

ROCK: A Causal Inference Framework for Reasoning about Commonsense Causality

Anonymous ACL submission

Abstract

Commonsense causality reasoning (CCR) aims at identifying plausible causes and effects in natural language descriptions that *deemed reasonable by an average person*. Although being of great academic and practical interest, this problem is still shadowed by the lack of a well-posed theoretical framework; existing work usually relies on various notions of correlation and is susceptible to *confounding co-occurrences*. This paper articulates the central question of CCR and develops a novel framework, ROCK, to Reason O(A)about Commonsense K(C)ausality based on classical causal inference principles. ROCK leverages temporal signals as incidental supervision, and makes use of *temporal propensities* that are analogous to propensity scores (Rosenbaum and Rubin, 1983) for balancing confounding effects. We implement a modular zero-shot pipeline, which is effective and demonstrates good potential for CCR on various datasets.

1 Introduction

Commonsense causality reasoning (CCR) has been an important yet non-trivial task in natural language processing (NLP) communities that exerts broad industrial and societal impacts (Kuipers, 1984; Gordon et al., 2012; Mostafazadeh et al., 2020; Sap et al., 2020). We articulate this task as

reasoning about plausible causes and effects of semantic meanings conveyed by natural languages that are accessible by an average reasonable person.

This definition naturally excludes questions that are beyond commonsense knowledge, such as those scientific in nature (e.g., does a surgery procedure reduce mortality?). Instead, it accommodates causal queries within the reach of an ordinary reasonable person. For example, in Figure 1, is Alice’s “entering a restaurant” (E_1) a plausible cause for her “ordering a pizza” (E_2)? Although the precedence

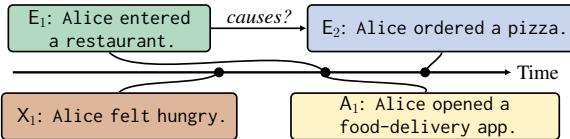


Figure 1: **An Example of CCR:** does E_1 cause E_2 ? The temporal order $E_1 \prec E_2$ does not necessitate causation due to confounding co-occurrences (e.g., X_1). Since when *conditioning on* X_1 , a *comparable* intervention A_1 of E_1 also precedes E_2 , the effect from E_1 to E_2 shrinks.

from E_1 to E_2 is logical, it might be less a “cause” compared with Alice’s “feeling hungry” (X_1).

We know that temporality certainly informs causation, but how to account for confounding co-occurrences such as X_1 ? Motivated by causal principles (Section 2), we formulate CCR as estimating the *change* in the likelihood of E_2 ’s occurrence due to intervening E_1 (denoted by $\neg E_1$):

$$\Delta = \mathbb{P}(E_1 \prec E_2) - \mathbb{P}(\neg E_1 \prec E_2) \quad (1)$$

where $\mathbb{P}(\cdot)$ can be estimated by language models (LMs) e.g., via mask language modeling (Section 4). If the occurrence of E_1 and $\neg E_1$ is purely random, Δ directly estimated informs CCR; however, due to confounding co-occurrences (X_1), one need to *balance* the covariates (events that precede E_1) to eliminate potential spurious correlations. In light of the capabilities and limitations of LMs, we propose *temporal propensity*, a surrogate propensity score for this purpose (Section 3). Indeed, we show in Section 5 that *although temporality is essential for CCR, it is vulnerable to spurious correlations without being properly balanced*.

Contributions. We articulate CCR from a completely new perspective using causal inference principles, and our contributions include (i) a novel causal framework for CCR that is the first of its kind to the best of our knowledge; (ii) mitigating confounding co-occurrences by matching temporal propensities; (iii) a modular pipeline for zero-shot CCR with demonstrated effectiveness.

071 2 Background

072 The problem of reasoning causal relationships, and
073 differentiating them from innocuous associations
074 has been contemplated and studied extensively in
075 human populations research spanning clinical trials,
076 epidemiology, political and social sciences, eco-
077 nomics, and many more (Fisher, 1958; Cochran and
078 Chambers, 1965; Rosenbaum, 2002; Imbens and
079 Rubin, 2015), among which causal practitioners
080 usually base their studies on the potential outcomes
081 framework (also known as the Rubin causal model,
082 see Neyman, 1923; Rubin, 1974; Holland, 1986)
083 and the structural equation model (Pearl, 1995;
084 Heckman, 2005; Peters et al., 2017) in observa-
085 tional studies; Granger causality (Granger, 1969)
086 for time series data, to name a few.

087 With the recent celebrated empirical success of
088 language models especially transformers (Devlin
089 et al., 2019; Radford et al., 2019), there is an in-
090 creasing interest in the NLP community on drawing
091 causal inference using textual data, where the ma-
092 jor focus treats textual data as either covariates or
093 study units (Keith et al., 2020; Feder et al., 2021)
094 on which causal queries are formed (e.g., does tak-
095 ing a medicine affects recovery, which are recorded
096 in textual medical records?). On the other hand,
097 CCR with natural language descriptions struggles
098 to fit in a causal inference framework: *textual data*
099 *in this case are just vehicles conveying semantic*
100 *meanings, not to be taken at the face value*, hence
101 it is difficult to draw the parallel between classical
102 causal inference that requires a clear definition of
103 study units, treatments, and outcomes.

104 2.1 Existing Approaches

105 Existing works related to CCR are usually grouped
106 under the umbrella term of commonsense reason-
107 ing (Rashkin et al., 2018; Ning et al., 2019a; Sap
108 et al., 2020) or event detection (O’Gorman et al.,
109 2016). Some of the notable progresses usually
110 come from leveraging explicit causal cues/links
111 (tokens such as “due to”) and use conditional prob-
112 abilities to measure “causality” (Chang and Choi,
113 2004; Do et al., 2011; Luo et al., 2016); leveraging
114 deep LMs via augmenting training datasets, design-
115 ing training procedures and/or loss functions (Sap
116 et al., 2019; Shwartz et al., 2020; Tamborrino et al.,
117 2020; Zhang et al., 2021; Staliunaite et al., 2021).

118 There are several works that share similarities
119 in perspective with ours, yet different in various
120 ways: Granger causality, which measures associa-

tion, is used by Kang et al. (2017) to detect event causes-and-effects; Bhattacharjya et al. (2020) studies events as point-processes, in a way arguably closer to association; Gerstenberg et al. (2021) uses a simulation model to reason physical causation. To the best of our knowledge, there is very little literature, if not none, on adopting a causal perspective in solving CCR.

121 2.2 Challenges of CCR

122 Many existing CCR methods (mostly supervised)
123 are based on ingenious designs and creative LM
124 engineering, theoretical justifications, however, are
125 sometimes desirable, as only then do we know how
126 general they can be. Indeed, recent studies reveal
127 that several supervised models may have exploited
128 certain artifacts in datasets to ace the evaluations
129 (Kavumba et al., 2019; Han and Wang, 2021).

130 This dilemma of constructing a well-founded
131 theoretical framework versus engineering models
132 to achieve excellent empirical performances is not
133 surprising, perhaps, given that the challenges of
134 CCR from causal perspectives are not trivial at all:
135 what is the study unit, treatment, and outcome in
136 this case? What does it mean to “intervene”, or
137 “manipulate” the treatment? Is treatment *stable*, or
138 is it desirable to consider multiple versions of it?

139 2.3 Principles of the ROCK Framework

140 In this paper, we attempt to these questions using,
141 among several causal principles, the following two
142 that are intuitive and directly appeal to human na-
143 ture (see e.g., Russell, 1912; Bunge, 1979):

- 144 • Precedence does not imply causation.
145 • Causation implies precedence.

146 Here the first principle warns us of *post-hoc* falla-
147 cies; the second informs us that the events must be
148 compared with those that are *in pari materia* (Mill,
149 1851; Hill, 1965), or having *balanced* covariates
150 (also called “pretreatments,” which we mean events
151 that occur prior to E_1 , cf. Rosenbaum, 1989). Our
152 CCR formulation in terms of temporality has sev-
153 eral benefits: (i) the intrinsic temporality of causal
154 principles characterizes its central role in CCR; (ii)
155 temporal signals bring about incidental supervision
156 (Roth, 2017; Ning et al., 2019a); (iii) although be-
157 ing a non-trivial question *per se*, reasoning tempo-
158 rarily has witnessed decent progresses lately, mak-
159 ing it more accessible than directly detecting causal
160 signals (Ning et al., 2017, 2018, 2019b; Zhou et al.,
161 2020; Vashishtha et al., 2020).

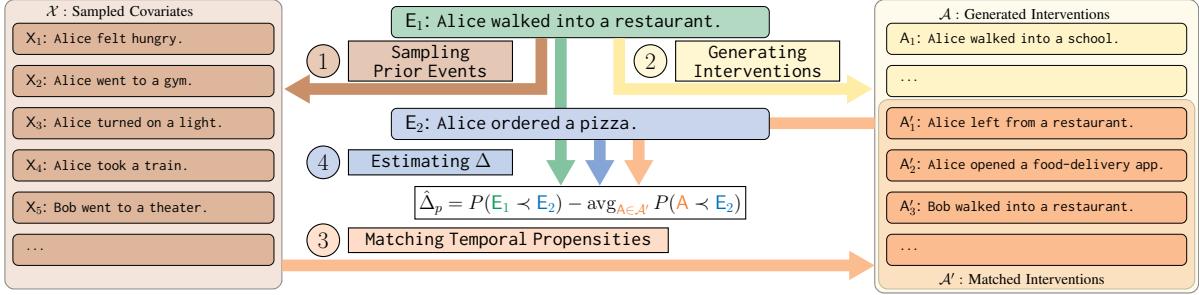


Figure 2: **Illustration of the ROCK framework.** Does E_1 cause E_2 ? To answer this query, ① the event sampler samples a set of covariates \mathcal{X} of events X_k that occur preceding E_1 . ② The intervention generator generates a set \mathcal{A} of interventions A_k on E_1 . ③ A subset $\mathcal{A}' \subset \mathcal{A}$ of interventions is selected whose temporal propensities $q(\mathbf{x}; \mathbf{A})$ is close to that of E_1 , $q(\mathbf{x}; E_1)$ (Equation (5)). ④ The temporal predictor uses \mathcal{A}' to estimate Δ .

3 The ROCK Framework

Notations. We use sans-serif letters for events, uppercase serif letters *indicators* of whether the corresponding event occurs¹, and lowercase serif letters the realizations of those indicators. For example, in Figures 1 and 2, E_1 : “Alice walked into a restaurant..,” $E_1 = \mathbb{1}\{E_1 \text{ occurs}\}$ and $e_{1,i} = \mathbb{1}\{E_1 \text{ occurs to the } i\text{-th study unit}\}$ ². We view the occurrence of events as point processes $E(t)$ on $t \in \{0, 1\}$ (e.g., present vs past). We write $E_1 \succ E_2$ (resp. \prec) to mean E_1 occurs succeeding (resp. preceding) E_2 . With a slight abuse of notation, we write $\mathbb{P}(E_1 \prec E_2) = \mathbb{P}(E_1(0), E_2(1))$ and $\mathbb{P}(E_2|E_1) = \mathbb{P}(E_2(1)|E_1(0))$. We write P for estimates of \mathbb{P} , and omit measure-theoretic details³.

Overview of the ROCK framework. We set the stage in this section and discuss implementation details in Section 4. Given E_1 and E_2 , as shown in Figure 2, ROCK samples the covariates set \mathcal{X} and interventions set \mathcal{A} , from which a matched subset \mathcal{A}' is selected via temporal propensities (Section 3.4). The score Δ is then estimated by Equation (5).

3.1 The Central Question of CCR

Given E_1 and E_2 , as aforementioned in Section 1, we articulate CCR as to estimate, the change of temporal likelihood *had* E_1 been *intervened*:

$$\Delta = \mathbb{P}(E_1 \prec E_2) - \mathbb{P}(\neg E_1 \prec E_2) \quad (2)$$

which assumes values in $[-1, 1]$ and is analogous to the *average treatment effect* in human populations research. Should we know the occurrence of E_1

¹By “occurs,” we mean “is observed.” We treat them interchangeable in the rest of our paper.

²Defined among other concepts in Section 3.2.

³Let \mathcal{E} be the set of commonsense events we consider, the probability space we are working on is $(\mathcal{E} \times \mathcal{E}, \sigma(\mathcal{E} \times \mathcal{E}), \mathbb{P})$.

is purely random, direct estimation is fair enough; however, since these probabilities are eventually estimated from training data, if there are confounding events X_k that always co-occurs with E_1 in the training data, they will bias this estimation. To this end, it is necessary to first clear out several key notions associated with this causal query, and then properly define the intervention $\neg E_1$.

3.2 Units, Covariates, and Interventions

One major challenge of framing a causal query for CCR is the ambiguity of the underlying mechanism. Unlike human populations research, where experiments and study units are obvious to define, it is not immediately clear what they are when faced with semantic meanings of languages. Yet, we can draw parallel again between semantic meanings and human subjects via the following thinking experiment: suppose each human-being keeps a journal/script detailing the complete timeline of her experiences since her conception, then we can treat each individual as a study unit where the temporal relations of events can be inferred from the journal. The treatment will then be the indicator $E_1 = \mathbb{1}\{E_1 \text{ occurs}\}$, the outcome $E_2 = \mathbb{1}\{E_2 \text{ occurs}\}$, and the covariates all events preceding E_1 . Hence a hypothetical observational study can be based on observations of a subpopulation of such journals, from which we form observed treatments $e_{1,1}, e_{1,2}, \dots, e_{1,k}$, the associated outcomes $e_{2,1}, e_{2,2}, \dots, e_{2,k}$, and covariates x_1, x_2, \dots, x_k (as vectors of indicators).

This identification naturally complies with the temporal nature of covariates (Rubin, 2005), since by definition they are *pretreatment* that take place *before* the treatment. We shall now address the issue of intervention (manipulation). Generally speaking, events are complex, and therefore inter-

vention in this case would be better interpreted in a broader sense than one particular type of manipulation such as negation. For example, with E_1 being “Alice walked into a restaurant,” suppose hypothetically, before E_1 , Alice did not walk into a restaurant ($\neg E_1$), we can thus compare $\mathbb{P}(E_1 \prec E_2)$ with $\mathbb{P}(\neg E_1 \prec E_2)$ to reason to what extent some event E_2 can be viewed as the effect due to E_1 . However, this is not the complete picture: Alice may have walked into somewhere else such as a bar; she may have, instead of walked into, but left a restaurant; instead of Alice, perhaps it was Bob who walked into a restaurant. The temporal information between these events and E_2 are also likely to inform causation between E_1 and E_2 , and they are no less interventions than negation. As such we interpret intervention in our framework in a broader sense, not necessarily only negation or the entailment of negations, but *any events that leads to plausible states of counterfactuality*. We will denote all possible interventions of E_1 as \mathcal{A} .

Remark. The well-celebrated *stable unit treatment value assumption* (SUTVA, Rubin, 1980) requires that for each unit there is only one version of the non-treatment. We can view the outcome of our framework in Equation (1) to be the temporal probability threshold at a certain level, $\mathbb{1}\{\mathbb{P}(\cdot \prec E_2) > 1 - \delta\}$, thus complying SUTVA.

3.3 Balancing Covariates

Direct estimation of Δ in Equation (1) is feasible only in an ideal world where those probabilities are estimated from randomized controlled trials (RCTs) such that the treatment (E_1) is assigned completely at random to study units. Due to confounding co-occurrences, events that precede E_1 need to be properly balanced (Mill, 1851; Rosenbaum and Rubin, 1983; Pearl and Mackenzie, 2018). Taking again as an example E_1 : “Alice walked into a restaurant,” and E_2 : “Alice ordered a pizza.” Suppose hypothetically, Alice’s twin sister Alicia, who has the exactly same life experiences up to the point when E_1 took place, but opted not walking into a restaurant, but opened a food delivery app on her phone ($\neg E_1$). Then we can reason that the cause-and-effect relationship from E_1 to E_2 is perhaps not large. On the other hand, if we know another irrelevant person, say Bob, undergone $\neg E_1$ and then E_2 , then perhaps we are not ready to give that conclusion since we do not know if Bob and Alice are comparable at the first place.

This example illustrates the importance of adjusting or balancing pretreatments. Hence we rewrite Equation (1) as conditional expectations among study units that are comparable, i.e.,

$$\mathbb{E}_{\mathbf{x}} [\mathbb{P}(E_1 \prec E_2 | \mathbf{x}) - \mathbb{P}(\neg E_1 \prec E_2 | \mathbf{x})], \quad (3)$$

where \mathbf{x} is the vector of covariates.

Remark. (i) We should define \mathbf{x} as events preceding E_1 , but *not* E_2 , which will potentially introduce posttreatment biases (Rosenbaum, 1984): if an X' that occurs between E_1 and E_2 is adjusted, Δ thus estimated quantifies the effect from E_1 to E_2 *without* passing through X' . (ii) Although \mathbf{x} should be those that are correlated with E_1 , adjusting for un-correlated effects does not introduce biases.

3.4 Matching Temporal Propensities

There are several techniques for balancing covariates such as sub-classification, matched sampling, and covariance adjustment, and via structural equations (Cochran and Chambers, 1965; Pearl, 1995). Rosenbaum and Rubin (1983) showed that such adjustment can be done using any balancing score, with the coarsest one being the propensity $p(\mathbf{x}) = \mathbb{P}(E_1(1) = 1 | \mathbf{x}(0))$, the probability of the occurrence of E_1 (at time 1) conditioning on the coavariates being \mathbf{x} (at time 0).

To properly identify what events constitute the covariate set is essential for our CCR framework. In the best scenario, it should include the real cause(s), which is, however, exactly what CCR solves. To circumvent this circular dependency, we use an LM to sample a large number of events preceding E_1 , which should provide a reasonable covariate set. In this case, directly computing $p(\mathbf{x})$ is less feasible; hence we propose to use a surrogate which we call *temporal propensities*:

$$q(\mathbf{x}) = q(\mathbf{x}; E_1) = (\mathbb{P}(E_1(1) = 1 | \mathbf{x}))_{\mathbf{x} \in \mathcal{A}} \quad (4)$$

with each coordinate measures the conditional probability of the event E_1 given an event in \mathbf{x} . Therefore, for some fixed threshold ϵ and $p \in \{1, 2\}$, we have the following estimating equation for the L_p -balanced score:

$$\hat{\Delta}_p = f(E_1, E_2) - \frac{1}{|\mathcal{A}'|} \sum_{A \in \mathcal{A}'} f(A, E_2), \quad (5)$$

$$\mathcal{A}' := \left\{ A \in \mathcal{A} : \frac{1}{|\mathcal{A}'|} \|q(\mathbf{x}; A) - q(\mathbf{x}; E_1)\|_p < \epsilon \right\}.$$

329 4 Implementation of ROCK

330 Having established a framework for CCR, we pro-
331 vide an exemplar implementation of ROCK in this
332 section. Our purpose is not to ace in every CCR
333 dataset, and we expect engineering efforts such as
334 prompt design can bring further improvements.

335 The core tool we will use is (fine-tuned) pre-
336 trained deep LMs. With the sheer amount of train-
337 ing data (e.g., over 800GB for the Pile dataset, [Gao et al. \(2020\)](#)), it is reasonable to assume that those
338 models would imitate responses of an average rea-
339 sonable person. On the other hand, it is hard for
340 generation models (masked or open-ended) to parse
341 information that are far from the mask tokens. It
342 is thus more feasible to use LMs to sample sum-
343 mary statistics of the relationships between a pair
344 of events, which is one of the main motivations for
345 using temporal propensities (Equation (4)).

347 4.1 Components of ROCK

348 For practical purposes, we represent an event as a
349 3-tuple $(\text{ARG0}, \text{V}, \text{ARG1})$. ROCK takes two events
350 E_1 and E_2 as inputs, and returns an estimate $\hat{\Delta}$
351 for Δ according to Equation (5). It contains four
352 components (cf. Figure 2): an event sampler that
353 samples a set \mathcal{X} of events that are likely to occur
354 preceding E_1 ; a temporal predictor whose output
355 $f(X_1, X_2)$ given two input events X_1 and X_2 is an
356 estimate of the temporal probability $\mathbb{P}(X_1 \prec X_2)$;
357 an intervention generator that generates a set \mathcal{A} of
358 events that are considered as interventions of the
359 event E_1 ; and finally a scorer that first forms the
360 temporal propensity vectors $q(\mathbf{x}; \mathbf{A}) \in \mathbb{R}^{|\mathcal{X}|}$ for
361 each sampled interventions $\mathbf{A} \in \mathcal{A}$, then estimates
362 Δ via Equation (5). We next discuss in greater
363 details our implementation of this pipeline.

364 4.2 Implementation Details

365 **Event Sampling.** Given an event E_1 (e.g., E_1 :
366 Alice walked into a restaurant.), we con-
367 struct the prompt by adding “Before that,” to
368 the sentence, forming “Alice walked into a
369 restaurant. Before that, ” as the final
370 prompt. We use the GPT-J model ([Wang and Ko-](#)
371 [matsuzaki, 2021](#)), which is pretrained on the Pile
372 dataset ([Gao et al., 2020](#)) for open-ended text gen-
373 eration where we set max length of returned se-
374 quences to be 30, temperature to be 0.9. We sample
375 $n = 100$ events, cropping at the first stop token of
376 the newly generated sentence to form \mathcal{X} .

377 **Temporal Prediction.** Given two events E_1 and
378 E_2 , we use mask language modeling to pre-
379 dict their temporal relation by form the prompt
380 $E_1 <\text{MASK}> E_2$ and collect the score $s_a(E_1, E_2)$ and
381 $s_b(E_1, E_2)$ for the tokens after and before. We
382 then symmetrize the estimates to form $s(E_1, E_2) =$
383 $\frac{1}{2}(s_a(E_1, E_2) + s_b(E_2, E_1))$. We can direct use
384 $s(E_1, E_2)$ for $f(E_1, E_2)$; we discuss possible nor-
385 malizations of this score in Section 5.

386 **Temporality Fine-Tuning.** Directly using a pre-
387 trained LM as the temporal predictor is likely to
388 suffer from low coverage, since the tokens before
389 and after usually are not among the top- k most
390 probable tokens. We can increase k but this does
391 not heuristically justify if the outputted scores are
392 meaningful. We thus use the New York Times
393 (NYT) corpus which contains NYT articles from
394 1987 to 2007 ([Sandhaus, 2008](#)) to fine-tune an LM.
395 Following the same procedure as [Zhou et al. \(2020\)](#),
396 we perform semantic role labeling (SRL) using Al-
397 lenNLP’s BERT SRL model ([Gardner et al., 2017](#))
398 to identify sentences with a temporal argument
399 (ARG-TMP) that starts with a temporal connective
400 tmp (either before or after). We then extract
401 the verb and its two arguments (V, ARG0, ARG1) as
402 well as of its temporal argument, thus forming an
403 event pair (E_1, E_2, tmp) . We are able to extract
404 397,174 such pairs and construct them into the
405 fine-tuning dataset consisting of “ $E_1 \text{ tmp } E_2$ ” and
406 “ $E_2 \neg \text{tmp } E_2$ ” for each extracted pair, where $\neg \text{tmp}$
407 is the reverse temporal connective (e.g., after if
408 tmp is before). We then fine-tune a pretrained
409 RoBERTa model (RoBERTa-BASE) using Hugging-
410 Face Transformers ([Wolf et al., 2020](#)) via mask lan-
411 guage modeling with masking probability $p = 0.1$
412 for each token. We choose a batch size of 500 and
413 a learning rate of 5×10^{-5} , and train the model to
414 convergence, which was around 135,000 iterations
415 when the loss converges to 1.37 from 2.02.

416 **Intervention Generator.** Given event E_1 , the in-
417 tervention generator generates a set \mathcal{A} of events
418 that are considered as interventions of the event A
419 in the sense of Section 3.2, which includes ma-
420 nipulating ARG0, V, and ARG1 respectively. We
421 achieve this goal by masking these components in-
422 dividually and filling in the masks using an LM.
423 There are several existing works on generating
424 interventions of sentences ([Feder et al., 2021](#)),
425 and we select PolyJuice ([Wu et al., 2021](#)) in
426 our pipeline due to its robustness. PolyJuice al-

427
428
429
430
431
432
433

lows conditional generation via control codes such as negation, lexical, resemantic, quantifier, insert, restructure, shuffle, and delete, each corresponds to different flavors on which the sentence is intervened. We drop the fluency-evaluation component of PolyJuice as they will be evaluated by the temporal predictor.

434
435
436
437
438
439
440
441
442
443
444

Score Estimation. Given the interventions \mathcal{A} and the sampled covariates \mathcal{X} , we can use the temporal predictor to estimate $\mathbb{P}(X \prec A)$ for all $X \in \mathcal{X}$ and $A \in \mathcal{A}$. To obtain temporal propensities $q(x; A)$ for all interventions, we need to estimate $\mathbb{P}(A = 1|X)$ for each X and A . Since by our sampling method, X occurs preceding E_1 , there is an implicit conditioning on E_1 , we may thus set $P(X(0)) = f(X, E_1)$ and $P(X(0), A(1)) = f(X, A)$; we will discuss possible normalizations in Section 5.2. We then form temporal propensity vectors as

$$445 \quad q(x; A) = \left(\frac{P(X(0))}{P(X(0), A(1))} \right)_{X \in \mathcal{X}}. \quad (6)$$

446 5 Empirical Studies

447 We put the ROCK framework into action, our findings reveal that *although temporality is essential*
448 *for CCR, without balancing covariates, it is prone*
449 *to spurious correlations.*

451 5.1 Setup and Details

452 **Evaluation Datasets.** We evaluate the ROCK
453 framework on the Choice of Plausible Alternatives dataset (COPA, Gordon et al., 2012) and a
454 self-constructed dataset of 153 instances using the
455 first dimension (cause-and-effect) of GLUCOSE
456 (GLUCOSE-D1, Mostafazadeh et al., 2020). Each
457 instance in COPA consists of a premise, two plau-
458 sible choices, and a question type asking which
459 choice is the choice (or effect) of the premise.
460 When asking for cause, we set the premise as E_1 ,
461 and two choices as E_2 respectively; otherwise we
462 take the premise as E_2 and two choices as E_1 respec-
463 tively. We choose the choice with the higher score.
464 We evaluate the development set of 100 instances
465 (COPA-DEV) and the test set of 500 instances
466 (COPA-TEST). To construct GLUCOSE-D1, we
467 take the test set and set the cause as premise, the
468 effect and another candidate event as two choices
469 then follow the same procedure.

471 **Baseline Scores and Variants.** To test the validity
472 and the effectiveness of ROCK, we compare

473 with the adjusted score $\hat{\Delta}_p$ with several other rea-
474 sonable scores that may be intuitive at first sight.

- 475 • L_1 -balanced score $\hat{\Delta}_1$: set $p = 1$ in (5).
- 476 • L_2 -balanced score $\hat{\Delta}_2$: set $p = 2$ in (5).
- 477 • Vanilla temporal score $\hat{\Delta}_{E_1} = \mathbb{P}(E_1 \prec E_2)$.
- 478 • Unadjusted score $\hat{\Delta}_{\mathcal{A}}$: set $\mathcal{A}' = \mathcal{A}$ in (5).
- 479 • Misspecified score $\hat{\Delta}_{\mathcal{X}}$: set $\mathcal{A}' = \mathcal{X}$ in (5).

480 Here the L_p -balanced scores are those balanced
481 using temporal propensities with L_p norm in Equa-
482 tion (5); the vanilla temporal score is perhaps the
483 most straightforward one, which treats temporal
484 precedence as causation; the unadjusted score is
485 obtained without balancing the covariates; the mis-
486 specified score mistakes the covariates for inter-
487 ventions. All these three have intuitive explana-
488 tions but are either insufficient for CCR or prone to
489 spurious correlations. Note that $\lim_{\epsilon \downarrow 0} \hat{\Delta}_p = \hat{\Delta}_{E_1}$
490 (when nothing is kept) and $\lim_{\epsilon \uparrow 1} \hat{\Delta}_p = \hat{\Delta}_{\mathcal{A}}$ (when
491 everything is kept).

492 5.2 Design Choices and Normalizations

493 We discuss several design choices and normaliza-
494 tions that might stabilize estimation procedures.
495 As shown in Section 5.4, although some of them
496 may benefit certain datasets, the improvements are
497 *marginal* compared with what temporal propensity
498 matching brings.

499 **Direct Matching (D).** In (6), we may directly
500 match the vectors of probabilities $(f(A, X))_{X \in \mathcal{X}}$.

501 **Temporality Pre-Filtering (F).** As the covariate
502 sampler and temporal predictor are two different
503 LMs, a sampled covariate might not be a pre-
504 ceding event judged by the temporal predictor. We
505 may filter the covariates before matching temporal
506 propensities such that $f(X, E_1) > f(X, E_2)$.

507 **Score Normalization (S).** In Section 4 we use
508 $s(E_1, E_2)$ for $f(E_1, E_2)$, we can also normalize it
509 and form $f(E_1, E_2)$ through

$$510 \quad f(E_1, E_2) = \frac{s(E_1, E_2)}{s(E_1, E_2) + s(E_2, E_1) + s(E_1, N) + s(N, E_1)} \quad (7)$$

511 where N represents the null event when no addi-
512 tional information is given, set as an empty string.

513 **Propensity Normalization (Q).** In Equa-
514 tion (6), we can also normalize the esti-
515 mates first before forming the q vectors via
516 $P(X(0)) = f(X, E_1) / \sum_{X' \in \mathcal{X}} f(X', E_1)$ and
517 $P(X(0), A(1)) = f(X, A) / \sum_{X' \in \mathcal{X}} f(X', A)$.

	Random Baseline	$\hat{\Delta}_1 \uparrow L_1$ -Balanced	$\hat{\Delta}_2 \uparrow L_2$ -Balanced	$\hat{\Delta}_{E_1} \uparrow$ Temporal	$\hat{\Delta}_A \uparrow$ Unbalanced	$\hat{\Delta}_x \uparrow$ Misspecified
COPA-DEV	0.5 ± 0.050	0.6900	0.7000	0.5800	0.5600	0.5300
COPA-TEST	0.5 ± 0.022	0.5640	0.5640	0.5200	0.5400	0.5240
GLUCOSE-D1	0.5 ± 0.040	0.6645	0.6968	0.5677	0.5742	0.6581
COPA-DEV (-T)	0.5 ± 0.050	0.6200	0.6300	0.5300	0.4800	0.5300
COPA-TEST (-T)	0.5 ± 0.022	0.5800	0.5740	0.4540	0.4600	0.4860
GLUCOSE-D1 (-T)	0.5 ± 0.040	0.6065	0.6194	0.5548	0.4387	0.3742

Table 1: **Best zero-shot results.** Shaded rows have temporal fine-tuning (T) disabled. (i) Estimators with temporal propensities balanced ($\hat{\Delta}_1$ and $\hat{\Delta}_2$) perform consistently better than the unbalanced and the temporal estimators. (ii) (2) In general, without temporality fine-tuning (“-T”, see Section 4), the performances degrade.

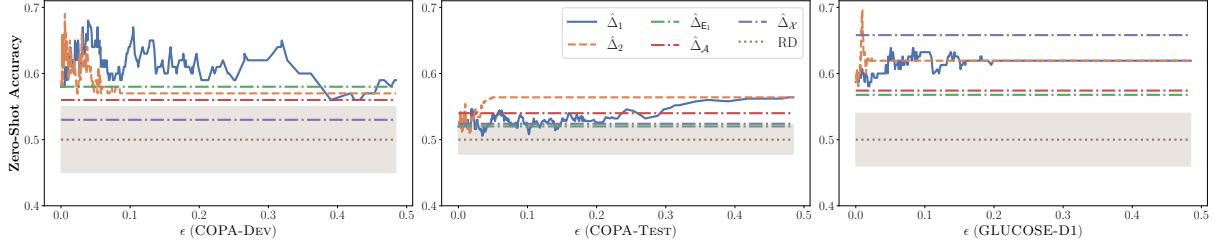


Figure 3: **Best zero-shot result vs ϵ .** Balanced estimators significantly outperform un-balanced and other variants for both COPA-DEV (left), COPA-TEST (middle) and GLUCOSE-D1 (right).

518
519
520
521
Co-occurrence Stabilization (C). The fine-tuned
522 temporal predictor may sometimes still fail to cover
523 the connectives. We can stabilize $\mathbb{P}(X \prec A)$ by
524 setting it to $(P(A(0), X(1)) + P(X(0), A(1)))/2$.

525
526
527
528
Estimand Normalization (E). We can normalize
529 the probability $\mathbb{P}(A \prec B)$ in the estimand Δ by
530 dividing $(P(A(0), B(1)) + P(B(0), A(1)))$.

531 5.3 Results

532 5.3.1 A Concrete Example

533 We first examine a particular example when the
vanilla temporal score $\hat{\Delta}_{E_1}$ fails but $\hat{\Delta}_1$ does not.

534 **Example 1** (Did $E_1^{(1)}$ or $E_1^{(2)}$ cause E_2 ?).

535 $E_1^{(1)}$: I was preparing to wash my hands.
536 $E_1^{(2)}$: I was preparing to clean the bathroom.
537 E_2 : I put rubber gloves on.
538 $A_{15}^{(1)}$: I was preparing to wash my feet.
539 $A_5^{(2)}$: Kevin was preparing to clean the bathroom.

540 This is the 63-rd instance in COPA-DEV
541 together a matched intervention (L_2 -balancing
542 with optimal ϵ) for each choice. The unad-
543 justed scores are $\hat{\Delta}_A(E_1^{(1)}, E_2) = 0.067894$
544 and $\hat{\Delta}_A(E_1^{(2)}, E_2) = 0.097539$ while the L_2 -
545 balanced scores are $\hat{\Delta}_2(E_1^{(1)}, E_2) = -0.010274$
546 and $\hat{\Delta}_2(E_1^{(2)}, E_2) = 0.001311$. The balanced score

547 selects the correct choice ($E_1^{(2)}$) with higher confi-
548 dence. More details are given in the Appendix.

549 5.3.2 Discussion

550 We show best zero-shot results over design choices
551 (and over ϵ) in Figure 3 and Table 1. As ROCK
552 tackles CCR from a completely new perspective,
553 there are no real baselines to compare with; our
554 goal is to demonstrate that *the causal inference*
555 *motivated method, temporal propensity matching,*
556 *mitigates spurious correlations* by comparing bal-
557 anced scores with unbalanced ones. We think this
558 perspective would also benefit the NLP community
559 at large for solving CCR and other related tasks.

560 **Comparison with existing methods.** Among un-
561 supervised baselines that we are able to reproduce,
562 the self-talk method (Shwartz et al., 2020) achieves
563 66% on COPA-DEV without external knowledge
564 and 69% when the CoMET-Net (Bosselut et al.,
565 2019) that contains commonsense knowledge is
566 used. Our L_2 -balanced score achieves 70% with-
567 out using any external commonsense knowledge.

568 **Temporal propensity matching is effective.** In
569 Table 1 (unshaded rows), we observe that balanced
570 scores have generally better performances on all
571 datasets compared with the temporal estimator and
572 the unadjusted estimator, implying that (i) tempo-
573 rality is important for CCR, yet they are susceptible

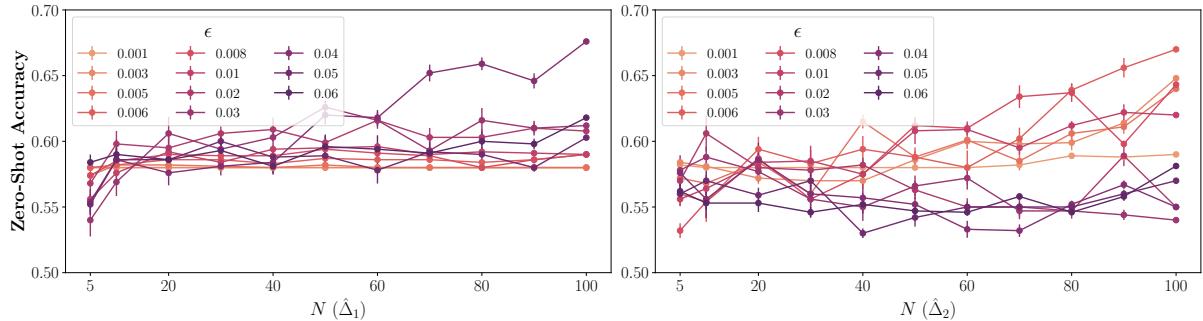


Figure 4: **Zero-shot result on COPA-DEV vs covariate set size $N = |\mathcal{X}|$ with 95%-confidence bands.** In general, using a larger N improves performances for both L_1 -balanced score ($\hat{\Delta}_1$, left) and L_2 -balanced score ($\hat{\Delta}_2$, right).

568 to spurious correlations; (ii) balancing covariates
569 via matching temporal propensities is effective.

570 **Rules-of-thumb for choosing ϵ .** The parameter ϵ
571 controls the threshold of covariates selection and p
572 controls its geometry (see e.g., [Hastie et al., 2015](#)).
573 Hinted by Figure 3, a general rule-of-thumb should
574 be $\epsilon < 0.1$. Table A.1 shows optimal ϵ values when
575 constrained to $[0, 0.1]$, where all are global optimal
576 except for COPA-TEST under L_1 -balanced score
577 (whose accuracy is 0.552). Hence we recommend
578 setting ϵ to be reasonably small ϵ such as within
579 $(0.01, 0.1)$ when $p = 1$ and relatively smaller such
580 as $(0.005, 0.05)$ when $p = 2$. The optimal value
581 depends on the implementation details of ROCK
582 components and domains of CCR to be performed,
583 yet these choices should result in a good start.

584 5.4 Ablation Studies

585 **Temporality Fine-Tuning.** Shaded rows in Ta-
586 ble 1 show that when we use the pretrained
587 RoBERTa-BASE without temporality fine-tuning
588 (we increase k to 30), almost all estimators do not
589 have decent performance. We conclude that (i)
590 pretrained LMs usually have poor ‘‘temporal aware-
591 ness,’’ and (ii) temporal fine-tuning helps LMs to
592 extract temporal knowledge essential to CCR.

593 **Covariate Set Size.** Figure 4 depicts zero-shot
594 results on COPA-TEST against the covariate set
595 size $N = |\mathcal{X}|$ together with 95%-confidence bands.
596 Here we only enable score normalizations (N)
597 among all six normalizations. We observe that
598 in general, increasing covariate set size improves
599 performances if ϵ is reasonable: if ϵ is too small,
600 added covariates may have little impacts while they
601 may introduce more noises if ϵ is too large.

602 **Normalizations.** In Section 5.2 we discussed six
603 possible normalizations. We report the best per-

	COPA-DEV		COPA-TEST		GLUCOSE-D1	
	$\hat{\Delta}_1 \uparrow$	$\hat{\Delta}_2 \uparrow$	$\hat{\Delta}_1 \uparrow$	$\hat{\Delta}_2 \uparrow$	$\hat{\Delta}_1 \uparrow$	$\hat{\Delta}_2 \uparrow$
Best	0.6900	0.7000	0.5640	0.5640	0.6645	0.6968
-S	0.01	0.06	-	-	0.08	0.11
-Q	0.01	-	-	-	0.03	-
-C	-	-	0.01	0.01	0.09	0.13
-E	0.01	0.01	-	-	0.03	-

Table 2: **Single-component ablations on normalizations.** Marked in red are percentage decreases compared with the best result (i.e., computed as $(a - b)/a$).

604 formance when each normalization is removed in
605 Table 2, where red marks the percentage decrease
606 compared with the best result (D and F not shown
607 as there is no change). Full ablations of all combi-
608 nations of normalizations and more discussions are
609 given in the Appendix. We observe that (i) certain
610 normalizations benefit certain datasets; (ii) in gen-
611 eral, improvements due to normalizations are only
612 *marginal*, so long as the framework is appropriate.

6 Discussions and Open Problems

613 We articulate the central question of CCR and in-
614 troduce ROCK, a novel framework for zero-shot
615 CCR based on classical causal inference principles.
616 ROCK sheds light on the CCR problem from new
617 perspectives that are arguably more well-founded
618 and demonstrates great potential for zero-shot CCR
619 as shown by empirical studies of various datasets.
620 There are several possible avenues for future works.
621 (i) **Prompt engineering** for better temporal predic-
622 tors and event sampler will likely benefit ROCK;
623 (ii) **Implicit events and reporting biases** in train-
624 ing data are likely to bias the LMs. How to account
625 for implicit events? (iii) **Contextual CCR** takes
626 the *contexts* into consideration. How to incorporate
627 contextual information for ROCK?

References

- Debarun Bhattacharjya, Tian Gao, and Dharmashankar Subramanian. 2020. [Order-dependent event models for agent interactions](#). In *Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence, IJCAI-20*, pages 1977–1983. International Joint Conferences on Artificial Intelligence Organization.
- Antoine Bosselut, Hannah Rashkin, Maarten Sap, Chaitanya Malaviya, Asli Celikyilmaz, and Yejin Choi. 2019. COMET: Commonsense Transformers for Automatic Knowledge Graph Construction. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics (ACL)*.
- Mario Bunge. 1979. *Causality and modern science*, 4 edition. Routledge.
- Du-Seong Chang and Key-Sun Choi. 2004. Causal relation extraction using cue phrase and lexical pair probabilities. In *IJCNLP*.
- William G Cochran and S Paul Chambers. 1965. The planning of observational studies of human populations. *Journal of the Royal Statistical Society. Series A (General)*, 128(2):234–266.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proc. of the Annual Conference of the North American Chapter of the Association for Computational Linguistics*.
- Quang Xuan Do, Yee Seng Chan, and Dan Roth. 2011. Minimally supervised event causality identification. In *Proceedings of the Conference on EMNLP*. Association for Computational Linguistics.
- Amir Feder, Katherine A. Keith, Emaad Manzoor, Reid Pryzant, Dhanya Sridhar, Zach Wood-Doughty, Jacob Eisenstein, Justin Grimmer, Roi Reichart, Margaret E. Roberts, Brandon M. Stewart, Victor Veitch, and Diyi Yang. 2021. [Causal inference in natural language processing: Estimation, prediction, interpretation and beyond](#).
- Ronald A Fisher. 1958. Cancer and smoking. *Nature*, 182(4635):596–596.
- Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason Phang, Horace He, Anish Thite, Noa Nabeshima, Shawn Presser, and Connor Leahy. 2020. The Pile: An 800gb dataset of diverse text for language modeling. *arXiv preprint arXiv:2101.00027*.
- Matt Gardner, Joel Grus, Mark Neumann, Oyyvind Tafjord, Pradeep Dasigi, Nelson F. Liu, Matthew Peters, Michael Schmitz, and Luke S. Zettlemoyer. 2017. [Allennlp: A deep semantic natural language processing platform](#).
- Tobias Gerstenberg, Noah D. Goodman, David A. Lagnado, and Joshua B. Tenenbaum. 2021. A counterfactual simulation model of causal judgments for physical events. *Psychological review*.
- Andrew Gordon, Zornitsa Kozareva, and Melissa Roemmele. 2012. [SemEval-2012 task 7: Choice of plausible alternatives: An evaluation of commonsense causal reasoning](#). In **SEM 2012: The First Joint Conference on Lexical and Computational Semantics – Volume 1: Proceedings of the main conference and the shared task, and Volume 2: Proceedings of the Sixth International Workshop on Semantic Evaluation (SemEval 2012)*, pages 394–398, Montréal, Canada. Association for Computational Linguistics.
- C. W. J. Granger. 1969. [Investigating causal relations by econometric models and cross-spectral methods](#). *Econometrica*, 37(3):424–438.
- Mingyue Han and Yinglin Wang. 2021. [Doing good or doing right? exploring the weakness of commonsense causal reasoning models](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*, pages 151–157, Online. Association for Computational Linguistics.
- Trevor Hastie, Robert Tibshirani, and Martin Wainwright. 2015. *Statistical Learning with Sparsity: The Lasso and Generalizations*. Chapman Hall.
- James J Heckman. 2005. Rejoinder: response to sobel. *Sociological Methodology*, 35(1):135–150.
- Austin Bradford Sir Hill. 1965. The environment and disease: Association or causation? *Journal of the Royal Society of Medicine*, 58:295 – 300.
- Paul W Holland. 1986. Statistics and causal inference. *Journal of the American statistical Association*, 81(396):945–960.
- Guido W Imbens and Donald B Rubin. 2015. *Causal inference in statistics, social, and biomedical sciences*. Cambridge University Press.
- Dongyeop Kang, Varun Gangal, Ang Lu, Zheng Chen, and Eduard Hovy. 2017. Detecting and explaining causes from text for a time series event. In *Conference on Empirical Methods on Natural Language Processing*.
- Pride Kavumba, Naoya Inoue, Benjamin Heinzerling, Keshav Singh, Paul Reisert, and Kentaro Inui. 2019. [When choosing plausible alternatives, clever hans can be clever](#). In *Proceedings of the First Workshop on Commonsense Inference in Natural Language Processing*, pages 33–42, Hong Kong, China. Association for Computational Linguistics.
- Katherine A Keith, David Jensen, and Brendan O’Connor. 2020. Text and causal inference: A review of using text to remove confounding from causal estimates. *arXiv preprint arXiv:2005.00649*.
- Benjamin Kuipers. 1984. Commonsense reasoning about causality: deriving behavior from structure. *Artificial intelligence*, 24(1-3):169–203.
- Zhiyi Luo, Yuchen Sha, Kenny Q Zhu, Seung-won Hwang, and Zhongyuan Wang. 2016. Commonsense causal reasoning between short texts. In *KR*, pages 421–431.
- John Stuart Mill. 1851. *A System of Logic, Ratiocinative and Inductive: Being a Connected View of the Principles of Evidence, and the Methods of Scientific Investigation*, volume 1 of *Cambridge Library Collection - Philosophy*. Cambridge University Press.
- Nasrin Mostafazadeh, Aditya Kalyanpur, Lori Moon, David Buchanan, Lauren Berkowitz, Or Biran, and Jennifer Chu-Carroll. 2020. [Glucose: Generalized and contextualized story explanations](#).

743	Jerzy S Neyman. 1923. On the application of probability theory to agricultural experiments. <i>Essay on principles. Section 9. Annals of Agricultural Sciences</i> , 10:1–51.	798
744		799
745		800
746	Qiang Ning, Zhili Feng, and Dan Roth. 2017. A Structured Learning Approach to Temporal Relation Extraction . In <i>Proc. of the Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 1038–1048, Copenhagen, Denmark. Association for Computational Linguistics.	801
747		802
748		
749		
750		
751		
752	Qiang Ning, Zhili Feng, Hao Wu, and Dan Roth. 2019a. Joint reasoning for temporal and causal relations. <i>arXiv preprint arXiv:1906.04941</i> .	803
753		804
754		805
755	Qiang Ning, Sanjay Subramanian, and Dan Roth. 2019b. An Improved Neural Baseline for Temporal Relation Extraction . In <i>Proc. of the Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> .	806
756		807
757		
758		
759	Qiang Ning, Hao Wu, Haoruo Peng, and Dan Roth. 2018. Improving temporal relation extraction with a globally acquired statistical resource . In <i>Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)</i> , pages 841–851, New Orleans, Louisiana. Association for Computational Linguistics.	810
760		811
761		
762		
763		
764		
765		
766		
767	Timothy J. O’Gorman, Kristin Wright-Bettner, and Martha Palmer. 2016. Richer event description: Integrating event coreference with temporal, causal and bridging annotation.	812
768		
769		
770	Judea Pearl. 1995. Causal diagrams for empirical research . <i>Biometrika</i> , 82(4):669–688.	813
771		814
772		
773	Judea Pearl and Dana Mackenzie. 2018. <i>The book of why: the new science of cause and effect</i> . Basic Books.	815
774		816
775		
776	Jonas Peters, Dominik Janzing, and Bernhard Schlkopf. 2017. <i>Elements of Causal Inference: Foundations and Learning Algorithms</i> . The MIT Press.	817
777		
778		
779	Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners.	818
780		819
781		
782	Hannah Rashkin, Maarten Sap, Emily Allaway, Noah A Smith, and Yejin Choi. 2018. Event2mind: Commonsense inference on events, intents, and reactions. <i>arXiv preprint arXiv:1805.06939</i> .	820
783		821
784	Paul R Rosenbaum. 1984. The consequences of adjustment for a concomitant variable that has been affected by the treatment. <i>Journal of the Royal Statistical Society: Series A (General)</i> , 147(5):656–666.	822
785		823
786		
787		
788	Paul R Rosenbaum. 1989. Optimal matching for observational studies. <i>Journal of the American Statistical Association</i> , 84(408):1024–1032.	824
789		825
790		
791	Paul R Rosenbaum. 2002. <i>Observational Studies</i> . Springer.	826
792		
793		
794	Paul R Rosenbaum and Donald B Rubin. 1983. The central role of the propensity score in observational studies for causal effects. <i>Biometrika</i> , 70(1):41–55.	827
795		828
796		
797	Dan Roth. 2017. Incidental Supervision: Moving beyond Supervised Learning . In <i>Proc. of the Conference on Artificial Intelligence (AAAI)</i> .	829
		830
	Donald B Rubin. 1974. Estimating causal effects of treatments in randomized and nonrandomized studies. <i>Journal of Educational Psychology</i> , 66(5):688.	831
		832
	Donald B Rubin. 1980. Bias reduction using mahalanobis-metric matching. <i>Biometrics</i> , pages 293–298.	833
		834
	Donald B Rubin. 2005. Causal inference using potential outcomes: Design, modeling, decisions. <i>Journal of the American Statistical Association</i> , 100(469):322–331.	835
	Bertrand Russell. 1912. On the notion of cause . <i>Proceedings of the Aristotelian Society</i> , 13:1–26.	836
		837
	E. Sandhaus. 2008. The New York Times Annotated Corpus. <i>Linguistic Data Consortium, Philadelphia</i> .	838
		839
	Maarten Sap, Hannah Rashkin, Derek Chen, Ronan Le Bras, and Yejin Choi. 2019. Social iqa: Commonsense reasoning about social interactions. In <i>EMNLP 2019</i> .	840
		841
	Maarten Sap, Vered Shwartz, Antoine Bosselut, Yejin Choi, and Dan Roth. 2020. Commonsense reasoning for natural language processing. In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: Tutorial Abstracts</i> , pages 27–33.	842
		843
	Vered Shwartz, Peter West, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. 2020. Unsupervised commonsense question answering with self-talk . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 4615–4629. Association for Computational Linguistics.	844
		845
	Ieva Staliunaite, Philip John Gorinski, and Ignacio Iacobacci. 2021. Improving commonsense causal reasoning by adversarial training and data augmentation. In <i>AAAI</i> .	846
		847
	Alexandre Tamborrino, Nicola Pellicanò, Baptiste Pannier, Pascal Voitot, and Louise Naudin. 2020. Pre-training is (almost) all you need: An application to commonsense reasoning. <i>ArXiv</i> , abs/2004.14074.	848
		849
	Siddharth Vashishta, Adam Poliak, Yash Kumar Lal, Benjamin Van Durme, and Aaron Steven White. 2020. Temporal reasoning in natural language inference. In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: Findings</i> , pages 4070–4078.	850
		851
	Ben Wang and Aran Komatsuzaki. 2021. GPT-J-6B: A 6 Billion Parameter Autoregressive Language Model . https://github.com/kingoflolz/mesh-transformer-jax .	852
		853
	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierrick Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. 2020. Transformers: State-of-the-art natural language processing . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations</i> , pages 38–45, Online. Association for Computational Linguistics.	854
		855
	Tongshuang Wu, Marco Tulio Ribeiro, Jeffrey Heer, and Daniel Weld. 2021. Polyjuice: Generating counterfactuals for explaining, evaluating, and improving models . In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 6707–6723, Online. Association for Computational Linguistics.	856
		857

858 Hongming Zhang, Yintong Huo, Xinran Zhao, Yangqiu Song,
859 and Dan Roth. 2021. Learning contextual causality be-
860 tween daily events from time-consecutive images. In *Pro-*
861 *ceedings of the IEEE/CVF Conference on Computer Vision*
862 *and Pattern Recognition (CVPR) Workshops*, pages 1752–
863 1755.

864 Ben Zhou, Qiang Ning, Daniel Khashabi, and Dan Roth. 2020.
865 Temporal common sense acquisition with minimal supervi-
866 sion. *arXiv preprint arXiv:2005.04304*.

A Additional Experiment Details

A.1 Rule-of-Thumb for Choosing ϵ

In Table A.1 we show the best ϵ values when constrained in $\epsilon \in [0, 0.1]$. Hence we recommend setting ϵ to be reasonably small ϵ such as within $(0.01, 0.1)$ when $p = 1$ and relatively smaller such as $(0.005, 0.05)$ when $p = 2$. The optimal value depends on the implementation details of ROCK components and domains of CCR to be performed, yet these choices should result in a good start.

A.2 Further Discussions on Temporality Fine-Tuning

In Figure 3, we observe that, counterintuitively, without temporality fine-tuning, the best performances of balanced estimators (0.58) are higher than those with temporality fine-tuning (0.564). Although this gap is within one standard deviation of the random baseline (0.022) thus no statistically significant conclusions can be drawn, but it might hint that pretrained LMs may have already been very aware of temporality. Is this really the case? A closer look at the full ablation table to be introduced shortly in Table A.5 reveals that the stellar performance is attributed to one particular normalization, estimand normalization (**E**), which was actually detrimental to another dataset (GLUCOSE-D1). Hence we think this normalization may favor certain dataset over others, thus we think it is not recommendable to include this normalization when dealing with a new dataset.

A.3 Full Ablation on Normalizations

Recall in Section 5.4 we discussed six possible normalizations that may stabilize the estimation procedure:

(D) Direct Matching: in (6), instead of forming the temporal propensity vectors q using conditional probabilities, we may directly match the vectors of probabilities $(f(A, X))_{X \in \mathcal{X}}$. This normalization is not well motivated but might be easier to compute under certain circumstances, hence we include it as a comparison.

(F) Temporality Pre-Filtering: as the covariate sampler and temporal predictor are two different LMs, a sampled covariate might not be a preceding event judged by the temporal predictor. Thus, we can filter the covariates \mathcal{X} before matching temporal propensities such

that we only keep covariates $X \in \mathcal{X}$ satisfying $f(X, E_1) > f(S, E_1)$.

(S) Score Normalization: in Section 4 we use $s(E_1, E_2)$ for $f(E_1, E_2)$. We can also normalize it and form $f(E_1, E_2)$ through

$$f(E_1, E_2) = \frac{s(E_1, E_2)}{s(E_1, E_2) + s(E_2, E_1) + s(E_1, N) + s(N, E_1)} \quad (\text{A.1})$$

where N represents the null event when no additional information is given, set as an empty string. In practice, this normalization does not differ much from the normalization

$$f(E_1, E_2) = \frac{s(E_1, E_2)}{s(E_1, E_2) + s(E_2, E_1)}, \quad (\text{A.2})$$

which does not involve N . However, using N has the benefit of stabilizing the estimate $f(\cdot, \cdot)$ as in rare scenarios $s(E_1, E_2)$ and $s(E_2, E_1)$ may both close to zero.

(Q) Propensity Normalization: in Equation (6), we can also normalize the estimates first before forming the q vectors via

$$\begin{aligned} P(X(0)) &= \frac{f(X, E_1)}{\sum_{X' \in \mathcal{X}} f(X', E_1)}, \\ P(X(0), A(1)) &= \frac{f(X, A)}{\sum_{X' \in \mathcal{X}} f(X', A)}, \end{aligned} \quad (\text{A.3})$$

where we estimate $P(X(0))$ as the relative frequency of $X(0)$ among all possible events in \mathcal{X} ; and $P(X(0), A(1))$ among all possible (X, A) pairs.

(C) Co-Occurrence Stabilization: on rare occasions, the fine-tuned temporal predictor may sometimes still fail to cover the connectives. We can stabilize $\mathbb{P}(X \prec A)$ by setting it to $(P(A(0), X(1)) + P(X(0), A(1))) / 2$. This in effect results in an alternative estimand based on co-occurrences of events (instead of precedence) and can be viewed as a weaker causation in CCR.

(E) Estimand Normalization: the score normalization (N) takes place at temporal propensity matching. We can normalize the temporal probability $\mathbb{P}(A \prec B)$ in the estimand Δ by dividing $(P(A(0), B(1)) + P(B(0), A(1)))$, thus setting

$$\mathbb{P}(A \prec B) = \frac{P(A(0), B(1))}{P(A(0), B(1)) + P(B(0), A(1))}. \quad (\text{A.4})$$

	COPA-DEV		COPA-TEST		GLUCOSE-D1	
	$\hat{\Delta}_1$	$\hat{\Delta}_2$	$\hat{\Delta}_1$	$\hat{\Delta}_2$	$\hat{\Delta}_1$	$\hat{\Delta}_2$
ϵ^*	0.043067	0.006029	0.059232	0.048837	0.046643	0.009374

Table A.1: Best choices of ϵ when $\epsilon < 0.1$.

	$\hat{\Delta}_1$	$\hat{\Delta}_2$	$\hat{\Delta}_{E_1}$	$\hat{\Delta}_{\mathcal{A}}$	$\hat{\Delta}_{\mathcal{X}}$
$(E_1, E_2^{(1)})$	-0.002	-0.002	0.106	0.002	0.106
$(E_1, E_2^{(2)})$	-0.001	-0.001	0.086	-0.012	0.086

Table A.2: Scores for Example A.1.

A.3.1 Ablation Results

We report ablations on all possible subset of normalizations together with temporality fine-tuning (-T, see Section 4 in Table A.5. Note that when \mathbf{D} is enabled, \mathbf{S} and \mathbf{Q} are not active and when \mathbf{C} is enabled, \mathbf{E} is not active, thus resulting in a total of $2^2(2^2 + 1)(2^1 + 1) = 30$ combinations

Ablations resulting in the best performances are highlighted in blue and those resulting in the worst the performances are highlighted in red. Shaded rows are results without temporal fine-tuning (using top $k = 30$ tokens in mask language modeling). We summarize our observations as follows.

Improvements due to normalizations are marginal. The gap between best and worst performance are marginal, except for the GLUCOSE-D1 dataset, which is mainly caused by enabling estimand normalization (\mathbf{E}). Without considering \mathbf{E} , the worst result is 0.594 (+Q or +FQ). Furthermore, we note the gap between the best results and the results under no normalizations (\emptyset) is also marginal, indicating that for CCR it is more important to have a well-established baseline and temporal signal extractors than exploring different normalizations.

Furthermore, the outliers are interesting: enabling estimand normalization (\mathbf{E}) has little or no effects on most datasets but can boost the performance on COPA-TEST under non fine-tuned temporal predictors (-T) while is detrimental to GLUCOSE-D1 under fine-tuned temporal predictors.

Rules-of-thumb for choosing normalizations. As a general rule-of-thumb, temporal score nor-

	$\hat{\Delta}_1$	$\hat{\Delta}_2$	$\hat{\Delta}_{E_1}$	$\hat{\Delta}_{\mathcal{A}}$	$\hat{\Delta}_{\mathcal{X}}$
$(E_1^{(1)}, E_2)$	0.056	-0.001	0.109	0.096	0.109
$(E_1^{(2)}, E_2)$	0.005	-0.010	0.279	0.118	0.279

Table A.4: Scores for Example A.3.

malization (\mathbf{S}) should be enabled and the q vectors should be properly formed (without direct matching \mathbf{D}); temporal pre-filtering (\mathbf{F}) and propensity normalization (\mathbf{Q}) in general do not affect the results significantly; co-occurrence stabilization (\mathbf{C}) has greater positive effect on datasets when a weaker notion of causation are desirable (e.g., GLUCOSE-D1 we constructed); while estimand normalization (\mathbf{E}) improves certain datasets (e.g., COPA-TEST without temporal fine-tuning), it has detrimental effects on some others (e.g., GLUCOSE-D1 with temporal fine-tuning), hence we recommend disabling it by default.

A.4 Full Examples

We also attach three full examples from our implementation of the ROCK. The problem instances are given below. For each instance, we tabulate 50 covariates sampled, all interventions generated, the corresponding $\|q(\mathbf{x}; \mathbf{A}) - q(\mathbf{x}; E_1)\|_p$, and the temporal probabilities $\mathbb{P}(\cdot \prec E_2)$.

Example A.1 (Did E_1 cause $E_2^{(1)}$ or $E_2^{(2)}$?).

E_1 : The teacher assigned homework to the students. 1006

$E_2^{(1)}$: The students passed notes. 1007

$E_2^{(2)}$: The students groaned. 1008

1009

This is the 72-nd instance of COPA-DEV, the full tables for inferring the causation from E_1 to $E_2^{(1)}$ and E_1 to $E_2^{(2)}$ are given in Table A.6 and Table A.7 respectively. Different scores are shown in Table A.2. Note that this example is not easy:

Example A.2 (Did $E_1^{(1)}$ or $E_1^{(2)}$ cause E_2 ?).

$E_1^{(1)}$: I was preparing to wash my hands. 1016

$E_1^{(2)}$: I was preparing to clean the bathroom. 1017

E_2 : I put rubber gloves on. 1018

1019

This is the 63-nd instance of COPA-DEV, the full tables for inferring the causation from $E_1^{(1)}$ to E_2 and $E_1^{(1)}$ to $E_2^{(2)}$ are given in Table A.8 and

	$\hat{\Delta}_1$	$\hat{\Delta}_2$	$\hat{\Delta}_{E_1}$	$\hat{\Delta}_{\mathcal{A}}$	$\hat{\Delta}_{\mathcal{X}}$
$(E_1^{(1)}, E_2)$	-0.010	-0.010	0.068	0.036	0.068
$(E_1^{(2)}, E_2)$	0.002	0.001	0.098	0.035	0.098

Table A.3: Scores for Example A.2.

1024 Table A.9 respectively. Different scores are shown
1025 in Table A.3.

Example A.3 (Did $E_1^{(1)}$ or $E_1^{(2)}$ cause E_2 ?).

1026 $E_1^{(1)}$: His pocket was filled with coins.

1027 $E_1^{(2)}$: He sewed the hole in his pocket.

1028 E_2 : The man's pocket jingled as he walked.

1029
1030
1031 This is the 79-th instance of COPA-DEV, the
1032 full tables for inferring the causation from $E_1^{(1)}$
1033 to E_2 and $E_1^{(1)}$ to $E_2^{(2)}$ are given in Table A.10 and
1034 Table A.11 respectively. Different scores are shown
1035 in Table A.2.

Dataset	Score	Best	Worst	\emptyset	+D	+F	+S	+Q	+C	+E	+DF	+DC	+DE	+FS	+FQ	+FC	+FE	+SQ	+SC	+SE	+QC	+OF	+DFC	+DFQ	+FSQ	+FSC	+FSOC	+FQE	+SOC	+SQE	+FSQC	+PSQE			
COPA-DEV	$\Delta_1 \uparrow$	0.690	0.620	0.670	0.600	0.670	0.650	0.650	0.650	0.650	0.650	0.650	0.650	0.650	0.650	0.650	0.650	0.650	0.650	0.650	0.660	0.670	0.670	0.670	0.670	0.670	0.670	0.670	0.670	0.670	0.670				
	$\Delta_2 \uparrow$	0.700	0.630	0.630	0.630	0.650	0.630	0.630	0.630	0.630	0.630	0.630	0.630	0.630	0.630	0.630	0.630	0.630	0.630	0.630	0.660	0.670	0.670	0.670	0.670	0.670	0.670	0.670	0.670	0.670	0.670				
COPA-TEST	$\Delta_1 \uparrow$	0.564	0.528	0.542	0.548	0.542	0.540	0.548	0.564	0.554	0.540	0.548	0.564	0.554	0.540	0.548	0.564	0.546	0.540	0.544	0.564	0.546	0.540	0.544	0.542	0.540	0.544	0.542	0.540	0.544	0.542	0.540	0.544	0.542	
	$\Delta_2 \uparrow$	0.564	0.526	0.534	0.532	0.534	0.532	0.534	0.534	0.534	0.534	0.534	0.534	0.534	0.534	0.534	0.534	0.534	0.534	0.534	0.539	0.538	0.538	0.538	0.538	0.538	0.538	0.538	0.538	0.538	0.538				
GLUCOSE-D1	$\Delta_1 \uparrow$	0.665	0.503	0.600	0.594	0.594	0.594	0.594	0.503	0.606	0.639	0.503	0.606	0.594	0.594	0.594	0.503	0.606	0.594	0.613	0.639	0.510	0.613	0.613	0.613	0.613	0.613	0.613	0.613	0.613	0.613	0.613			
	$\Delta_2 \uparrow$	0.697	0.503	0.594	0.600	0.594	0.600	0.594	0.600	0.600	0.639	0.503	0.606	0.594	0.594	0.594	0.503	0.606	0.594	0.613	0.639	0.510	0.613	0.613	0.613	0.613	0.613	0.613	0.613	0.613	0.613	0.613			
COPA-DEV (-T)	$\Delta_1 \uparrow$	0.620	0.550	0.590	0.550	0.580	0.580	0.580	0.570	0.620	0.550	0.560	0.610	0.580	0.570	0.570	0.620	0.580	0.560	0.560	0.620	0.560	0.610	0.610	0.610	0.610	0.610	0.610	0.610	0.610	0.610	0.610	0.610		
	$\Delta_2 \uparrow$	0.630	0.530	0.610	0.530	0.610	0.600	0.580	0.600	0.620	0.530	0.550	0.600	0.580	0.600	0.600	0.620	0.580	0.560	0.560	0.620	0.580	0.610	0.610	0.610	0.610	0.610	0.610	0.610	0.610	0.610	0.610	0.610		
COPA-TEST (-T)	$\Delta_1 \uparrow$	0.574	0.580	0.484	0.494	0.486	0.484	0.484	0.484	0.484	0.486	0.486	0.486	0.486	0.486	0.484	0.484	0.484	0.484	0.484	0.484	0.484	0.484	0.484	0.484	0.484	0.484	0.484	0.484	0.484	0.484	0.484	0.484	0.484	0.484
	$\Delta_2 \uparrow$	0.606	0.510	0.568	0.555	0.568	0.568	0.568	0.568	0.568	0.568	0.568	0.568	0.568	0.568	0.568	0.568	0.568	0.568	0.568	0.561	0.619	0.555	0.561	0.561	0.619	0.555	0.561	0.561	0.561	0.561	0.561			
GLUCOSE-D1 (-T)	$\Delta_1 \uparrow$	0.619	0.503	0.568	0.555	0.568	0.568	0.568	0.568	0.568	0.568	0.568	0.568	0.568	0.568	0.568	0.568	0.568	0.568	0.568	0.587	0.587	0.587	0.587	0.587	0.587	0.587	0.587	0.587	0.587	0.587				
	$\Delta_2 \uparrow$	0.619	0.503	0.568	0.555	0.568	0.568	0.568	0.568	0.568	0.568	0.568	0.568	0.568	0.568	0.568	0.568	0.568	0.568	0.568	0.581	0.581	0.581	0.581	0.581	0.581	0.581	0.581	0.581	0.581	0.581				

Table A.5: Full ablation studies on normalizations. Ablations resulting in the best performances are highlighted in blue and those resulting in the worst performances are highlighted in red. Shaded rows are results without temporal fine-tuning (using top $k = 30$ tokens in mask language modeling). (i) The gap between best and worst performance are marginal, except for the GLUCOSE-D1 dataset, which is mainly caused by enabling estimand normalization (**E**). Without considering **E**, the worst result is 0.594 (**+Q** or **+FQ**). (ii) In general, temporal fine-tuning helps. The only exception on COPA-TEST is caused by the estimand normalization (**E**). (iii) As a general rule-of-thumb, it does not hurt to start with no normalizations enabled.

Sampled Covariates \mathcal{X}	$\ q(\mathbf{x}; \mathbf{A}) - q(\mathbf{x}; \mathbf{E}_1)\ _p$	\mathbf{E}_1 and Interventions \mathcal{A}
X_1 : He had written a brief book summary of the book and, using a set of questions.	0	
X_2 : There was homework help, help desk, and online support.	0.035	
X_3 : No one did the work on line, and no one received good grades for it.	0.0508	
X_4 : The kids had to do their school homework online.	0.084	
X_5 : This was the norm.	0.0279	
X_6 : The class would sit quietly and listen to their teacher talk.	0.1053	
X_7 : It was free time.	0.1291	
X_8 : Homework was only assigned when the teacher had a class with a lot of work for the students to.	0.1535	A_1 : The teacher assigned homework to the students.
X_9 : He did not give homework to his students.	0.1591	A_2 : The professor assigned homework to the students.
X_{10} : There was a long period of time when nobody ever did any homework.	0.1595	A_3 : The tourists van or the teacher assigned homework to the students.
X_{11} : He had been teaching them during class for weeks.	0.1635	A_4 : The teacher took homework to the students.
X_{12} : It was just a fun afternoon with the kids, and then it turned into a time of dr.	0.1999	A_5 : The teacher was assigning Justin with the homework to the students.
X_{13} : Every night, a student would get to work on their homework.	0.2011	A_6 : The teacher replaced the carpet for the library last night because the carpet was old homework to the students.
X_{14} : The students had to listen to music and watch a video, respectively, before they could do that.	0.2011	A_7 : The teacher assigned tests to the students.
X_{15} : The students had to do the homework themselves.	0.2011	A_8 : No one was assigned homework to the students.
X_{16} : The children used to go to school in the morning and study their books until the evening.	0.2011	A_9 : Unless the senior performed, the teacher assigned homework to the students.
X_{17} : The teacher assigned homework to the entire class.	0.2011	A_{10} : While Leong on the other hand assigned homework to the students.
X_{18} : Each student was given a piece of paper with some number on it.	0.1524	A_{11} : The teacher didn't give homework to the students.
X_{19} : They could play without limits.	0.1349	A_{12} : The teacher didn't assign homework to the students.
X_{20} : They thought that the homework was just a part of my study in each class.	0.1099	A_{13} : The teacher didn't tell anyone homework to the students.
X_{21} : Students who have not completed their homework will not be allowed to go to the next class.	0.0468	A_{14} : The teacher assigned nothing to the students.
X_{22} : The assignment was simple, they were just to read the assigned reading.	0.0485	A_{15} : The teacher assigned no children to the students.
X_{23} : However, he handed out the following set of questions, which the teacher posed one by one to.	0.0362	A_{16} : The teacher assigned no class to the students.
X_{24} : Students were not given much homework.	0.0315	A_{17} : The student assigned homework to the students.
X_{25} : Students were only encouraged to work on assignments and were not explicitly told to do extra.	0.0188	A_{18} : The professor assigned homework to the students.
X_{26} : Only the school's teacher did so.	0.0512	A_{19} : The teacher wrote on the algebraic homework to the students.
X_{27} : He would just talk to them or read articles or give his own opinion on the subject.	0.0515	A_{20} : The teacher read homework to the students.
X_{28} : There was no homework.	0.0201	A_{21} : The teacher assigned tests to the students.
X_{29} : The students had to read the textbook and test their knowledge of the material.	0.0647	A_{22} : The teacher assigned to the classroom stopped to the students.
X_{30} : They were assigned to do some homework.	0.0298	A_{23} : The teacher assigned anger to the students.
X_{31} : Teachers would typically assign the work to the students, but this teacher assigned it to the students and.	0.0301	
X_{32} : There were no homework assignments at all.	0.0301	
X_{33} : The students would go to the internet and download games.	0.0301	
X_{34} : The assignment had already been completed.	0.0301	
X_{35} : He asked his students on the first day of class to write down on A4 paper any questions.	0.0301	
X_{36} : No homework was assigned.	0.0301	
X_{37} : The students were all in the classroom, sitting in rows like the soldiers in the First World War.	0.0301	
X_{38} : I just gave them a paper with one page written on it.	0.0301	
X_{39} : There was no homework.	0.0301	
X_{40} : The students were told the homework, and the students were to do the homework on their own.	0.0301	

Table A.6: Example 1a: the first plausible pair of the 72-th instance in COPA-DEV, matched interventions are highlighted. Here E_1 : The teacher assigned homework to the students . and E_2 : The students passed notes .

Sampled Covariates \mathcal{X}	$\ q(\mathbf{x}; \mathcal{A}) - q(\mathbf{x}; \mathcal{E}_1)\ _p$	\mathcal{E}_1 and Interventions \mathcal{A}	$\mathbb{P}(\cdot \prec \mathcal{E}_2)$
X_{1i} : He had written a brief book summary of the book and, using a set of questions.	0		
X_{2i} : There was homework help, help desk, and online support.	0.0135		0.5308
X_{3i} : No one did the work on time, and no one received good grades for it.	0.0508		0.5263
X_{4i} : The kids had to do their school homework online.	0.0894		0.5207
X_{5i} : This was the norm.	0.0279		0.5260
X_{6i} : The class would sit quietly and listen to their teacher talk.	0.1053		0.5340
X_{7i} : It was free time.	0.1291		0.5396
X_{8i} : Homework was only assigned when the teacher had a class with a lot of work for the students to.	0.0591	E _{1i} : The teacher assigned homework to the students. A _{1i} : The professor assigned homework to the students.	0.5015
X_{9i} : He did not give homework to his students.	0.0365	A _{2i} : The professor supported that tourists assigned homework to the students.	0.5249
X_{10i} : There was a long period of time when nobody ever did any homework.	0.0201	A _{3i} : The tourists ran, or the teacher assigned homework to the students.	0.5249
X_{11i} : She had been teaching them during class for weeks.	0.0870	A _{4i} : The teacher took homework to the students.	0.5249
X_{12i} : It was just a fun afternoon with the kids, and then it turned into a time of dr.	0.0820	A _{5i} : The teacher was assigning Justin with the homework to the students.	0.5249
X_{13i} : Every night, a student would get to work on their homework.	0.0485	A _{6i} : The teacher replaced the car for the library last night because the carpet was old homework to the students.	0.5249
X_{14i} : The students had to listen to music and watch a video, respectively, before they could do their.	0.0485	A _{7i} : The teacher assigned tests to the students.	0.5249
X_{15i} : The students had to do the homework themselves.	0.1521	A _{8i} : No one was assigned homework to the students.	0.5249
X_{16i} : The children used to go to school in the morning and study their books until the evening.	0.0820	A _{9i} : Unless the senator performed, the teacher assigned homework to the students.	0.5249
X_{17i} : The teacher assigned homework to the entire class.	0.0524	A _{10i} : Noelle along on the other hand assigned homework to the students.	0.5249
X_{18i} : Each student was given a piece of paper with some number on it.	0.1349	A _{11i} : The teacher didn't give homework to the students.	0.5249
X_{19i} : They could play without limits.	0.1999	A _{12i} : The teacher didn't assign homework to the students.	0.5249
X_{20i} : I thought that the homework was just a part of my study in each class.	0.0468	A _{13i} : The teacher didn't tell anyone homework to the students.	0.5249
X_{21i} : Students who have not completed their homework will not be allowed to go to the next class.	0.0485	A _{14i} : The teacher assigned nothing to the students.	0.5249
X_{22i} : The assignment was simple, they were just to read the assigned reading.	0.0362	A _{15i} : The teacher assigned no children to the students.	0.5249
X_{23i} : However, he handed out the following set of questions, which the teacher posed one by one to.	0.0301	A _{16i} : The teacher assigned no class to the students.	0.5249
X_{24i} : Students were not given much homework.	0.0135	A _{17i} : The student assigned homework to the students.	0.5249
X_{25i} : Students were only encouraged to work on assignments and were not explicitly told to do extra.	0.0488	A _{18i} : The professor assigned homework to the students.	0.5249
X_{26i} : Only the school's teacher did so.	0.0512	A _{19i} : The teacher worked on the algebraic homework to the students.	0.5249
X_{27i} : He would just talk to them or read articles or give his own opinion on the subject.	0.0515	A _{20i} : The teacher wrote homework to the students.	0.5249
X_{28i} : There was no homework.	0.0647	A _{21i} : The teacher read homework to the students.	0.5249
X_{29i} : The students had to read the textbook and test their knowledge of the material.	0.0298	A _{22i} : The teacher assigned tests to the students.	0.5249
X_{30i} : Teachers would typically assign the work to the students, but this teacher assigned it to the students and.		A _{23i} : The teacher assigned anger to the students.	0.5249
X_{31i} : There were no homework assignments at all.			
X_{32i} : The students would go to the Internet and download games.			
X_{33i} : The assignment had already been completed.			
X_{34i} : He asked his students on the first day of class to write down on A4 paper any questions.			
X_{35i} : No homework was assigned.			
X_{36i} : The students were all in the classroom, sitting in rows like the soldiers in the First World War.			
X_{37i} : I just gave them a paper with one page written on it.			
X_{38i} : There was no homework.			
X_{39i} : The students were told the homework, and the students were to do the homework on their own.			

Table A.7: Example 1b: the second plausible pair of the 72-th instance in COPA-DEV, matched interventions are highlighted. Here E_1 : The teacher assigned homework to the students . and E_2 : The students groaned.

Sampled Covariates \mathcal{X}	$\ q(\mathbf{x}; \mathbf{A}) - q(\mathbf{x}; \mathbf{E}_1)\ _p$	\mathbf{E}_1 and Interventions \mathcal{A}	$\mathbb{M}^{\mathbf{C}} \prec \mathbf{E}_2$
$X_{\mathbf{i}}: I$ had scrubbed my face, arms, and chest; using a baby shampoo called "San.	0		
$X_{\mathbf{i}}: I$ was preparing to wash my hands.	0.2485		
$X_{\mathbf{i}}: I$ was running low and got my hands wet but not the shoes because the hands were preparing to wash my hands.	0.2485		
$X_{\mathbf{i}}: I$ was standing close to the sink because the sink was well lit preparing to wash my hands.	0.1792		
$X_{\mathbf{i}}: I$ wanted to get rid of the smell of bleach and use water instead because the water was clean and preparing to wash my hands.	0.2153		
$X_{\mathbf{i}}: I$ had put on a new pair of latex gloves; I'm very careful about hand cleaning.	0.2153		
$X_{\mathbf{i}}: I$ wanted to take out my medicine and check all my symptoms.	0.0752		
$X_{\mathbf{i}}: I$ had been brushing the sand from my clothes.	0.1014		
$X_{\mathbf{i}}: I$ had put a couple of paper towels in their drawer by the sink.	0.1210		
$X_{\mathbf{i}}: I$ had been sitting in the armchair in their living room.	0.0677		
$X_{\mathbf{i}}: I$ had decided to make a cup of tea.	0.124		
$X_{\mathbf{i}}: I$ had been playing with my son, watching and old video on YouTube, and I.	0.3054		
$X_{\mathbf{i}}: I$ had been brushing the sand from my clothes.	0.070		
$X_{\mathbf{i}}: I$ went through the washing ceremony to check the level of purity in my body, I washed my face.	0.0000		
$X_{\mathbf{i}}: I$ washed my face.	0.0586		
$X_{\mathbf{i}}: I$ had just finished eating my breakfast.	0.0912		
$X_{\mathbf{i}}: I$ prepared a simple salad and some rolls on the table.	0.1014		
$X_{\mathbf{i}}: I$ scrubbed my hands with a little bit of soap.	0.057		
$X_{\mathbf{i}}: I$ turned to the side of the mirror, and I had a look.	0.0308		
$X_{\mathbf{i}}: I$ always take my shoes off.	0.0324		
$X_{\mathbf{i}}: I$ However, I removed some leftover food from the table, where the two men had been eating.	0.1373		
$X_{\mathbf{i}}: I$ took a few deep breaths and had a conversation with my heart.	0.1263		
$X_{\mathbf{i}}: I$ had just changed into the outfit I was wearing, a pretty, pale pink T-shirt and.	0.1601		
$X_{\mathbf{i}}: I$ was clearing away breakfast things.	0.3710		
$X_{\mathbf{i}}: I$ used to take a shower, and now it was time to do that again.	0.3322		
$X_{\mathbf{i}}: I$ I washed my face.	0.1219		
$X_{\mathbf{i}}: I$ needed to check my phone.	0.1294		
$X_{\mathbf{i}}: I$ I took the time to put on another pair of socks, and the socks for that matter.	0.4810		
$X_{\mathbf{i}}: I$ washed my hands more than a thousand times.	0.0910		
$X_{\mathbf{i}}: I$ was putting on a gown and cap, and to check the medications I had received for.	0.0700		
$X_{\mathbf{i}}: I$ had been sitting at my desk, answering emails and making phone calls.	0.2772		
$X_{\mathbf{i}}: I$ I'd touched the wall for some reason.	0.2068		
$X_{\mathbf{i}}: I$ used to dry them properly.	0.0970		
$X_{\mathbf{i}}: I$ I made sure my hands were clean.	0.0000		
$X_{\mathbf{i}}: I$ as a last resort, I would always scrub the top of my hands with a nail brush to.	0.1389		
$X_{\mathbf{i}}: I$ had looked into their bathroom mirror.	0.034		
$X_{\mathbf{i}}: I$ was talking to him.	0.0324		
$X_{\mathbf{i}}: I$ I was wiping my pants on the sides of shorts and shirt, like a dirty secret.	0.2359		
$X_{\mathbf{i}}: I$ had been holding a glass of orange juice, which I had drained, and a bowl of.	0.0919		
$X_{\mathbf{i}}: I$ was standing in their room, with the window open.	0.0919		
$X_{\mathbf{i}}: I$ of course, I had put my coat on.	0.0966		
$X_{\mathbf{i}}: I$ used to wash my hands.	0.0000		
$X_{\mathbf{i}}: I$ I had rolled a towel off the rack and was drying my hair.	0.0082		
$X_{\mathbf{i}}: I$ was taking a shower.	0.0000		
$X_{\mathbf{i}}: I$ I was talking to him.	0.1424		
$X_{\mathbf{i}}: I$ I had been wiping my pants on the sides of shorts and shirt, like a dirty secret.	0.0325		
$X_{\mathbf{i}}: I$ I had been holding a glass of orange juice, which I had drained, and a bowl of.	0.0000		
$X_{\mathbf{i}}: I$ I brushed my teeth and brushed my hair. I even dried my hair, and then I went.	0.0559		
$X_{\mathbf{i}}: I$ I would take a breath.	0.0784		
$X_{\mathbf{i}}: I$ I brushed my teeth, put on deodorant, shared, dried my hair.	0.0489		
$X_{\mathbf{i}}: I$ I had brushed my teeth, applied makeup, and removed my contacts.	0.0783		
$X_{\mathbf{i}}: I$ I removed all my jewelry.	0.0000		
$X_{\mathbf{i}}: I$ I had to take my shoes off.	0.0487		

Table A.8: Example 2a: the first plausible pair of the 63-th instance in COPA-DEV, matched interventions are highlighted. Here $E_1 : I$ was preparing to wash my hands and $E_2 : I$ put rubber gloves on.

Sampled Covariates X	$\ q(x; \Lambda) - q(x; E_1)\ _p$	E_1 and Interventions \mathcal{A}	$\mathbb{P}(\cdot < E_2)$
X_1 : I had scrubbed the kitchen floor and the sink and, uh, that kind of thing.	0.2191		0.5023
X_2 : There was the need to remove the rubbish from the front garden.		E_1 : I was preparing to clean the bathroom. A_1 : I was building the car instead of the house since the house was incompatible with cleanliness. preparing to clean the bathroom.	0.3765
X_3 : I had been walking up and down the hall, running my hands over the wood-paneled.	0.1228	A_2 : I was so bad at the job that I was preparing to clean the bathroom.	0.4935
X_4 : I had put a load of clothes, some books and some DVD's into the washing machine.	0.0598	A_3 : I was doing a job preparing to clean the bathroom.	0.5171
X_5 : I wanted to take out the trash and empty all the containers.	0.0703	A_4 : A woman was preparing to clean the bathroom.	0.5083
X_6 : I had put a couple of washcloths on the bathroom counter.	0.0325	A_5 : Kevin was preparing to clean the bathroom.	0.4889
X_7 : I had to deal with the dirty clothes hamper and the clothes on the floor.	0.0651	A_6 : Emily was preparing to clean the bathroom.	0.4522
X_8 : I had decided to clean the stove.	0.1247	A_7 : I was going to take a bath instead of to clean the bathroom.	0.3744
X_9 : I decided to take a short break and get some ice-cream.	0.1402	A_8 : I was able to do the cleaning in time, and was pretty good at it, although the to clean the bathroom.	0.3202
X_{10} : I had vacuumed the room.	0.0952	A_9 : I was going to clean the bathroom.	0.4437
X_{11} : I went down to the living room and turned on the TV.	0.0000	A_{10} : I was preparing to clean the bathroom.	0.5023
X_{12} : I needed to find a bottle opener.	0.0078	A_{11} : Bill was preparing to clean the bathroom.	0.4967
X_{13} : I prepared a simple salad and some crackers and cheeses in our very cute wooden bowl.	0.0513	A_{12} : Emily was preparing to clean the bathroom.	0.4727
X_{14} : I scribbled down the sink, the counter, the tile and my hands.	0.1011	A_{13} : I was preparing to cook dinner the bathroom.	0.5102
X_{15} : I turned on the TV to catch my favorite news programme, when the host came on and said.	0.0583	A_{14} : I was preparing to sleep the bathroom.	0.4998
X_{16} : I always take my shower.	0.1601	A_{15} : I was preparing to cook dinner for my family, the bathroom.	0.4288
X_{17} : I took a bath and washed my hair.	0.0437	A_{16} : I was preparing to clean the bathroom.	0.5098
X_{18} : I had to clean the roses.	0.0549	A_{17} : I was preparing to clean the kitchen table.	0.5073
X_{19} : I showered, put on a clean pair of pants and a shirt.	0.0798	A_{18} : I wasn't preparing to clean the bathroom.	0.4573
X_{20} : I did my normal routine.	0.1496	A_{19} : I didn't want to wash either the towels or the sponge. I was preparing to clean the bathroom.	0.3829
X_{21} : I washed the coffee table, and before that, I'd vacuumed the floor.	0.0645	A_{20} : I was not preparing to clean the bathroom.	0.4836
X_{22} : I needed to flush the toilet.	0.2527	A_{21} : When I was done, I was preparing to clean the bathroom.	0.2748
X_{23} : I went to the kitchen to put away another load of dishes.	0.0717	A_{22} : I was preparing to clean the bathroom.	0.4897
X_{24} : I needed to put on a pair of rubber gloves.	0.0505	A_{23} : no one was preparing to clean the bathroom.	0.4948
X_{25} : I'd already wiped the kitchen counter.	0.1418	A_{24} : I was not rushing too much and didn't get to clean the bathroom.	0.6000
X_{26} : I would take out the trash.	0.1548	A_{25} : I was not able to do the dishes, so I had to do the dishes. I had no problem washing to clean the bathroom.	0.4008
X_{27} : I used to clean the living room and the kitchen, and even the bathroom sometimes.	0.0765	A_{26} : I was not going to clean the bathroom.	0.4419
X_{28} : I made sure the refrigerator was stocked.	0.1146	A_{27} : I was not supposed to was preparing to clean the bathroom.	0.5107
X_{29} : As a last resort, I would always open the medicine cabinet and remove any expired birth control pills.	0.0505	A_{28} : No one was preparing to clean the bathroom.	0.4948
X_{30} : Of course, I had put my makeup on.	0.0503	A_{29} : No one was preparing to clean the bathroom.	0.4911
X_{31} : I had to clear away the table, then rinse dishes and clean the table.	0.1011	A_{30} : I was preparing to cook dinner the bathroom.	0.5102
X_{32} : I checked on the kids, who were doing what they normally did.	0.0773	A_{31} : I was preparing to skip the whole the bathroom.	0.4682
X_{33} : I was taking a shower.	0.1399	A_{32} : I was preparing to wash my hands, but I forgot to take the bathroom.	0.4426
X_{34} : I decided to organize the cabinets in the kitchen, because organizing might calm me down.	0.0549	A_{33} : I was not able to clean the kitchen table.	0.5073
X_{35} : I was doing laundry.	0.1130	A_{34} : I was preparing to clean the bathroom all the time. I was not able to clean the bathroom all of the time and it would take forever.	0.4897
X_{36} : I would remove the dishes stored in the kitchen, and then, at last, I would clean.	0.0430	A_{35} : No one was preparing to clean the bathroom.	0.5059
X_{37} : I turned on the radio.	0.0814	A_{36} : I was preparing to clean the kitchen counter.	0.5023
X_{38} : Though, I'd go to the kitchen and make a small snack for myself.	0.1086	A_{37} : I was preparing to clean the bathroom.	0.4850
X_{39} : I was taking a shower.	0.0551	A_{38} : I smelled the cookies and wondered if mom was preparing to clean the bathroom.	0.4470
X_{40} : I made myself a green smoothie.	0.0890	A_{39} : I drew it could be used by anyone, just a stranger, preparing to clean the bathroom.	0.4522
X_{41} : I would remove the dishes stored in the kitchen, and then, at last, I would clean.	0.0454	A_{40} : Emily was preparing to clean the bathroom.	0.5137
X_{42} : I brushed my teeth and took a shower, but I was not very interested in either task.	0.3539	A_{41} : She was preparing to clean the bathroom.	0.3144
X_{43} : I would take a good look at all of the objects in the room, cataloging and organizing.	0.0952	A_{42} : I was hoping someone would fill in on the details, so I tried to write to clean the bathroom.	0.4437
X_{44} : I'd brushed my teeth, put on deodorant and makeup, and made sure.	0.1586	A_{43} : I was going through the old photos and couldn't decide which mirror to use for my mirror. The Mirror was too old to clean the bathroom.	0.4813
X_{45} : I was making some soup.	0.0000	A_{44} : I was preparing to clean the bathroom.	0.5023
X_{46} : I had to take the covers off of my bed and the bedspread.	0.0293	A_{45} : I was preparing to put the clothes on the bathroom.	0.5046
X_{47} : I had to clear all the stuff out of my room.	0.0931	A_{46} : I was preparing to put in the sewing machine the bathroom.	0.4958
X_{48} : I took a shower and had a leisurely breakfast.	0.1011	A_{47} : I was preparing to cook dinner the bathroom.	0.5102
X_{49} : I'd had a shower, washed my face, dried it with a towel then put.		A_{48} : I was going to get us a drink of water and maybe some dinner.	

Table A.9: **Example 2b:** the second plausible pair of the 63-th instance in COPA-DEV, matched interventions are highlighted. Here $E_1 : I$ was preparing to clean the bathroom, and $E_2 : I$ put rubber gloves on.

Sampled Covariates \mathcal{X}	$\ q(\mathbf{x}; \mathbf{A}) - q(\mathbf{x}; \mathbf{E}_1)\ _p$	\mathbf{E}_1 and Interventions \mathcal{A}	$\mathbb{P}[\cdot \prec \mathbf{E}_2]$
X_1 : He'd had nothing in his pockets but his father's pocket watch, and some old coins he'd had.			
X_2 : There had been the time he'd been a little boy, about four years old, and...			
X_3 : He had been a very small fish in a very small pond.			
X_4 : He had contained a knife and a set of keys—but there was nothing in it now.			
X_5 : It seemed only to constrain its breath and blood.			
X_6 : He was a slave, and his owner used him roughly when displeased.	0.1069		0.2980
X_7 : He had put a couple of coins in his pouch.	0.0706	A_1 : His pocket was filled with coins.	0.1394
X_8 : His wallet had been empty.		A_2 : His pocket however had been filled with coins.	0.3912
X_9 : It was filled with crumpled-up bills.		A_3 : His pocket contained nine corgage filled with coins.	0.3936
X_{10} : He had been a slave-hunter to the west, a murderer in the service of his city.	0.0620	A_4 : His pocket had a large amount of space filled with coins.	0.0620
X_{11} : It had been full of notes.	0.1457	A_5 : His pocket was lost when he accidentally took some with coins.	0.2749
X_{12} : It had been half-filled with tobacco and a dirty handkerchief.	0.0626	A_6 : His pocket was empty with coins.	0.2597
X_{13} : He had been working as a waiter at the Côte du Rhône at the Hotel Mont.	0.0938	A_7 : His pocket was filled with sandals and at least one with coins.	0.2331
X_{14} : He couldn't even have given the goldfish.	0.0747	A_8 : A cowboy on the back of a wagon was filled with coins.	0.1933
X_{15} : It was empty, but he could think of no problem more urgent than collecting them.	0.0364	A_9 : A pocket iron was filled with coins.	0.1925
X_{16} : He'd been just like them.	0.1158	A_{10} : Mark was filled with coins.	0.0994
X_{17} : He'd been a beggar.	0.2276	A_{11} : His pocket did not work as well as the shirt, because the shirt was filled with coins.	0.1870
X_{18} : He'd been a farmer.	0.1331	A_{12} : His pocket wasn't filled with coins.	0.2477
X_{19} : The contents of his pockets would have been only a pair of pants, a shirt and, if it,	0.1170	A_{13} : His pocket was not filled with coins.	0.1345
X_{20} : He'd ridden them at the old folk's home, where he lived.	0.1539	A_{14} : His pocket was empty but his pocket had nine with coins.	0.0834
X_{21} : He had only had the money to keep him going.	0.1590	A_{15} : His pocket was not holding the lotion with coins.	0.1852
X_{22} : He didn't feel too well.	0.1126	A_{16} : His pocket was not touching with coins.	0.0209
X_{23} : He'd been just like them.	0.3496	A_{17} : No matter how you feel about country music [...] the fact that it featured the really catchy John Denver does not appeal to me was filled with coins.	0.1971
X_{24} : The police had come and taken the man's wallet.	0.1112	A_{18} : No pocket book was filled with coins.	0.2345
X_{25} : He had been wearing a thick bracelet with a chain and gold links, a gift from the wife of...	0.0807	A_{19} : Nothing was filled with coins.	0.2980
X_{26} : He'd been making his way through a series of shops.	0.0000	A_{20} : His pocket had been filled with coins.	0.3352
X_{27} : He'd been standing with his back to the wall, holding an umbrella over his head.	0.0820	A_{21} : His pocket was heavily filled with coins.	0.2535
X_{28} : He'd been trying to buy himself with the money that came from the sale of his books.	0.0239	A_{22} : His pocket was snuffed and he decided to find a wallet instead because the pocket might have coins with coins.	0.0913
X_{29} : Of course, he had been a slave.	0.2621	A_{23} : His pocket was full with coins.	0.3068
X_{30} : He used to have some.	0.0481	A_{24} : Her bag was filled with coins.	0.1762
X_{31} : He'd carried his entire life in his pockets.	0.0753	A_{25} : Someone was filled with coins.	0.2306
X_{32} : He was wearing a shapeless dark coat with a dark salt underneath, a black hat, and...	0.0730	A_{26} : Her wallet was filled with coins.	0.2127
X_{33} : He'd been a poor kid that the state had taken from his mother and put in one home after...			
X_{34} : A few years back, he had an old steel frame.			
X_{35} : He would have sold his coat, and his shirt, and then the little shoes upon his feet.			
X_{36} : He must have been a poor man, and probably very pious.			

Table A.10: **Example 3a:** the first plausible pair of the 79-th instance in COPA-DEV, matched interventions are highlighted. Here E_1 : His pocket was filled with coins. and E_2 : The man's pocket jingled as he walked.

Sampled Covariates X	$\ q(\mathbf{x}; \mathbf{A}) - q(\mathbf{x}; \mathbf{E}_1)\ _p$	\mathbf{E}_1 and Interventions \mathcal{A}	$\mathbb{P}(\cdot < \mathbf{E}_2)$
X_1 : He'd had nothing in his pocket to worry about.			
X_2 : There had been no hole, just a thin line of cloth.			
X_3 : He cut off all his fingers, including those on his left hand.			
X_4 : He had put a knife in the pocket of his coveralls.			
X_5 : He was a young, handsome fellow with dark blue eyes and black curly hair.	0		
X_6 : He had put a couple of sticks in his pouch.	0.1075		
X_7 : His wallet had been in his hand, and he had thrown the wallet on the ground as.	0.2156		
X_8 : He had been a stranger to himself.	0.0682		
X_9 : It had been in his mouth.	0.0660		
X_{10} : He went down to the river to check the damage.	0.0875		
X_{11} : He stuffed a paper bag in place of his wallet.	0.0849		
X_{12} : It was a secret, what with the police and all.	0.2149		
X_{13} : He had been afraid to kill anyone.	0.0830		
X_{14} : He'd hidden his heart, but not in a safe.	0.0707		
X_{15} : He'd thought she'd just pulled it out of thin air, but there it was.	0.1181		
X_{16} : He didn't feel too bad about taking the wallet from your wallet, but after he se.	0.2376		
X_{17} : He'd been just like any other.	0.1161		
X_{18} : The hole had been on the inside of his coat.	0.1128		
X_{19} : He had been wearing his uniform cap with the rank insignia, and he put it on.	0.1479		
X_{20} : He hadn't even looked at it.	0.0869		
X_{21} : He'd seen up the hole in his leg.	0.1461		
X_{22} : There was nothing in it, not even the thorns, and there was nothing there but.	0.3286		
X_{23} : No one had seen it.	0.0871		
X_{24} : He had not wanted to talk about it.	0.0764		
X_{25} : The two men had argued.	0.0456		
X_{26} : He'd been thinking of the boy's parents.	0.0715		
X_{27} : He had been trying to reach the hospital.	0.1560		
X_{28} : He had put a small basket of powder in her shoe, had used a hairpin to.	0.0614		
X_{29} : He'd been a good kid that the others seemed to like.	0.0555		
X_{30} : He had hidden the bullet in his leg.	0.0588		
X_{31} : As a little girl, before she had even known what it meant to be.	0.0311		
X_{32} : He had hidden it in the bottom of pot of ointment.	0.0499		
X_{33} : He had done nothing but, he could do nothing now but lie and wait, and be.			
X_{34} : It was not easy to find a place to hide the gun and if he was asked about.			
X_{35} : He'd had a small knife to cut his clothes, but he didn't.			
X_{36} : He always bought new jeans every time he went to Kaart.			
X_{37} : He'd had nothing.			
X_{38} : He had seen the other pocket.			
X_{39} : The pocket was for the book.			
X_{40} : He had to get the tire.			
E_1 : He sewed the hole in his pocket.			
A_1 : Then he sewed the hole in his pocket.			
A_2 : The boy was grumpy in high school, but happy at school, so the teacher taught him sewed the hole in his pocket.			
A_3 : He cut pieces from the plate sewed the hole in his pocket.			
A_4 : He pulled the toy gun with the hole in his pocket.			
A_5 : He burnt the hole by pulling the hole in his pocket.			
A_6 : He pulled off the blanket and got a hole in his pocket.			
A_7 : He sewed the quilt better than a teepee because the teepee was a sloppy job in his pocket.			
A_8 : He sewed with a towel more in his pocket.			
A_9 : He sewed the hole with a wire instead of a plier in his pocket.			
A_{10} : He couldn't sew the hole in his pocket.			
A_{11} : No matter how you feel about country music I for one can't stand it despite my Houston roots), this only instilled sewed the hole in his pocket.			
A_{12} : He never knew how to sewed the hole in his pocket.			
A_{13} : He never filled the hole in his pocket.			
A_{14} : He couldn't bend the iron rod and instead tied the hole in his pocket.			
A_{15} : He had the hole in his pocket.			
A_{16} : He sewed no better than the Ielaxo which cut off his eye in his pocket.			
A_{17} : He sewed no better with the machine than with the method, because the machine was not precise in his pocket.			
A_{18} : He sewed not only the hole but also the whole ball inside the hole in his pocket.			
A_{19} : Jack sewed the hole in his pocket.			
A_{20} : Someone sewed the hole in his pocket.			
A_{21} : He screened up sewed the hole in his pocket.			
A_{22} : He flunked out of high school, ended up in a strange town, and started writing about the weird the hole in his pocket.			
A_{23} : He stabbed Wisbech with a rope in the hole in his pocket.			
A_{24} : He sewed not only the hole but also the whole ball inside the hole in his pocket.			
A_{25} : He sewed the turkey with a T-shirt in his pocket.			
A_{26} : He sewed the bell necklace in his pocket.			
A_{27} : He sewed the rope with a chisel in his pocket.			

Table A.11: **Example 3b:** the second plausible pair of the 79-th instance in COPA-DEV, matched interventions are highlighted. Here E_1 : He sewed the hole in his pocket and E_2 : The man's pocket jingled as he walked.