPT-4 Turbo, we were able to create a dataset across all three
181 problem categories. The text problems were augmented by automatically generating additional
182 question perspectives and ground truths based on the existing CLADDER scenarios. The coding
183 and mathematical problems were scaled up by having GPT-4 Turbo create new scenarios containing
184 causal relationships and generate the corresponding questions and ground truths. This scale-up
185 process resulted in a more comprehensive and diverse dataset, enabling a thorough evaluation of
186 causal reasoning abilities in large language models across various domains.

187 **3.4 Quality Control**

188 **3.4.1 Causal Inference Engine Design**

189 To ensure the accuracy and consistency of the generated questions and answers, we developed a
190 causal inference engine. This engine utilizes causal diagrams and conditional probabilities associated
191 with each question to compute the answers for all questions. The causal inference engine serves as a
192 verification layer, comparing the answers generated by the language model. If the answer generated
193 by the language model differs from the answer generated by the causal inference engine, the case
194 will be manually inspected, and the ground truth will be generated by human experts. Here are the
195 Causal Inference Engine design details:

196 **Input**

- 197 • A causal scenario described in natural language, code, or mathematical equations, including
198 causal relationships among variables, known conditions, etc.
- 199 • A causal query, which is a question based on causal scenario

200 **Steps**

201 **1. Causal Graph Extraction:**

202 (a) For natural language scenarios, we identify variables and causal relationships, and
203 construct causal graphs ($G := (V, E)$) by implementing a pipeline consisting of semantic
204 parsing and coreference resolution modules. The semantic parsing module first uses
205 the Stanford Parser [12] to perform syntactic parsing and obtain the sentence structure.
206 Then, it applies Compositional Semantics [13] to recursively map the syntactic parse
207 tree to a logical form, based on the principle of compositionality. The coreference
208 resolution module uses techniques such as the mention-pair model [14] to determine
209 which mentions refer to the same entity, and merges the variables corresponding
210 to coreferent mentions. From the outputs of the semantic parsing and coreference
211 resolution modules, the pipeline automatically extracts variables from nouns and
212 noun phrases, and identifies causal relationships indicated by verbs and conjunctions
213 expressing causality [15]. Finally, the causal graph construction module takes the
214 extracted variables as nodes (V) and causal relationships as directed edges (E) to
215 automatically build the causal graph [1].

216 (b) For code scenarios, we identify variables and their dependencies, and construct causal
217 graphs by implementing a pipeline that analyzes the code structure, control flow, and
218 data flow. The pipeline first uses a code parser, such as the ast module [17] in Python, to
219 generate an abstract syntax tree (AST). It then performs control flow analysis using tech-
220 niques like control flow graphs (CFGs) [18] and program dependence graphs (PDGs)
221 [21], and data flow analysis using def-use chains [19] and static single assignment
222 (SSA) form [20], to identify execution paths, dependencies between statements, and
223 variable dependencies. These analyses help automatically extract variables and their
224 relationships from the code structure. Finally, the causal graph construction module
225 takes the extracted variables as nodes (V) and their dependencies as edges (E) to build
226 the causal graph based on the code semantics [1], capturing the causal relationships
227 between variables and enabling further reasoning and analysis.

228 (c) For math scenarios, we identify variables and their functional relationships, and con-
229 struct causal graphs by implementing a pipeline that parses and analyzes the mathe-
230 matical equations. The pipeline first uses a math expression parser, such as the SymPy
231 library [22] in Python, to convert the equations into an abstract syntax tree (AST)
232 representation. It then traverses the AST to identify variables and their functional
233 relationships, such as dependencies and algebraic operations, using techniques like
234 symbolic differentiation [23] and expression simplification [24]. These analyses help
235 automatically extract variables and their relationships from the equation structure.
236 Finally, the causal graph construction module takes the extracted variables as nodes (V)
237 and their functional relationships as directed edges (E) to build the causal graph based
238 on the equation semantics, similar to the approach in [1]. The resulting causal graph
239 captures the causal relationships between variables in the mathematical equations,
240 enabling further reasoning and analysis.

241 2. **Query Classification:** Classify the causal query into one of the three levels of the Ladder
242 of Causation (Association, Intervention, Counterfactuals). Formalize the query into the
243 corresponding causal language, as discussed in [4].

244 3. **Estimand Derivation:**

245 (a) For text and math scenarios, we construct a module that uses causal inference algorithms
246 (e.g., do-calculus [25], counterfactual inference formulas [26]) to derive the estimand
247 based on the causal graph and query type.

248 (b) For code scenarios, we use program analysis techniques (e.g., symbolic execution, data
249 dependency analysis, control flow analysis) to derive the estimand based on the code
250 structure and query type. This involve simulating interventions on code variables and
251 analyzing the resulting program behavior.

252 4. **Data Matching:** Match the terms in the estimand with the available data or constraints
253 in the scenario to obtain a computable estimand expression. Check the completeness and
254 consistency of the data. Raise warnings or errors if critical data is missing. For code
255 scenarios, this involve executing the code with specific inputs and observing the outputs.
256 This step is similar to the data matching phase in [4].

257 5. **Causal Effect Estimation:**

258 (a) Calculate the causal effect value based on the estimand expression and the available
259 data, yielding the answer to the query.

260 (b) For scenarios with unobserved confounders, use instrumental variable estimation [27]
261 or front-door adjustment [25].

262 (c) For code scenarios, this involve comparing program behaviors under different interven-
263 tions.

264 This step is inspired by causal effect estimation phase in [4].

265 **Output**

- 266 • Answer to the causal query, including the estimated causal effect, confidence interval, and
267 key assumptions.

268 In a summary, our Causal Inference Engine extends the original design presented in [4] by incorporat-
269 ing domain-specific graph extraction and estimand derivation techniques to handle causal inference
270 problems in text, code, and math scenarios. The overall pipeline remains consistent with the one
271 described in [4], but the internal methods are adapted to the specific structures and semantics of each
272 domain.

273 **3.4.2 Quality Control Process**

274 After expansion with GPT4-Turbo, we obtained around 10000 x 4 text-based questions, 2000 x 4
275 math questions, and 2000 x 4 coding questions, along with their GPT-4 Turbo generated answers.
276 To ensure the accuracy of the ground truth of each questions, we employed a strict quality control
277 process as showing below:

278 We used the causal inference engine introduced above to independently solve the problems and
279 generate its own set of answers. We compared the answers generated by GPT-4 Turbo and the causal
280 inference engine. If two answers were the same, we updated the answer as ground truth. If any of the
281 answers were inconsistent, we conducted a manual analysis of the question and answers to determine
282 the correct answer and update ground truth accordingly.

283 This multi-step quality control process, involving the use of causal inference engine and human
284 expert check, ensures that the final dataset contains accurate and reliable questions and answers. The
285 manual review of inconsistent answers further enhances the quality of the dataset by addressing any
286 discrepancies or edge cases that the models may encounter.

287 **4 Benchmark Results**

288 **4.1 Baseline of Mainstream LLMs**

289 We tested several state-of-the-art large language models, including GPT-4, Claude-3, LLAMA-3, and
290 others, on our CausalBench. The evaluation metrics included: Four-Type Questions Group Correction
291 Rate, Overall Correction Rate (Ignore Question Type), From Cause to Effect without Intervention
292 Correction Rate, From Effect to Cause without Intervention Correction Rate, From Cause to Effect
293 with Intervention Correction Rate, and From Effect to Cause with Intervention Correction Rate. For
294 each causal scenario, there are four questions: cause-to-effect without intervention, effect-to-cause
295 without intervention, cause-to-effect with intervention, and effect-to-cause with intervention. The
296 Four-Type Questions Group Correction Rate represents the proportion of scenario cases where all
297 four types of questions of one scenario are all answered correctly by the large language models.
298 If any of the four questions of a scenario is answered incorrectly, the scenario is considered to be
299 answered incorrectly by the LLM. The Overall Correction Rate (Ignore Question Type) is calculated
300 by dividing the total number of correctly answered questions by the total number of questions, without
301 categorizing the questions by type and scenario. The From Cause to Effect without Intervention
302 Correction Rate is calculated by dividing the number of correctly answered "From Cause to Effect
303 without Intervention" type questions by the total number of this type of questions. Similarly, the
304 From Effect to Cause without Intervention Correction Rate is calculated by dividing the number of
305 correctly answered "From Effect to Cause without Intervention" type questions by the total number of
306 this type of questions. The remaining two metrics, From Cause to Effect with Intervention Correction
307 Rate and From Effect to Cause with Intervention Correction Rate, follow the same calculation method
308 as the previous two metrics, focusing on their respective question types.

309 Here are the tables showing LLMs' performance on text, math, and code problems.

Model	Four-Type Questions Group Correction Rate (%)	Overall Correction Rate (Ignore Question Type) (%)	From Cause to Effect without Intervention Correction Rate (%)	From Effect to Cause without Intervention Correction Rate (%)	From Cause to Effect with Intervention Correction Rate (%)	From Effect to Cause with Intervention Correction Rate (%)
GPT-4 Turbo	36.9	73.3	74.4	71.2	73.8	73.7
Claude3-Opus	36.8	72.6	74.1	70.9	73.2	72.2
Mistral-7B	25.5	63.6	58.7	66.5	64.2	65.0
Llama3-70B	21.8	61.5	62.6	59.6	63.8	60.1
Llama2-7B	20.7	62.1	62.8	64.0	56.4	65.4
GPT-3.5	16.7	57.8	57.6	58.5	56.2	58.7
Gemma-7b-it	12.8	50.7	50.0	46.9	53.6	52.1
Bloomz	4.2	41.7	41.0	40.7	41.7	43.6
AquilaChat	1.9	31.1	28.7	32.4	33.1	30.4

Table 1: LLM Performance on Text Problems

Model	Four-Type Questions Group Correction Rate (%)	Overall Correction Rate (Ignore Question Type) (%)	From Cause to Effect without Intervention Correction Rate (%)	From Effect to Cause without Intervention Correction Rate (%)	From Cause to Effect with Intervention Correction Rate (%)	From Effect to Cause with Intervention Correction Rate (%)
Mistral-7B	62.0	87.2	78.9	85.6	85.3	98.9
GPT-4 Turbo	61.4	88.7	78.6	88.3	91.7	96.0
Claude3-Opus	54.6	85.9	74.7	87.1	86.5	95.4
Llama3-70B	40.8	80.7	56.8	86.8	82.0	97.1
Gemma-7b-it	38.3	79.2	50.4	82.8	91.1	92.0
AquilaChat	25.3	68.1	57.0	67.8	69.2	78.3
Bloomz	23.9	69.2	53.3	76.8	67.3	79.7
GPT-3.5	15.9	63.3	47.1	71.5	48.6	86.1
Llama2-7B	2.8	42.3	45.3	54.2	17.5	52.4

Table 2: LLM Performance on Math Problems

Model	Four-Type Questions Group Correction Rate (%)	Overall Correction Rate (Ignore Question Type) (%)	From Cause to Effect without Intervention Correction Rate (%)	From Effect to Cause without Intervention Correction Rate (%)	From Cause to Effect with Intervention Correction Rate (%)	From Effect to Cause with Intervention Correction Rate (%)
Llama3-70B	43.8	77.0	82.0	75.7	73.9	76.0
Claude3-Opus	39.6	71.3	78.6	71.3	68.7	66.5
GPT-4 Turbo	37.2	71.0	80.6	67.5	73.2	62.5
Gemma-7b-it	32.3	68.4	74.1	67.7	66.0	65.4
Mistral-7B	31.4	66.8	67.5	68.3	61.3	70.2
GPT-3.5	25.0	64.5	71.9	65.4	59.8	60.6
Llama2-7B	22.6	61.9	79.0	45.5	76.3	46.8
Bloomz	17.5	52.4	49.6	56.8	46.4	56.8
AquilaChat	14.7	47.3	36.8	56.4	38.9	57.2

Table 3: LLM Performance on Code Problems

310 4.2 Test Result Summary

311 The evaluation results of state-of-the-art large language models on CausalBench provide valuable
 312 insights into their causal reasoning capabilities across textual, mathematical, and coding problem
 313 domains:

314 Overall, the models achieved higher correction rates on mathematical problems compared to textual
 315 and coding problems. For instance, GPT-4 achieved an 88.7% overall correction rate on math

316 problems, while scoring 73.3% and 71.0% on text and code problems, respectively. This suggests
 317 that causal reasoning in mathematical contexts is relatively easier for LLMs compared to natural
 318 language and programming domains.

319 The Four-Type Questions Group Correction Rate, which measures the proportion of scenarios where
 320 all four reasoning perspectives are correctly answered, was consistently lower than the Overall
 321 Correction Rate (Ignore Question Type) across all problem types. For example, GPT-4 achieved a
 322 61.4% Four-Type Questions Group Correction Rate on math problems, compared to an 88.7% Overall
 323 Correction Rate. This indicates that LLMs often struggle to maintain a comprehensive understanding
 324 of causal relationships when questioned from multiple perspectives.

325 The introduction of interventions in the causal scenarios led to mixed results in correction rates across
 326 models and problem types. In the text domain, the correction rates slightly decreased for most models
 327 when interventions were introduced. However, in the math domain, the correction rates generally
 328 improved with interventions. For instance, GPT-4’s performance increased from 78.6% to 91.7% on
 329 cause-to-effect questions with intervention in math problems. In the coding domain, the impact of
 330 interventions varied across models, with some showing improvements and others exhibiting a decline
 331 in performance.

332 Among the tested models, GPT-4 and Claude-3 consistently outperformed other large language models
 333 (LLMs) across most problem types and reasoning dimensions, achieving the highest correction rates.
 334 Mistral demonstrated strong performance in mathematical problems but exhibited shortcomings in
 335 code-related tasks. Conversely, LLAMA-3 showed robust performance in code-related problems but
 336 faced challenges with text and mathematical tasks.

337 5 Correlation with Hallucination

338 To analyze the correlation between LLMs’ causal reasoning ability and their hallucination rate, we
 339 referred to the LLMs’ performance on hallucination datasets. The hallucination evaluation results
 340 were obtained from the Hallucination Leaderboard, developed by Vectara [30]. This leaderboard
 341 provides a comparison of LLM performance in maintaining a low hallucination rate and ensuring
 342 factual consistency when summarizing a set of facts.

Model	Hallucination Rate	Factual Consistency Rate	Answer Rate	Average Summary Length (Words)
GPT 4 Turbo	2.5 %	97.5 %	100.0 %	86.2
Llama3-70B	4.5 %	95.5 %	99.2 %	68.5
Mistral 7B Instruct-v0.2	4.5 %	95.5 %	100.0 %	106.1
Llama2-7B	5.6 %	94.4 %	99.6 %	119.9
Claude3-Opus	7.4 %	92.6 %	95.5 %	92.1
Google Gemma-7b-it	7.5 %	92.5 %	100.0 %	113.0

Table 4: Performance of LLMs on the Hallucination Dataset

343 The hallucination evaluation process involves measuring the hallucination rate, factual consistency
 344 rate, answer rate, and average summary length. These metrics provide a comprehensive understanding
 345 of each model’s tendency to hallucinate and its ability to maintain factual accuracy [30].

346 After comparing the LLMs’ performance on CausalBench with their performance on the Hallucination
 347 evaluation leaderboard provided by Vectara on Huggingface [30], we found that models with stronger
 348 causal reasoning abilities tend to exhibit lower hallucination rates. For instance, GPT-4 Turbo,
 349 LLAMA-3-70B, and Mistral-7B, which demonstrated superior performance on causal reasoning
 350 tasks, also had low hallucination rates. In contrast, models like Google Gemma-7b-it and LLAMA-2-
 351 7B, which showed weaker performance on our CausalBench, had higher hallucination rates of 7.5%
 352 and 5.6%, respectively.

353 This trend indicates a potential link between a model’s ability to understand and reason about causal
 354 relationships and its likelihood of not producing hallucinations. Further research is required to explore
 355 this correlation in more depth and to understand the underlying mechanisms driving this relationship.

356 **6 Impact and Limitations**

357 **6.1 Impact**

358 For the first time, we innovatively propose four types of questioning approaches for the same causal
359 scenario: cause-to-effect, effect-to-cause, cause-to-effect with intervention, and effect-to-cause with
360 intervention. We also calculate the proportion of cases where large language models correctly answer
361 all four types of questions for a given causal scenario. This effectively avoids the situation where
362 large language models coincidentally answer causal questions correctly without understanding the
363 causal relationships embedded in the causal scenario, thereby improving the accuracy of the dataset’s
364 test results. By providing causal reasoning problems spanning multiple domains(text, code, math), it
365 addresses the limitations of existing causal datasets and offers a more comprehensive and robust tool
366 for assessing the causal reasoning abilities of language models. The findings in this paper suggest that
367 models with stronger causal reasoning capabilities tend to exhibit lower hallucination rates, providing
368 a new perspective on exploring the relationship between causal reasoning and reducing hallucinations.
369 CausalBench has the potential to become a benchmark for driving progress in causal reasoning in
370 artificial intelligence.

371 **6.2 Limitations**

372 CausalBench has several limitations that need to be addressed in future work. These include the need
373 for further expanding the domain coverage, increasing the scale of the dataset, incorporating causal
374 discovery tasks and exploring the intrinsic mechanisms between causal reasoning and hallucinations
375 through more empirical studies.

376 **7 Conclusion**

377 In this paper, we present CausalBench, a comprehensive benchmark dataset for evaluating the causal
378 reasoning capabilities of large language models. CausalBench innovatively proposes four types of
379 questioning approaches for each causal scenario: cause-to-effect, effect-to-cause, cause-to-effect with
380 intervention, and effect-to-cause with intervention. By calculating the proportion of cases where
381 models correctly answer all four question types, CausalBench effectively assesses whether LLMs
382 truly understand the underlying causal relationships, mitigating the impact of models coincidentally
383 providing correct answers without causal comprehension.

384 The dataset encompasses a diverse set of problems spanning textual, mathematical, and coding
385 domains, addressing the limitations of existing causal reasoning benchmarks. Evaluated on Causal-
386 Bench, state-of-the-art LLMs demonstrate stronger performance on mathematical problems compared
387 to textual and coding tasks. Notably, models with superior causal reasoning abilities tend to exhibit
388 lower hallucination rates, suggesting a potential link between the two capabilities.

389 Despite its contributions, CausalBench has several limitations, including the need for expanded
390 domain coverage and deeper exploration of the intrinsic mechanisms connecting causal reasoning and
391 hallucination reduction. Future work will focus on addressing these limitations, further refining the
392 evaluation metrics, and providing insights to advance the development of causal reasoning abilities in
393 large language models. CausalBench serves as a robust tool and an important step towards achieving
394 artificial general intelligence.

395 **References**

- 396 [1] Judea Pearl. 2009. Causality: Models, Reasoning, and Inference. *Cambridge University Press*.
- 397 [2] Bernhard Schölkopf. 2019. Causality for machine learning. *arXiv preprint arXiv:1911.10500*.
- 398 [3] Lex Fridman and Judea Pearl. 2022. Causal Reasoning, Counterfactuals, and the Path to AGI.
399 *Miniature Brain Machinery Webinar Review*.
- 400 [4] Zhijing Jin, Jiarui Liu, Zhiheng Lyu, Spencer Poff, Mrinmaya Sachan, Rada Mihalcea, Mona
401 T. Diab, and Bernhard Schölkopf. 2023. Can large language models infer causation from
402 correlation? CoRR, abs/2306.05836.

- 403 [5] Divyansh Kaushik, Eduard Hovy, and Zachary C Lipton. 2020. Learning the difference that
404 makes a difference with counterfactually-augmented data. In *8th International Conference on*
405 *Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenRe-
406 view.net.
- 407 [6] Judea Pearl. 1995. Causal diagrams for empirical research. *Biometrika*, 82(4):669–688.
- 408 [7] Judea Pearl et al. 2000. Causality: Models, reasoning and inference. Cambridge University
409 Press.
- 410 [8] Jörg Frohberg and Frank Binder. 2022. CRASS: A novel data set and benchmark to test
411 counterfactual reasoning of large language models. In *Proceedings of the Thirteenth Language*
412 *Resources and Evaluation Conference*, pages 2126–2140, Marseille, France. European Language
413 Resources Association.
- 414 [9] Matthew Ho, Aditya Sharma, Justin Chang, Michael Saxon, Sharon Levy, Yujie Lu, and William
415 Yang Wang. 2022. Wikiwhy: Answering and explaining cause-and-effect questions. *arXiv*
416 *preprint arXiv:2210.12152*.
- 417 [10] Leshem Choshen, Paarth Neekhara, Kyle Richardson, Lisa Xue, Madian Hou, Shehzaad
418 Neekhara, Yao Chen, and Heike Adel. 2022. CLADDER: A Causal Language Model for Causal
419 Reasoning. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language*
420 *Processing*, pages 6205–6224.
- 421 [11] Zhijing Jin, Jiarui Liu, Zhiheng Lyu, Spencer Poff, Mrinmaya Sachan, Rada Mihalcea, Mona
422 T. Diab, and Bernhard Schölkopf. 2023. Can large language models infer causation from
423 correlation? *CoRR*, abs/2306.05836.
- 424 [12] Dan Klein and Christopher D. Manning. 2003. Accurate Unlexicalized Parsing. In *Proceedings*
425 *of the 41st Annual Meeting of the Association for Computational Linguistics*, pages 423–430.
- 426 [13] Luke S. Zettlemoyer and Michael Collins. 2005. Learning to Map Sentences to Logical Form:
427 Structured Classification with Probabilistic Categorical Grammars. In *Proceedings of the 21st*
428 *Conference on Uncertainty in Artificial Intelligence*, pages 658–666.
- 429 [14] Wee Meng Soon, Hwee Tou Ng, and Daniel Chung Yong Lim. 2001. A Machine Learning Ap-
430 proach to Coreference Resolution of Noun Phrases. *Computational Linguistics*, 27(4):521–544.
- 431 [15] Pengfei Li and Kezhi Mao. 2019. Knowledge-oriented Convolutional Neural Network for
432 Causal Relation Extraction from Natural Language Texts. *Expert Systems with Applications*,
433 115:512–523.
- 434 [16] Emre Kıcıman, Robert Ness, Amit Sharma, and Chenhao Tan. 2023. Causal reasoning and large
435 language models: Opening a new frontier for causality. *arXiv preprint arXiv:2305.00050*.
- 436 [17] Python Software Foundation. 2023. ast — Abstract Syntax Trees. <https://docs.python.org/3/library/ast.html>. Accessed: 2023-06-05.
- 437
- 438 [18] Frances E. Allen. 1970. Control flow analysis. *ACM SIGPLAN Notices*, 5(7):1-19.
- 439 [19] Mary Jean Harrold and Gregg Rothermel. 1994. Performing data flow testing on classes. In
440 *Proceedings of the 2nd ACM SIGSOFT Symposium on Foundations of Software Engineering*
441 *(SIGSOFT’94)*, pages 154-163.
- 442 [20] Ron Cytron, Jeanne Ferrante, Barry K. Rosen, Mark N. Wegman, and F. Kenneth Zadeck. 1991.
443 Efficiently computing static single assignment form and the control dependence graph. *ACM*
444 *Transactions on Programming Languages and Systems (TOPLAS)*, 13(4):451-490.
- 445 [21] Jeanne Ferrante, Karl J. Ottenstein, and Joe D. Warren. 1987. The program dependence graph
446 and its use in optimization. *ACM Transactions on Programming Languages and Systems*
447 *(TOPLAS)*, 9(3):319–349.
- 448 [22] Aaron Meurer et al. 2017. SymPy: Symbolic computing in Python. *PeerJ Computer Science*,
449 3:e103.

- 450 [23] Andreas Griewank and Andrea Walther. 2008. *Evaluating Derivatives: Principles and Tech-*
451 *niques of Algorithmic Differentiation* (2nd ed.). SIAM.
- 452 [24] Joel Moses. 1971. Algebraic simplification: A guide for the perplexed. *Communications of the*
453 *ACM*, 14(8):527-537.
- 454 [25] Judea Pearl. 1995. Causal diagrams for empirical research. *Biometrika*, 82(4):669–688.
- 455 [26] Judea Pearl et al. 2000. *Causality: Models, reasoning and inference*. Cambridge University
456 Press.
- 457 [27] Joshua D. Angrist, Guido W. Imbens, and Donald B. Rubin. 1996. Identification of causal
458 effects using instrumental variables. *Journal of the American Statistical Association* 91, no.
459 434: 444-455.
- 460 [28] Jörg Frohberg and Frank Binder. 2022. CRASS: A novel data set and benchmark to test
461 counterfactual reasoning of large language models. In *Proceedings of the Thirteenth Language*
462 *Resources and Evaluation Conference*, pages 2126–2140, Marseille, France. European Language
463 Resources Association.
- 464 [29] Judea Pearl and Dana Mackenzie. 2018. *The Book of Why: The New Science of Cause and*
465 *Effect*. Basic Books.
- 466 [30] Simon Hughes and Minseok Bae. 2023. *Vectara Hallucination Leaderboard*. Vectara, Inc.
467 <https://github.com/vectara/hallucination-leaderboard>.
- 468 [31] Matej Žečević, Moritz Willig, Devendra Singh Dhama, and Kristian Kersting. 2023. Causal
469 parrots: Large language models may talk causality but are not causal. *Transactions on Machine*
470 *Learning Research*.
- 471 [32] Cheng Zhang, Stefan Bauer, Paul Bennett, Jiangfeng Gao, Wenbo Gong, Agrin Hilmkil, Joel
472 Jennings, Chao Ma, Tom Minka, Nick Pawlowski, et al. 2023. Understanding causality with
473 large language models: Feasibility and opportunities. *arXiv preprint arXiv:2304.05524*.
- 474 [33] Hector Geffner. 2018. Model-based vs. Model-free Reinforcement Learning: A Tutorial. In
475 *European Summer School on Logic, Language and Information*.
- 476 [34] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. 2016. *Deep Learning*. MIT Press.
477 <http://www.deeplearningbook.org>.
- 478 [35] Diederik P. Kingma and Max Welling. 2013. Auto-Encoding Variational Bayes. *arXiv preprint*
479 *arXiv:1312.6114*.
- 480 [36] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei Rusu, Joel Veness, Marc Bellemare,
481 Alex Graves, Martin Riedmiller, Andreas Fidjeland, Georg Ostrovski, Stig Petersen, Charles
482 Beattie, Amir Sadik, Ioannis Antonoglou, Helen King, Dhharshan Kumaran, Daan Wierstra,
483 Shane Legg, and Demis Hassabis. 2015. Human-level control through deep reinforcement
484 learning. *Nature*, 518(7540):529–533.
- 485 [37] David Silver, Aja Huang, Chris J. Maddison, Arthur Guez, Laurent Sifre, George van den
486 Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot,
487 Sander Dieleman, Dominik Grewe, John Nham, Nando de Freitas, Shakir Mohamed, Thor
488 Graepel, Timothy P. Lillicrap, Martin Riedmiller, and Demis Hassabis. 2016. Mastering the
489 game of Go with deep neural networks and tree search. *Nature*, 529(7587):484–489.
- 490 [38] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. 2015. Deep learning. *Nature*,
491 521(7553):436–444.
- 492 [39] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training
493 of Deep Bidirectional Transformers for Language Understanding. In *Proceedings of the 2019*
494 *Conference of the North American Chapter of the Association for Computational Linguistics:*
495 *Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186.

- 496 [40] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal,
 497 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel
 498 Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M.
 499 Ziegler, Jeff Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin,
 500 Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford,
 501 Ilya Sutskever, and Dario Amodei. 2020. Language Models are Few-Shot Learners. In *Advances*
 502 *in Neural Information Processing Systems*, volume 33, pages 1877–1901.
- 503 [41] Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. 2018. Improving
 504 Language Understanding by Generative Pre-Training. *OpenAI Blog*.

505 Checklist

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

540 **8 Appendix A: CausalBench Dataset Link**

541 <https://huggingface.co/datasets/CCLV/CausalBench>

542 **9 Appendix B: Text Question Example**

543 **Causal Scenario:**

544 Imagine a self-contained, hypothetical world with only the following conditions, and without any
545 unmentioned factors or causal relationships: Parents' intelligence has a direct positive effect on
546 parents' social status and child's intelligence. Other unobserved factors has a positive direct effect on
547 parents' social status and child's intelligence.

548 **Question 1:**

549 If a child is intelligent, would it be more likely that this child had intelligent parents?

550 **Question Type:**

551 Inference from Effect to Cause without Intervention

552 **Ground Truth:**

553 Yes

554 **Explanation:** The probability of the child being intelligent given that the parents are intelligent is
555 higher than the probability of the child being intelligent given that the parents are unintelligent, so
556 the ground truth is yes.

557 **Question 2:**

558 If the parents are intelligent, is the child more likely to be intelligent?

559 **Question Type:**

560 Inference from Cause to Effect without Intervention

561 **Ground Truth:**

562 Yes

563 **Explanation:** The probability of the child being intelligent given that the parents are intelligent is
564 higher than the probability of the child not being intelligent given that the parents are intelligent,
565 since parent's intelligence has positive effect on child's intelligence.

566 **Question 3:**

567 If we intervene to make the parents intelligent (e.g., through education or training), is the child more
568 likely to be intelligent?

569 **Question Type:**

570 Inference from Cause to Effect with Intervention

571 **Ground Truth:**

572 Yes

573 **Explanation:** By intervening to increase the parents' intelligence, the child's intelligence is more
574 likely to increase due to the causal chain. Although other unobserved factors also affect the child's
575 intelligence, the direct positive effect of parents' intelligence still exists.

576 **Question 4:**

577 If we observe a child is intelligent, and then intervene to make the child unintelligent (e.g., through
578 some kind of impairment), does this make it less likely that the child's parents are intelligent?

579 **Question Type:**

580 Inference from Effect to Cause with Intervention

581 **Ground Truth:**

582 No

583 **Explanation:**

584 The child's intelligence is the result of the combined effects of parents' intelligence and other factors.
585 Even if we intervene to decrease the child's intelligence, it does not change the parents' level of
586 intelligence. Therefore, in this case, the change in the child's intelligence does not affect our judgment
587 of whether the parents are intelligent or not.

588 **10 Appendix C: Code Question Example**

589 **Causal Scenario:**

```
590 class SalesData {
591     int totalSales , newSubscriptions;
592     double pricePerSubscription;
593
594     public SalesData(int newSubscribers , double price) {
595         this.newSubscriptions = newSubscribers;
596         this.pricePerSubscription = price;
597         updateSales ();
598     }
599
600     public void updatePrice(double newPrice) {
601         this.pricePerSubscription = newPrice;
602         updateSales ();
603     }
604
605     public void addSubscriptions(int additionalSubs) {
606         this.newSubscriptions += additionalSubs;
607         updateSales ();
608     }
609
610     private void updateSales () {
611         totalSales = (int) (newSubscriptions * pricePerSubscription);
612     }
613
614     public int getTotalSales () {
615         return totalSales;
616     }
617 }
618
619 SalesData monthlyReport = new SalesData(100, 10.0);
620 monthlyReport.addSubscriptions(50);
621 monthlyReport.updatePrice(15.0);
```

622 **Question 1:**

623 If the number of new subscriptions increases, will total sales also increase, assuming no other
624 changes?

625 **Question Type:**

626 From cause to effect without intervention

627 **Ground Truth:**

628 Yes

629 **Explanation:**

630 The method 'addSubscriptions' adds new subscriptions and then immediately calls 'updateSales',
631 which recalculates total sales based on the new number of subscriptions and the current price per

632 subscription. Therefore, with no other changes, increasing the number of new subscriptions directly
633 leads to an increase in total sales.

634 **Question 2:**

635 Does an increase in total sales imply an increase in the price per subscription?

636 **Question Type:**

637 From effect to cause without intervention

638 **Ground Truth:**

639 No

640 **Explanation:**

641 An increase in total sales can occur either from an increase in the price per subscription or from
642 an increase in the number of new subscriptions due to the calculation in 'updateSales'. Hence, an
643 increase in total sales does not necessarily imply that the price per subscription has increased; it could
644 also be due to an increase in the number of subscriptions.

645 **Question 3:**

646 If we manually increase the price per subscription, will this result in an increase in total sales?

647 **Question Type:**

648 From cause to effect with intervention

649 **Ground Truth:**

650 Yes

651 **Explanation:**

652 Increasing the price per subscription using 'updatePrice' method causes 'updateSales' to be called,
653 calculating the new total sales using the increased price. Assuming the number of subscriptions
654 remains constant, this intervention in price directly causes an increase in total sales.

655 **Question 4:**

656 If total sales decrease after an intervention, does this mean we decreased the number of new subscrip-
657 tions?

658 **Question Type:**

659 From effect to cause with intervention

660 **Ground Truth:**

661 No

662 **Explanation:**

663 A decrease in total sales after an intervention could be due to either a decrease in the number of new
664 subscriptions or a decrease in the price per subscription. As these two factors multiply to compute
665 total sales, the decrease could be attributed to either factor independently or both. Thus, a decrease in
666 total sales does not definitively determine that the intervention was a decrease in the number of new
667 subscriptions.

668 **11 Appendix D: Math Question Example**

669 **Causal Scenario:**

670 Investigate the influence of a linear operator transformation $z = L(x)$ on a vector field x governed by
671 $\frac{d}{dt}x = Mx$, where M is a constant matrix. The transformation L represents another linear operator
672 with a constant matrix.

673 **Question 1:**

674 If the transformation $z = L(x)$ is applied immediately at $t = 0$ to a vector x_0 , followed by evolution
675 under $\frac{d}{dt}x = Mx$ without further intervention, does the state $z(t)$ at $t = T$ exactly replicate the result
676 of evolving x_0 directly under $\frac{d}{dt}x = Mx$ until $t = T$?

677 **Question Type:**

678 From cause to effect without intervention

679 **Ground Truth:**

680 No

681 **Explanation:**

682 Applying the transformation $z = L(x)$ modifies the initial conditions. The trajectory of $z(t)$ and $x(t)$
683 would differ unless L commutes with the exponential of M , which generally is not the case. Hence,
684 the state transformations under L can produce a distinct evolutionary path in comparison to the direct
685 evolution of x_0 .

686 **Question 2:**

687 Can the original vector x_0 be reliably determined at $t = 0$ after observing the vector $z(t)$ at $t = T$,
688 without knowing if the transformation L was applied?

689 **Question Type:**

690 From effect to cause without intervention

691 **Ground Truth:**

692 No

693 **Explanation:**

694 Without information on the application of L , reconstructing the exact initial state x_0 from $z(t)$ is
695 not straightforward. The application of L can alter the vector in ways that are not easily reversible,
696 especially if L and M are not designed to reveal their effects straightforwardly.

697 **Question 3:**

698 If an additional linear transformation H is applied at time t_1 as an intervention, will the final state
699 $z(T)$ at $T > t_1$ be independent of the initial transformation L and solely determined by M and H ?

700 **Question Type:**

701 From cause to effect with intervention

702 **Ground Truth:**

703 No

704 **Explanation:**

705 The final state $z(T)$ will depend on L , M , and H . The transformations imposed by L initially, and
706 H later, both play critical roles. These factors, combined with the dynamics driven by M , contribute
707 to a state at T that relies on all three matrices, affected by their interaction and properties.

708 **Question 4:**

709 Based on knowing only the vector $z(T)$ at time T , is it feasible to precisely identify the transforma-
710 tions (L , H , or both) that were previously applied?

711 **Question Type:**

712 From effect to cause with intervention

713 **Ground Truth:**

714 No

715 **Explanation:**

716 Determining which transformations were applied based on the final vector $z(T)$ alone is challenging
717 due to the overlapping effects matrices may have in transforming the state space. The interactions of L

718 and H with the matrix exponential of M can result in equivalent states from different transformation
719 sequences.