

SSL FRAMEWORK FOR CAUSAL INCONSISTENCY BETWEEN STRUCTURES AND REPRESENTATIONS

Anonymous authors

Paper under double-blind review

ABSTRACT

1 The cross-pollination of deep learning and causal discovery has catalyzed a burgeon-
 2 ing field of research, seeking to elucidate causal relationships within non-statistical
 3 data forms like images, videos, and text. Such data, often being named ‘indefinite
 4 data’, exhibit unique challenges—inconsistency between causal structure and rep-
 5 resentation, which are not common in conventional data forms. To tackle this issue,
 6 we theoretically develop intervention strategies suitable for indefinite data and
 7 derive causal consistency condition (CCC). Moreover, we design a self-supervised
 8 learning (SSL) framework that considers interventions as ‘views’ and CCC as a
 9 ‘philosophy’ with two implement examples on Supervised Specialized Models
 10 (SSMs) and Large Language Models (LLMs), respectively. To evaluate pure in-
 11 consistency manifestations, we have prepared the first high-quality causal dialogue
 12 dataset- *Causalogue*. Evaluations are also performed on three other downstream
 13 tasks. Extensive experimentation has substantiated the efficacy of our methodology,
 14 illuminating how CCC could potentially play an influential role in various fields.
 15 Our code is available in url of anonymous code and data.

16 1 INTRODUCTION

17 With the integration of deep learning and causal inference (Sauter et al., 2023; Balashankar et al.,
 18 2023; Lu et al., 2023), increasingly extensive non-statistical data forms, involving images (Jerzak
 19 et al., 2022; Ribeiro et al., 2023), text (Zhang et al., 2023b), and videos (Bagi et al., 2023), have
 20 been drawn into the field of causal discovery. Numerous causality-related studies (Chen et al.,
 21 2023c;a) suggested that these non-statistical data present two fundamental differences with traditional
 22 statistical data: the **representation** and the **structure** differences. Specifically, 1) non-statistical data
 23 (such as images, text, videos), which requires deep representations (such as matrices, embeddings,
 24 optical flow), to participate in causal inference (Schölkopf et al., 2021); but statistical data (like body
 25 temperature, blood pressure, age) inherently exists in a numerical format. 2) unlike statistical data
 26 originating from a fixed causal structure, non-statistical data is drawn from various underlying causal
 27 structures (Löwe et al., 2022). These studies further categorize the data requiring deep representations
 28 and accommodating multiple causal structures to “**indefinite data**”.

29 We observe that indefinite data introduces an inconsistency between structure and representation,
 30 which has not been encountered in other data forms yet. The process of learning causal representations
 31 creates divergence from the process of learning causal structures due to incorporating additional
 32 non-linear parameters. This causal inconsistency manifests in conflicting causal conclusions being
 33 drawn from the structures and representations, which can lead to poor outcomes in high-level causal
 34 models including identifying shortcuts (Wu et al., 2022; Fan et al., 2022; Feng et al., 2023), predicting
 35 incorrect spans (Zhao et al., 2023; Chen et al., 2020; Zhao et al., 2022b), and domain generalization
 36 (Magliacane et al., 2018; Yue et al., 2021; Chen & Bühlmann, 2021).

37 Nevertheless, existing research has overlooked causal inconsistency on indefinite data. In conventional
 38 data forms, multiple methods pivoting on interventions (Ahuja et al., 2023; Lyle et al., 2023), transfer
 39 entropy (Zhou et al., 2022; Silini & Masoller, 2021) and covariance matrix (Kong et al., 2023; Peña,
 40 2023) have naturally satisfied the causal consistency. However, these methods are impeded by other
 41 representational or structural conflicts when extended to indefinite data. Particular methods available
 42 for indefinite data, such as Yu et al. (2019), Chen et al. (2023b), and Löwe et al. (2022), only focus
 43 on how to achieve causal identifiability with various hypotheses.

44 Therefore, to step-by-step achieve causal consistency between the structure and representation of
45 indefinite data, the main contributions of this paper are as follows:

46 In Section 2, we review the background of indefinite data, covering the involved definitions, exam-
47 ples, assumptions, and related works. Following that, we delve into an in-depth analysis of what
48 distinguishes indefinite data, and why causal inconsistency arises in this data form.

49 In Section 3, we introduce a general definition of interventions that enables computing the relationship
50 strength of two target variables without backdoor paths and known distributions. This allows a
51 range of consistency theories, such as exact transformation (Rubenstein et al., 2017) and causal
52 abstraction (Beckers & Halpern, 2019), to be expanded to indefinite data potentially. Inspired by these
53 theories, we further propose a causal consistency condition (CCC). It describes that if the strength
54 sets of two causal models are equivalent given an equivalent intervention set, then the two causal
55 models are consistent.

56 In Section 4, we design a self-supervised learning (SSL) framework that utilizes the CCC as a
57 philosophy, where the causal structure and representation are allocated to separate causal models,
58 whose consistency needs to be verified. Different interventions can be regarded as different “views”,
59 and the measures to gauge causal strength are treated as “augments”. The strength sets are label-
60 agnostic, making the assurance of equivalent strength sets the learning goal of this SSL framework.
61 Additionally, we offer two implement examples — one embodies a trainable module for supervised
62 specialized models (SSMs) and the other executes a prompt instruction for unsupervised large
63 language models (LLMs).

64 In Section 5, we introduce an innovatively constructed dataset— “*Causalogue*” for testing causal
65 consistency, along with the description of its generation process. The dataset comprises 1638
66 dialogues generated by GPT-4 (OpenAI, 2023), with a strategic design that explicates which previous
67 utterances were known when generating each utterance.

68 In Section 6 and Section 7, we conduct experiments on both the *Causalogue* and real-world datasets,
69 validating the accuracy of identifying causal models, the improvements of our proposed SSL frame-
70 work to causal consistency, and effectiveness in three downstream tasks. Additionally, we discuss the
71 potentially crucial role of the CCC within broader research fields.

72 In summary, this paper contributes: insights into general intervention, the causal consistency condition,
73 an SSL framework for testing causal consistency, two corresponding implementation examples, a
74 new causal dataset, and extensive experimentation.

75 2 BACKGROUNDS AND RELATED WORKS

76 2.1 CAUSAL DATA AND INDEFINITE DATA

77 **Definition 1** (Causal Data). ¹ *The causal relationships exist in a dataset $\mathbf{D} = \{X_s\}_{s=1}^S$ which has S
78 samples and M ($M \geq 1$) causal structures ($\mathcal{G} = \{\mathcal{E}_m, \mathcal{V}_m\}_{m=1}^M$). Each structure \mathcal{G}_m corresponds
79 to several samples separately. Hence, each sample $X_{s,m} \in \mathbb{R}^{N \times D}$ belongs to a causal structure
80 $\mathcal{G}_m = \{\mathcal{E}_m, \mathcal{V}_m\}$ and consists of N_m variables: $X_s = \{x_{s,m,n}\}_{n_m=1}^{N_m}$. $\hat{x}_{s,m,n} \in \mathbb{R}^{1 \times D}$ ($D \geq 1$)
81 represents the causal representation of a variable $x_{s,m,n}$ where D denotes the dimension of the
82 causal representation. We assume that the number of causal skeletons is equal to the number of
83 causal structures. Based on the above datasets, we define three data paradigms:*

- 84 • **Definite Data:** *The causal structure is single-skeleton ($M = 1$) and the causal variable is*
85 *single-value ($D = 1$).*
- 86 • **Semi-Definite Data:** *The causal structure is single-skeleton ($M = 1$) and the causal*
87 *variable is multi-value ($D > 1$), or the causal structure is multi-skeleton ($M > 1$) and the*
88 *causal variable is single-value ($D = 1$).*

¹The skeleton M and variable dimension D serve to broaden perspectives on causal data, hence introducing certain conflicts with traditional cognition of causal model. This caused previous reviewers to struggle with conceiving what indefinite data looks like, and why we distinguish indefinite data from other 2 paradigms via skeleton and dimension. Therefore, we dedicatedly established Appendix A, which elucidates these questions through abundant data examples and details the preliminaries including SCMs.

- 89 • **Indefinite Data:** The causal structure is multi-skeleton ($M > 1$) and the causal variable is
90 multi-value ($D > 1$).

91 Definition 1 redefines 3 types of data paradigms from the 2 perspectives of structure and representation.
92 The latter two paradigms often carry incomplete or ambiguous causal labels. Therefore, we adopt the
93 concept of “skeleton” to stand for the causal structure due to the unclear structure labels. Moreover,
94 given the prevalent indefinite datasets mostly include modals like textual conversations and video
95 sources, we propose two hypotheses compatible with these modals:

96 **Hypothesis 1** (Causal Identifiability). The natural order (e.g., time-order) w.r.t. $\{x_{s,m,n}\}_{n_m=1}^{N_m}$ is
97 defined as a linear order $\prec_{X_{s,m}}$. Given that causal order w.r.t. $\{x_{s,m,n}\}_{n_m=1}^{N_m}$ is defined as a partial
98 order $\preceq_{X_{s,m}}$, $\forall \langle x_1, x_2 \rangle \in \prec_{X_{s,m}}$ (i.e., $x_1 \prec_{X_{s,m}} x_2$), there must be $\langle x_1, x_2 \rangle \in \preceq_{X_{s,m}}$.

99 **Hypothesis 2** (Causal Emergence). The causal generative process of multi-value representation
100 is composed of non-autonomous modules that inform or influence each other, meaning that the
101 representation is causally entangled over all dimensions, that is, $E(\hat{x}_{s,m,n}) \doteq x_{s,m,n}$.

102 Hypothesis 1 illustrates the natural linear order of indefinite data (e.g., $\{U_1, U_2, U_3, U_4\}$, where U_1
103 to U_4 respectively represent 4 utterances appearing in time-series, and $U_i \prec U_j$ indicates that U_i
104 precedes U_j in time) belongs to the causal partial order. Consequently, the adjacency matrix of the
105 natural linear order is a triangular matrix, which naturally corresponds to a DAG. Thus, there is no
106 need for measures such as acyclic constraints (Zheng et al., 2018) to ensure causal identifiability.

107 Hypothesis 2 can provide insights into causal representation from the perspective of the law of large
108 numbers. For statistical data, such as temperature, we need enough samples to grasp its characteristics
109 (or distributions) in a particular environment. However, for non-statistical data, this is unnecessary.
110 For instance, any sentence is enough to express its semantics, a single image can be read for its
111 content. This hypothesis releases the limitation of insufficient samples, allowing us to achieve the
112 causal consistency condition through strength sets.

113 2.2 WHY CAUSAL INCONSISTENCY ARISES?

114 Figure 1 visualizes evaluation results of causal
115 consistency via tested 5 methods: PC (Kalisch
116 & Bühlman, 2007), ACD (Löwe et al., 2022),
117 DAG-GNN (Yu et al., 2019), CAE (Chen et al.,
118 2023b), and biCD (Chen et al., 2023c). They
119 represent prevalent methods in specific data
120 forms, respectively. Two conclusions can be
121 obtained from Figure 1: 1) The strongest causal
122 inconsistency is found in indefinite data forms
123 ($M > 1, D > 1$), while definite data ($M = 1,$
124 $D = 1$) performs the weakest causal inconsis-
125 tency. 2) When existing methods are applied to
126 non-default data forms (hollow markers), their
127 consistency performance is always inferior to
128 the native methods for that data form.

129 Either $M > 1$ or $D > 1$ contributes to a rise
130 in inconsistency. In general, when $M > 1$, the
131 optimization for causal strength changes from
132 f to $\sum \alpha_m f_m$, which leads to a lower accuracy
133 of causal structure than the ones of $M = 1$, due
134 to the existence of Pareto Optimality (Censor,
135 1977). Meanwhile, $D > 1$ introduces deep rep-
136 resentations, resulting in an inexact transfor-
137 mation. When both $M > 1$ and $D > 1$ are present, we assume f_m can be decoupled from \tilde{p}_φ . The causal
138 structure learning reads $\mathcal{G}_m = h(X, \varphi)$, and the causal representation learning is $\hat{X}_m = h(X, \varphi, \theta)$.
139 The different learning processes with additional parameter θ intrinsically increase more inconsistency
140 than other data paradigms. Details are elaborated on in Appendix A.3.

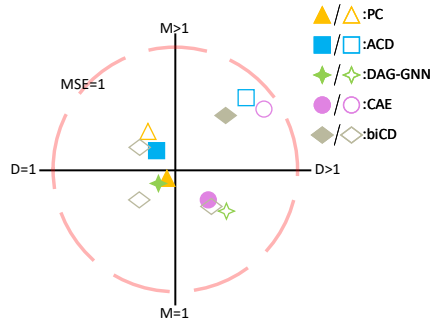


Figure 1: We compared the consistency of different methods in 3 data paradigms (if available). The consistency was represented by the MSE of the similarity matrices for structure and representation. The filled markers represent methods being in their default data forms, while the hollow markers signify that they are in extendable but non-default data forms.

141 2.3 RELATED WORKS: CAUSAL METHODS ON DIFFERENT DATA PARADIGMS

142 **M=1 and D=1:** Methods applied to definite data (Janzing et al., 2012; Ramsey et al., 2017; Cai et al.,
143 2020) can ensure consistency due to their effective utilization of statistical advantages.

144 **M>1 and D=1:** To address the challenge of multi skeletons, some methods try to learn an invariance,
145 such as ACD (Löwe et al., 2022). This range of approaches (Huang et al., 2020a; Dhir & Lee,
146 2020; Huang et al., 2020b; 2019) uses neural networks to automate the learning the distribution of
147 causal structure $\sum_{m=1}^M \log(P(X_m))$. Despite the promising performance of the whole structures,
148 the accuracy of individual structures decreases thereby increasing inconsistency.

149 **M=1 and D>1:** Similarly, for the more favorable causal representation, assumption of the stationary
150 structure becomes the most prevailing choice involved interpretability (Fan et al., 2022; Wu et al.,
151 2022), relationship analysis (Chen et al., 2023b; Zhao et al., 2022a), and domain generalization (Lv
152 et al., 2022; Jiang & Veitch, 2022). Among these methods, the critical cause of inconsistency stems
153 from inaccuracies in recovering relations. Specifically, the complexity of the data makes it challenging
154 to identify exact causal relationships so that none can achieve the same perfect level as that in definite
155 data.

156 **M>1 and D>1:** The inconsistency is more a result of additional parameter in the learning processes
157 of structures and representations. Especially, indefinite data is in its infancy, and current methods focus
158 more on enhancing causal accuracy, with explorations into causal consistency yet to be conducted.
159 For example, while DAG-GNN (Yu et al., 2019) can be applied to indefinite data, it does not yield
160 satisfactory results compared to its default data form (M=1 and D=1). Methods such as CAE (Chen
161 et al., 2023b) and biCD (Chen et al., 2023c), although they have improved causal accuracy, display
162 significant causal inconsistency, particularly in advanced tasks with incomplete labels.

163 3 HOW TO CHECK CAUSAL CONSISTENCY

164 An observed variable x_t of indefinite data does not satisfy $P(y|x_t = t_1), P(y|x_t = t_2), \dots$ without
165 adequate samples $t = t_1, t = t_2, \dots$ so that the distribution is not clear for intervention. Hence, we
166 define a general intervention, intending to bypass distribution assumptions to obtain interventions.

167 **Definition 2** (General Intervention). *General intervention is represented by the do_g operator with*
168 *the objective of setting the parent set of the observed variable to \emptyset .*

$$do_g(x_t) := Pa(x_t) = \emptyset \quad (1)$$

169 where $Pa(x)$ represents a parent set of x .

170 Benefitted from Hypothesis 2, effects of $do_g(x_t)$ are equivalent to effects of the set of perfect
171 interventions: $\{do(x_t = t_1), do(x_t = t_2), \dots\}$. (For simplicity, unless specially stated, the term do
172 in the rest of this paper represent either do or do_g .) Definition 2 introduces feasible intervention for
173 indefinite data to allow us to draw inspiration from Definition 6 in Appendix A.1: the consistency of
174 two causal model can be verified under any reasonable intervention. To ensure this idea, the causal
175 models need to include intervention sets and strength sets.

176 **Definition 3** (Causal Model). *Let causal model $M_X = \langle S_X, I_X, F_X^{I_X} \rangle$, where S_X represents an SCM*
177 *for the model with the variable set $X = (x_i : i \in \mathbb{I}_x)$, \mathbb{I}_x is the index of causal partial order over X ,*
178 *$I_X := (do(i, j) : i, j \in \mathbb{I}_x, \preceq_{do})$ represents a set of all reasonable bi-variable perfect interventions*
179 *satisfying partial order, $F_X^{I_X} := (f_X^{do(i,j)} : i, j \in \mathbb{I}_x, \preceq_{do})$ represents the causal strength of set X*
180 *under corresponding interventions.*

181 Considering the presence of front-door paths, it is complicated to directly calculate the causal
182 relationship between any two variables by intervention on just one variable. Consequently, we form
183 the intervention set using bi-variable interventions. For example, in a binary definite data set, the
184 intervention set could be $I_X = \{\emptyset, do(x_1 = 0, x_2 = 0), do(x_2 = 0, x_3 = 1), \dots\}$. In the indefinite
185 data, the intervention set could be $I_X = \{\emptyset, do_g(x_1, x_2), do_g(x_2, x_3), \dots\}$. \preceq_{do} represents that for
186 any pair $\langle x_i, x_j \rangle$ in I_X , where x_I signifies the x_i in the pair $\langle x_i, x_j \rangle$ and all pairs previous to
187 $\langle x_i, x_j \rangle$, it is always in the causal partial order that x_i does not follow x_j . The strength set $F_X^{I_X}$
188 would have the same partial order \preceq_{do} , representing any causal strength of model M_X corresponding

189 to the perfect intervention of I_X . In definite data, according to the causal factorization mentioned in
 190 Appendix A.1, the strength set can be equated to the distribution set $\mathbb{P}_X^{I_X}$. Finally, we would like to
 191 introduce the causal consistency condition:

192 **Theorem 1** (Causal Consistency Condition (CCC)). Let $\mathcal{U}_X = (S_X, I_*, F_X^{I_*})$ and $\mathcal{V}_Y =$
 193 $(S_Y, I_*, F_Y^{I_*})$ be two causal models. The intervention set I_* denotes that there is an identity mapping
 194 between X and Y . If any term $f_{y_1, y_2}^{do(i, j)}$ in $F_Y^{I_*}$ satisfies:

$$f_{y_1, y_2}^{do(i, j)} = f_{x_1, x_2}^{do(i, j)} \quad (2)$$

195 the \mathcal{U}_X is consistent with \mathcal{V}_Y (Proof is given in Appendix B).

196 **Example 1.** In indefinite data, we assume that the causal structure belongs to a definite causal model
 197 \mathcal{U} , and the causal representation belongs to an indefinite causal model \mathcal{V} . Since \mathcal{U} and \mathcal{V} have the
 198 same causal variables, there exists an order-preserving bijection $\omega := I_X \Leftrightarrow I_Y$. If \mathcal{U} and \mathcal{V} satisfy
 199 the CCC, let $F_\circlearrowleft = (f_{a, b} : a, b \in \mathfrak{I}_x, (a \preceq_X i, b = i) \text{ or } (a \preceq_X j, b = j), \preceq_{do})$ w.r.t. $do(i, j)$. Any
 200 term $f_{y_1, y_2}^{do(i, j)}$ in $F_Y^{I_*}$ satisfies $f_{y_1, y_2}^{\omega(do(i, j))} = f_{x_1, x_2}^{do(i, j)}$ and if the factorization of $f_{y_1, y_2}^{\omega(do(i, j))}$ includes
 201 $f_{a, b} \in F_\circlearrowleft$, it satisfies $f_{y_1, y_2}^{\omega(do(i, j))} = f_{x_1, x_2}^{do(i, j)} = 0$. The conclusion is also satisfied on I_Y .

202 The SMS hypothesis proposed in Schölkopf et al. (2021) elucidated Example 1. Simply put, if
 203 the causal structure is robust, the interventionally-affected conditional probability can not influence
 204 the interventionally-unaffected conditional probability in causal factorization $P(x_1, x_2, \dots, x_s) =$
 205 $\prod_{s=1}^S P(x_s | X_{Pa_{x_s}})$. On the contrary, the interventionally-unaffected conditional probability could
 206 not maintain stability if causal model is unrobust. Therefore, we not only require the strengths of the
 207 intervention nodes to be consistent ($f_{y_1, y_2}^{\omega(do(i, j))} = f_{x_1, x_2}^{do(i, j)} = 0$), but also that the strengths of nodes
 208 without induced paths to the intervention nodes stay consistent ($f_{y_1, y_2}^{\omega(do(i, j))} = f_{x_1, x_2}^{do(i, j)} \neq 0$).

209 4 SSL FRAMEWORK

210 4.1 FORMULATION ARCHITECTURE

211 The causal structure and causal representation of each sample in the indefinite data can be viewed
 212 as belonging to two individual models. Specifically, the causal structure $\mathcal{G}_{s, m} = h(X_{s, m}, \varphi)$
 213 can be seen as part of a definite data causal model $\mathcal{U} = (S_{X_{s, m}}, I_*, F_{X_{s, m}}^{I_*})$, while the causal
 214 representation $\hat{X}_{s, m} = h(X_{s, m}, \varphi, \theta)$ can be considered part of an indefinite data causal model
 215 $\mathcal{V} = (S_{\hat{X}_{s, m}}, I_*, F_{\hat{X}_{s, m}}^{I_*})$ ($\mathfrak{I}_{\hat{x}} = \mathfrak{I}_x, \preceq_{\hat{X}_{s, m}} = \preceq_{X_{s, m}}$). Therefore, there are two causal models (\mathcal{U}
 216 and \mathcal{V}) correspondingly with two causal structures ($\mathcal{G}_{s, m}$ and $\hat{\mathcal{G}}_{s, m}$) and two causal representation
 217 ($X_{s, m}$ and $\hat{X}_{s, m}$). The aim of SSL framework is to establish equivalent strength set $F_{X_{s, m}}^{I_*} = F_{\hat{X}_{s, m}}^{I_*}$,
 218 thereby achieving causal consistency $\mathcal{U} = \mathcal{V}$. We elucidate the roles of “view”, “augment”, and
 219 “philosophy” within our framework as follows.

220 **View:** We define an intervention ($do(i, j) \in I_*$) as a view. For example, $do(x_1, x_2)$, $do(x_2, x_3)$, and
 221 $do(x_1, x_3)$ could be 3 individual views of \mathcal{U} ; $do_g(x_1, x_2)$, $do_g(x_2, x_3)$, and $do_g(x_1, x_3)$ could be
 222 ones of \mathcal{V} .

223 **Augment:** We define the specific measures for obtaining the $F_{X_{s, m}}^{I_*}$ under intervention as the augments.

224 e.g., $f_{\hat{X}}^{do_g(i, j)} = \text{augment}_{do_g(i, j)}(\hat{X})$. $F_{X_{s, m}}^{I_*}$ can be directly obtained via checking $\mathcal{G}_{s, m}$ or $S_{X_{s, m}}$
 225 because it belongs to the definite-data causal model.

226 **Philosophy:** We define the causal consistency between \mathcal{U} and \mathcal{V} as the philosophy. e.g., \mathcal{U} and \mathcal{V}
 227 should satisfy the CCC (Theorem 1): for any view, $f_{\hat{X}}^{do_g(i, j)} = f_X^{do(i, j)}$.

228 Overall, as illustrated in Figure 2. The views of causal structure and causal representation correspond
 229 with each other. The “augment” process derives strength sets separately under these views, which are
 230 then evaluated for consistency according to the philosophy of causal consistency. Within this, both
 231 the “Augment” and “Consistency check” modules in Figure 2 depend on specific implementation.
 232 The “Augment” should not introduce any new parameters, otherwise the optimization would be:

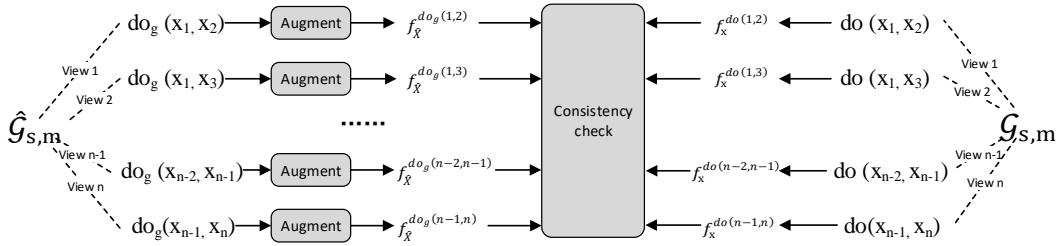


Figure 2: The SSL framework for causal consistency. The grey rectangular boxes represent modules that require specific implementation. From left to center, the process describes how causal representations are transformed into causal strengths. From right to center, the process illustrates how causal structures are converted into strengths.

Table 1: Number of the samples in *Causalogue* Dataset

| Versions | Structure Types | | | | | | | | | | Total |
|----------|-----------------|----------|-----------|----------|--------|---------|----------|---------|----------|-----------|-------|
| | Chain_I | Chain_II | Chain_III | Chain_IV | Fork_I | Fork_II | Fork_III | Fork_IV | Hybrid_I | Hybrid_II | |
| Small | 276 | 84 | 141 | 44 | 257 | 237 | 251 | 67 | 185 | 77 | 1638 |
| Large | 0 | 524 | 508 | 513 | 1215 | 645 | 501 | 372 | 499 | 635 | 5412 |

233 $\min_{\theta, \delta_1, \dots, \delta_K} \sum_{s=1}^S \sum_{m=1}^M \sum_{k=1}^K \mathcal{L}_k(\hat{\mathcal{G}}_{s,m}, \delta_k, \theta)$. The “consistency check” module shares parameters and does not introduce intervention-unique parameters. Therefore, the optimization of our entire
 234 SSL process is written as:
 235

$$\min_{\theta} \sum_{s=1}^S \sum_{m=1}^M \mathcal{L}_k(\hat{\mathcal{G}}_{s,m}, \theta) \tag{3}$$

236 4.2 TWO EXAMPLES FOR IMPLEMENTATION

237 We provide two implementation examples. The first one is on a supervised specialized model
 238 (SSM) generating high-level causal representations and structures. The second is implemented on an
 239 unsupervised large language model (LLM), which can be used to directly infer causal relationships
 240 between utterances (as dialogues are typically indefinite data). The consistency check modules in
 241 both examples are accomplished through similarity matrices, yet the augment modules are completely
 242 different: the first example computes strength by modifying adjacency matrices, while the second
 243 example offers two approaches: prompts and pre-trained models. Detailed implementation specifics
 244 are thoroughly described in Appendix C.

245 5 NEW SIMULATION DATASET-*Causalogue*

246 Existing indefinite datasets suffer from issues including incomplete labeling and insufficient sam-
 247 ples. These numerous entangled problems make it challenging to achieve pure evaluation for the
 248 inconsistency. Additionally, the challenge of manual annotation is considerable, as the presence
 249 of numerous ambiguous samples could make classification boundaries unstable. Fortunately, the
 250 powerful human-computer conversation abilities of LLMs, such as GPT-4, have made automated
 251 annotation possible. Thus, to provide a high-quantity dataset for checking consistency, we have
 252 made an endeavor to generate controlled, causal dialogues via GPT-4 ending up with a new dataset,
 253 *causalogue*. This is the first dialogue dataset that includes comprehensive causal relationship labels
 254 for indefinite data. Besides causality-related tasks, the dataset is available for all tasks related to
 255 dialogue relationships (e.g., dialogue generation, relation extraction, and text classification).

256 The dataset incorporates 10 types of causal structures ($M = 10$), each with several samples (Detailed
 257 number are presented in Table 1, “Small” signifies samples that have been manually checked as
 258 correctly labeled, while “large” refers to all samples generated by GPT-4 without manual verification).

Table 2: Summarization of datasets and baselines

| Ours | Tasks | Datasets | Baselines | Metrics |
|---------------------|-------|----------------------------|--|-----------------------|
| | CD | <i>Causalogue</i> | ACD, DAG-GNN, ACCD, biCD, DisC, DIR | AUROC, HD, F1, MSE |
| Ours _{SSM} | ECPE | RECCON | ACCD, biCD, EDKA-GM, seF | F1 |
| | ERC | MELD, EmoryNLP, DD, IEM | ACCD, biCD, DAG-ERC, DualGAT, MultiEMO | F1 |
| | TAS | GTEA, 50salads, Breakfast | MS-TCN++, ASRF, CETNet, C2F | acc, Edit, F1@k,C-Dis |
| Ours _{LLM} | CD | <i>Causalogue</i> , RECCON | Zero-shot, Zero-shot-Cot, Auto-Cot | F1 |

259 All samples consist of 4 causal variables. In each sample, binary causal relationships have been
 260 labeled between any two utterances. A detailed exposition of the dataset’s attributes and creation
 261 process can be found in Appendix D.

262 6 EXPERIMENTS

263 6.1 DATASETS, BASELINES, AND METRICS

264 Including the Causal Discovery (CD) task on *Causalogue* dataset, we also evaluate our method on
 265 real-world datasets (RECCON (Poria et al., 2021), MELD (Poria et al., 2019), EmoryNLP (Zahiri &
 266 Choi, 2018), DD (Li et al., 2017), IEM (Busso et al., 2008), GTEA (Fathi et al., 2011), 50salads (Stein
 267 & McKenna, 2013), and Breakfast (Kuehne et al., 2014).) spanning three downstream tasks (Emotion-
 268 cause Pair Extraction (ECPE) task, Emotion Recognition in Conversation (ERC) task, and Temporal
 269 Action Segmentation (TAS) task) involving both text and video.

270 The experiments also incorporate a variety of baselines. For the supervised specialized models
 271 (SSMs), it encompasses causal deep models such as ACD (Löwe et al., 2022), DAG-GNN (Yu et al.,
 272 2019), ACCD (Chen et al., 2023b), biCD (Chen et al., 2023c), and intervention deep models like
 273 DisC (Fan et al., 2022), DIR (Wu et al., 2022), and our example (*Ours_{SSM}* in Appendix C.1).
 274 Moreover, we evaluate the downstream tasks with additional SOTA work pertinent to each task,
 275 such as EDKA-GM (Li et al., 2023a), seF (Li et al., 2023b) for ECPE task, DAG-ERC (Shen
 276 et al., 2021), DualGAT (Zhang et al., 2023a), MultiEMO (Shi & Huang, 2023) for ERC task, and
 277 MS-TCN++ (Li et al., 2020), ASRF (Ishikawa et al., 2021), CETNet (Wang et al., 2023), and
 278 C2F (Singhania et al., 2021) for TAS task. In terms of LLMs, we compared prompt-based baselines:
 279 Zero-Shot (Kojima et al., 2022), Zero-Shot-Cot (Kojima et al., 2022), Auto-Cot (Zhang et al., 2022),
 280 and our example (*Ours_{LLM}* in Appendix C.2). on public GPT-4 of the gpt-4-32k-0314 version.

281 For the CD task, we employed an array of metrics: Area Under the Receiver Operating Characteristic
 282 curve (AUROC), F1 score, Hamming Distance (HD), and Mean Squared Error (MSE) to comprehen-
 283 sively evaluate both the precision and consistency of causality. For different downstream tasks, we
 284 utilized their prevalent metrics for evaluations.

285 The criss-cross relationships between these datasets, tasks, metrics, and baselines have been summa-
 286 rized in Table 2 and the details of them are shown in Appendix E).

287 6.2 IMPLEMENTATION DETAILS

288 For the SSM, different pre-training models and implementation parameters were adopted for different
 289 downstream tasks. Detailed descriptions of these variants are provided in Appendix E.4. As for LLM,
 290 we solely implemented a prompt instruction method without adjusting the model parameters. The
 291 models include gpt-4-32k-0314.

292 6.3 SIMULATED DATASET

293 6.3.1 RESULTS OF THE SSMs

294 We conducted experiments for causal accuracy and consistency of the SSMs on the *Causalogue* dataset.
 295 Causal accuracy was evaluated through the performance of causal graphs and causal representations.
 296 Causal consistency, was assessed by measuring the distance within the similarity matrix between the
 297 graphs and representations. As demonstrated by Table 3, our method significantly improved causal
 298 consistency, which correspondingly led to an enhancement of causal accuracy. This also elaborates

Table 3: Results of SSMS on *Causalogue* Dataset. 95% confidence interval shown. All evaluation metrics (except HD) were normalized to the range $[0, 1]$. A value closer to 1 indicates better performance. HD (Hamming Distance) measure the corrected edges from results to labels.

| Methods | Causal Structure | | Causal Representation | | Causal Consistency | |
|---------------------|------------------------|------------------------|------------------------|------------------------|------------------------|------------------------|
| | AUROC | HD | AUROC | F1 | AUROC | 1-MSE |
| ACD | 0.84 \pm 0.02 | 0.89 \pm 0.21 | 0.85 \pm 0.02 | 0.88 \pm 0.01 | 0.51 \pm 0.01 | 0.49 \pm 0.01 |
| DAG-GNN | 0.56 \pm 0.04 | 1.51 \pm 0.34 | 0.90 \pm 0.01 | 0.88 \pm 0.02 | 0.50 \pm 0.01 | 0.49 \pm 0.02 |
| DAG-DisC | 0.68 \pm 0.27 | 1.40 \pm 0.42 | 0.88 \pm 0.02 | 0.87 \pm 0.02 | 0.52 \pm 0.00 | 0.50 \pm 0.01 |
| DAG-DIR | 0.67 \pm 0.38 | 1.36 \pm 0.36 | 0.89 \pm 0.03 | 0.86 \pm 0.03 | 0.51 \pm 0.01 | 0.50 \pm 0.01 |
| ACCD | 0.79 \pm 0.11 | 1.02 \pm 0.15 | 0.93 \pm 0.01 | 0.92 \pm 0.03 | 0.60 \pm 0.05 | 0.59 \pm 0.11 |
| biCD | 0.91 \pm 0.04 | 0.56 \pm 0.10 | 0.86 \pm 0.03 | 0.89 \pm 0.02 | 0.64 \pm 0.04 | 0.59 \pm 0.07 |
| Ours _{SSM} | 0.94 \pm 0.01 | 0.29 \pm 0.05 | 0.94 \pm 0.01 | 0.95 \pm 0.01 | 0.95 \pm 0.01 | 0.92 \pm 0.01 |

299 the fact that, until now, the causal consistency in indefinite data has been often overlooked though it
300 is a crucial problem.

301 Our findings also indicates some additional conclusions. As shown between three ‘DAG’-related conclu-
302 sions. As shown between three ‘DAG’-related base-
303 lines, the intervention methods proposed by DisC and
304 DIR could enhance the causal graph identification ca-
305 pability. This improvement is attributed to that inter-
306 ventions can underlyingly adapt models to cross i.i.d.
307 environment. However, their interventions introduced
308 bias that lies in forming negative samples by combin-
309 ing the causal pattern with the background from
310 other samples in the batch, when CCC sets it as \emptyset .
311 Their intervention concepts do enhance the model’s
312 discriminative capacity for causal patterns and short-
313 cuts, but the sparsity of indefinite data samples has
314 been introduced as bias into contrastive learning. In
315 addition, ACD and biCD are methods specifically
316 targeted at multi-value and multi-skeleton data, re-
317 spectively. Therefore, they have been particularly
318 emphasized in our experimental results.

319 To evaluate the volume of the interventions to results,
320 we tested the performance under conditions ranging from an empty intervention set (no interventions
321 carried out) to the maximum intervention set (all interventions carried out). Figure 3 demonstrates
322 that interventions can significantly enhance causal consistency, thereby improving causal accuracy.
323 Moreover, the size of the complete intervention set is close to 70% of the maximum intervention set,
324 which is also reflected in the Figure 3 as a notable stability after 70%.

325 Lastly, we conducted an ablation study to deter-
326 mine the contributions of each mechanism.
327 Table 4 demonstrates that the specific implemen-
328 tation of the module contributes little, but the
329 causal mechanism (-adj) and causal identifiability
330 (-mask) are the primary contributors to the
331 causal discovery in indefinite data. However, the
332 causal representation benefits more from fitting
333 capacity, hence the decreasing induced by the
334 causal mechanism is not as evident as the other
335 two measures.

336 6.3.2 RESULTS OF THE LLMs

337 Our_{LLM}, was tested on the *Causalogue* and RECCON datasets, with a simple set of experiments
338 detailed in Appendix F. Specifically, we assessed the the upper bound of accuracy of varying methods
339 to calculate Sim^r , and recorded the performance of our model as it approaches this bound. Finally,
340 we illustrated specific question-answer content through a case analysis.

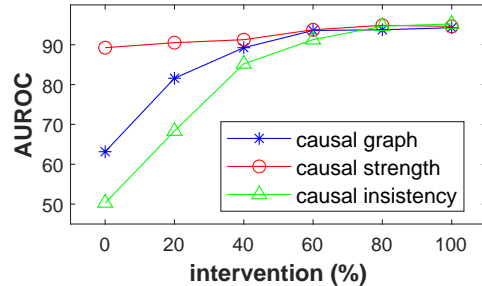


Figure 3: Performance of the Ours_{SSM} under different sizes of intervention. ‘‘Intervention 20%’’ refers to an intervention set composed by randomly selecting 20% *do* operators from the maximum intervention set.

Table 4: Ablation Results of AUROC on three measures. ‘-cos.sim’: replacing cosine similarity with MSE, ‘-adj’ Removing matrix A , ‘-mask’: replacing Hypothesis 1 with no acyclic constraints.

| Model | Structure | Representation | Consistency |
|----------|-----------|----------------|-------------|
| -cos_sim | ↓ 0.02 | ↓ 0.01 | ↓ 0.00 |
| -adj | ↓ 0.34 | ↓ 0.12 | ↓ 0.43 |
| -mask | ↓ 0.46 | ↓ 0.11 | ↓ 0.47 |

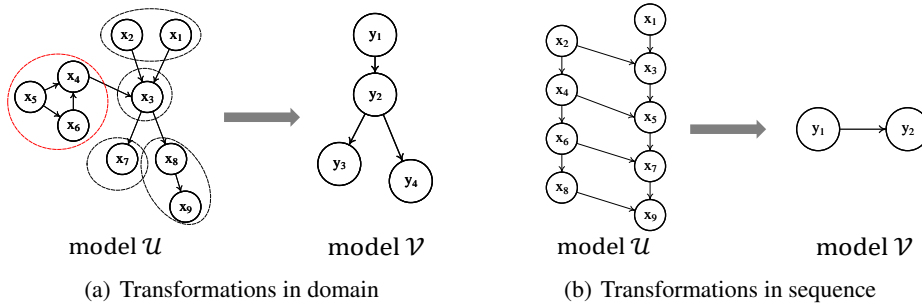


Figure 4: Two potential insights of general causal consistency condition. Black elliptical frames represents the clustering of similar nodes, and red one represents marginalization of irrelevant nodes.

341 6.4 REAL-WORLD DATASETS

342 We analyzed performance on three total downstream tasks, of which two are text modal: ECPE and
 343 ERC, and one is video modal: TAS. Comparative results against corresponding SOTA baselines,
 344 along with visualization cases, can be found in Appendix G. The results collectively suggest that
 345 Ours_{SSM} performs well not only under experiment-environment causal models, but also enables
 346 more appropriate inference under high-level causal models.

347 7 DISCUSSION: MORE INSIGHTS AND CONTRIBUTIONS

348 Focusing on addressing the SSL framework for causal consistency of indefinite data, this paper
 349 contributes to general intervention measures for distribution-unknown data, novel dataset, and
 350 reasoning on LLMs. However, considering the evolving progression of causal abstraction, our study
 351 only represents a beginning. The causal consistency condition (CCC) can only be applied to causal
 352 models with identical variable sets, whereas causal abstraction in definite data can already process
 353 models featuring two distinct variable sets— marking a significant discrepancy in contribution.
 354 Expanding causal abstraction to indefinite data (turning CCC into general CCC) could accurately
 355 search simplified models for intractable and complex models, thereby providing insights for many
 356 fields of deep learning.

357 Figure 4 describes two applications of general CCC-in domains and sequences, illustrating a new func-
 358 tion of transformations from complex to simple models which are not confined to mere consistency
 359 checks. Specifically, transformations in domains primarily include clustering and marginalization
 360 in Figure 4 a: $\tau : (x_1, x_2) \rightarrow (y_1), (x_3) \rightarrow (y_2), (x_7) \rightarrow (y_3), (x_8, x_9) \rightarrow (y_4)$. This could solid-
 361 ify theoretical backing for established methods within various fields if a ω mapping exists on the
 362 intervention sets to make the strength sets equal. One notable application is interpretability of graph
 363 neural networks (Fan et al., 2022; Wu et al., 2022), where the goal is distinguishing between causal
 364 and other patterns amidst numerous nodes. Actually, we have explored the mapping from micro
 365 to macro model in the TAS task: a frame corresponding to a variable in the micro model \mathcal{U} , while
 366 a segment being a variable in the macro model \mathcal{V} , striving to construct equivalent mappings from
 367 frame to segment. The concept that the simplified model \mathcal{V} can be viewed as the emergence of macro
 368 relationships in complex system \mathcal{U} , holding potential in researching areas such as meta-learning
 369 and domain generalization. Transformations in sequence, as shown in Figure 4 b are closely tied to
 370 temporal causal discovery, expressing dynamic processes through stable behaviors—a crucial aspect
 371 of sequence-to-sequence models and time-series forecasting.

372 However, there are significant challenges in implementing general CCC both theoretically and
 373 practically. To accomplish the case of “identity mapping” in causal abstraction on indefinite data,
 374 we introduced innovative intervention measures, frameworks, and validation methods. The general
 375 solution mentioned above requires much more theoretical research and experimental evaluation
 376 than this paper offers, and we eagerly anticipate sharing these findings in the future. In summary,
 377 preserving causal consistency in indefinite data marks an inventive and promising initiation in aligning
 378 causal theory with the deep learning.

379 REFERENCES

- 380 Kartik Ahuja, Divyat Mahajan, Yixin Wang, and Yoshua Bengio. Interventional causal representation
381 learning. In *International Conference on Machine Learning*, pp. 372–407. PMLR, 2023.
- 382 Shayan Shirahmad Gale Bagi, Zahra Gharaee, Oliver Schulte, and Mark Crowley. Generative causal
383 representation learning for out-of-distribution motion forecasting. *arXiv preprint arXiv:2302.08635*,
384 2023.
- 385 Ananth Balashankar, Srikanth Jagabathula, and Lakshmi Subramanian. Learning conditional granger
386 causal temporal networks. In *2nd Conference on Causal Learning and Reasoning*, 2023.
- 387 Sander Beckers and Joseph Y Halpern. Abstracting causal models. In *Proceedings of the aai
388 conference on artificial intelligence*, volume 33, pp. 2678–2685, 2019.
- 389 Carlos Busso, Murtaza Bulut, Chi-Chun Lee, Abe Kazemzadeh, Emily Mower, Samuel Kim, Jean-
390 nette N Chang, Sungbok Lee, and Shrikanth S Narayanan. Iemocap: Interactive emotional dyadic
391 motion capture database. *Language resources and evaluation*, 42(4):335–359, 2008.
- 392 Ruichu Cai, Jincheng Ye, Jie Qiao, Huiyuan Fu, and Zhifeng Hao. Fom: Fourth-order moment based
393 causal direction identification on the heteroscedastic data. *Neural Networks*, 124:193–201, 2020.
- 394 Joao Carreira and Andrew Zisserman. Quo vadis, action recognition? a new model and the kinetics
395 dataset. In *proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp.
396 6299–6308, 2017.
- 397 Yair Censor. Pareto optimality in multiobjective problems. *Applied Mathematics and Optimization*, 4
398 (1):41–59, 1977.
- 399 Hang Chen, Bingyu Liao, Jing Luo, Wenjing Zhu, and Xinyu Yang. Learning a structural causal
400 model for intuition reasoning in conversation. *arXiv preprint arXiv:2305.17727*, 2023a.
- 401 Hang Chen, Jing Luo, Xinyu Yang, and Wenjing Zhu. Affective reasoning at utterance level in
402 conversations: A causal discovery approach. *EMNLP-main*, 2023b.
- 403 Hang Chen, Xinyu Yang, and Qing Yang. Learning to recover causal relationship from indefinite
404 data in the presence of latent confounders. *arXiv preprint arXiv:2305.02640*, 2023c.
- 405 Xinhong Chen, Qing Li, and Jianping Wang. Conditional causal relationships between emotions and
406 causes in texts. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language
407 Processing (EMNLP)*, pp. 3111–3121, 2020.
- 408 Yuansi Chen and Peter Bühlmann. Domain adaptation under structural causal models. *The Journal
409 of Machine Learning Research*, 22(1):11856–11935, 2021.
- 410 Anish Dhir and Ciarán M Lee. Integrating overlapping datasets using bivariate causal discovery. In
411 *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pp. 3781–3790, 2020.
- 412 Aleksandr Ermolov, Aliaksandr Siarohin, Enver Sangineto, and Nicu Sebe. Whitening for self-
413 supervised representation learning. In *International Conference on Machine Learning*, pp. 3015–
414 3024. PMLR, 2021.
- 415 Shaohua Fan, Xiao Wang, Yanhu Mo, Chuan Shi, and Jian Tang. Debiasing graph neural networks
416 via learning disentangled causal substructure. *Advances in Neural Information Processing Systems*,
417 35:24934–24946, 2022.
- 418 Alireza Fathi, Xiaofeng Ren, and James M Rehg. Learning to recognize objects in egocentric
419 activities. In *CVPR 2011*, pp. 3281–3288. IEEE, 2011.
- 420 Qizhang Feng, Ninghao Liu, Fan Yang, Ruixiang Tang, Mengnan Du, and Xia Hu. Degree: De-
421 composition based explanation for graph neural networks. *arXiv preprint arXiv:2305.12895*,
422 2023.

- 423 H.A. Guvenir, B. Acar, G. Demiroz, and A. Cekin. A supervised machine learning algorithm for
424 arrhythmia analysis. In *Computers in Cardiology 1997*, pp. 433–436, 1997. doi: 10.1109/CIC.
425 1997.647926.
- 426 Yaojie Hu and Jin Tian. Neuron dependency graphs: A causal abstraction of neural networks. In
427 *International Conference on Machine Learning*, pp. 9020–9040. PMLR, 2022.
- 428 Biwei Huang, Kun Zhang, Pengtao Xie, Mingming Gong, Eric P Xing, and Clark Glymour. Spe-
429 cific and shared causal relation modeling and mechanism-based clustering. *Advances in Neural*
430 *Information Processing Systems*, 32, 2019.
- 431 Biwei Huang, Kun Zhang, Mingming Gong, and Clark Glymour. Causal discovery from multiple
432 data sets with non-identical variable sets. In *Proceedings of the AAAI Conference on Artificial*
433 *Intelligence*, volume 34, pp. 10153–10161, 2020a.
- 434 Biwei Huang, Kun Zhang, Jiji Zhang, Joseph D Ramsey, Ruben Sanchez-Romero, Clark Glymour,
435 and Bernhard Schölkopf. Causal discovery from heterogeneous/nonstationary data. *J. Mach. Learn.*
436 *Res.*, 21(89):1–53, 2020b.
- 437 Yuchi Ishikawa, Seito Kasai, Yoshimitsu Aoki, and Hirokatsu Kataoka. Alleviating over-segmentation
438 errors by detecting action boundaries. In *Proceedings of the IEEE/CVF winter conference on*
439 *applications of computer vision*, pp. 2322–2331, 2021.
- 440 Dominik Janzing, Joris Mooij, Kun Zhang, Jan Lemeire, Jakob Zscheischler, Povilas Daniušis,
441 Bastian Steudel, and Bernhard Schölkopf. Information-geometric approach to inferring causal
442 directions. *Artificial Intelligence*, 182:1–31, 2012.
- 443 Connor T Jerzak, Fredrik Johansson, and Adel Daoud. Image-based treatment effect heterogeneity.
444 *arXiv preprint arXiv:2206.06417*, 2022.
- 445 Yibo Jiang and Victor Veitch. Invariant and transportable representations for anti-causal domain
446 shifts. *arXiv preprint arXiv:2207.01603*, 2022.
- 447 Markus Kalisch and Peter Bühlman. Estimating high-dimensional directed acyclic graphs with the
448 pc-algorithm. *Journal of Machine Learning Research*, 8(3), 2007.
- 449 Diederik P Kingma and Max Welling. Auto-encoding variational bayes, 2022.
- 450 Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large
451 language models are zero-shot reasoners. *Advances in neural information processing systems*, 35:
452 22199–22213, 2022.
- 453 Insung Kong, Yuha Park, Joonhyuk Jung, Kwonsang Lee, and Yongdai Kim. Covariate balancing
454 using the integral probability metric for causal inference. *arXiv preprint arXiv:2305.13715*, 2023.
- 455 Hilde Kuehne, Ali Arslan, and Thomas Serre. The language of actions: Recovering the syntax and
456 semantics of goal-directed human activities. In *Proceedings of the IEEE conference on computer*
457 *vision and pattern recognition*, pp. 780–787, 2014.
- 458 Min Li, Hui Zhao, Tiquan Gu, and Di Ying. Experiencer-driven and knowledge-aware graph model
459 for emotion–cause pair extraction. *Knowledge-Based Systems*, pp. 110703, 2023a.
- 460 Min Li, Hui Zhao, Tiquan Gu, Di Ying, and Bin Liao. Class imbalance mitigation: A select-then-
461 extract learning framework for emotion-cause pair extraction. *Expert Systems with Applications*,
462 pp. 121386, 2023b.
- 463 Shi-Jie Li, Yazan AbuFarha, Yun Liu, Ming-Ming Cheng, and Juergen Gall. Ms-tcn++: Multi-stage
464 temporal convolutional network for action segmentation. *IEEE transactions on pattern analysis*
465 *and machine intelligence*, 2020.
- 466 Yanran Li, Hui Su, Xiaoyu Shen, Wenjie Li, Ziqiang Cao, and Shuzi Niu. DailyDialog: A manually
467 labelled multi-turn dialogue dataset. In *Proceedings of the Eighth International Joint Conference on*
468 *Natural Language Processing (Volume 1: Long Papers)*, pp. 986–995, Taipei, Taiwan, November
469 2017. Asian Federation of Natural Language Processing. URL <https://aclanthology.org/I17-1099>.

- 471 Sindy Löwe, David Madras, Richard Zemel, and Max Welling. Amortized causal discovery: Learning
472 to infer causal graphs from time-series data. In *Conference on Causal Learning and Reasoning*,
473 pp. 509–525. PMLR, 2022.
- 474 Wang Lu, Wang Wang, Jindong Yidong, and Xing Xie. Towards optimization and model selection for
475 domain generalization: A mixup-guided solution. In *The KDD'23 Workshop on Causal Discovery,*
476 *Prediction and Decision*, pp. 75–97. PMLR, 2023.
- 477 Fangrui Lv, Jian Liang, Shuang Li, Bin Zang, Chi Harold Liu, Ziteng Wang, and Di Liu. Causality
478 inspired representation learning for domain generalization. In *Proceedings of the IEEE/CVF*
479 *Conference on Computer Vision and Pattern Recognition*, pp. 8046–8056, 2022.
- 480 Clare Lyle, Arash Mehrjou, Pascal Notin, Andrew Jesson, Stefan Bauer, Yarin Gal, and Patrick
481 Schwab. Discobax-discovery of optimal intervention sets in genomic experiment design. 2023.
- 482 Sara Magliacane, Thijs Van Ommen, Tom Claassen, Stephan Bongers, Philip Versteeg, and Joris M
483 Mooij. Domain adaptation by using causal inference to predict invariant conditional distributions.
484 *Advances in neural information processing systems*, 31, 2018.
- 485 OpenAI. Gpt-4 technical report, 2023.
- 486 Jose Peña. Factorization of the partial covariance in singly-connected path diagrams. In *Conference*
487 *on Causal Learning and Reasoning*, pp. 814–849. PMLR, 2023.
- 488 Soujanya Poria, Devamanyu Hazarika, Navonil Majumder, Gautam Naik, Erik Cambria, and Rada
489 Mihalcea. MELD: A multimodal multi-party dataset for emotion recognition in conversations.
490 In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp.
491 527–536, Florence, Italy, July 2019. Association for Computational Linguistics. doi: 10.18653/v1/
492 P19-1050. URL <https://aclanthology.org/P19-1050>.
- 493 Soujanya Poria, Navonil Majumder, Devamanyu Hazarika, Deepanway Ghosal, Rishabh Bhard-
494 waj, Samson Yu Bai Jian, Pengfei Hong, Romila Ghosh, Abhinaba Roy, Niyati Chhaya, et al.
495 Recognizing emotion cause in conversations. *Cognitive Computation*, 13:1317–1332, 2021.
- 496 Joseph Ramsey, Madelyn Glymour, Ruben Sanchez-Romero, and Clark Glymour. A million variables
497 and more: the fast greedy equivalence search algorithm for learning high-dimensional graphical
498 causal models, with an application to functional magnetic resonance images. *International journal*
499 *of data science and analytics*, 3:121–129, 2017.
- 500 Fabio De Sousa Ribeiro, Tian Xia, Miguel Monteiro, Nick Pawlowski, and Ben Glocker. High fidelity
501 image counterfactuals with probabilistic causal models. 2023.
- 502 Paul K Rubenstein, Sebastian Weichwald, Stephan Bongers, Joris M Mooij, Dominik Janzing, Moritz
503 Grosse-Wentrup, and Bernhard Schölkopf. Causal consistency of structural equation models. *arXiv*
504 *preprint arXiv:1707.00819*, 2017.
- 505 Andreas WM Sauter, Erman Acar, and Vincent François-Lavet. A meta-reinforcement learning
506 algorithm for causal discovery. In *Conference on Causal Learning and Reasoning*, pp. 602–619.
507 PMLR, 2023.
- 508 Bernhard Schölkopf, Francesco Locatello, Stefan Bauer, Nan Rosemary Ke, Nal Kalchbrenner,
509 Anirudh Goyal, and Yoshua Bengio. Toward causal representation learning. *Proceedings of the*
510 *IEEE*, 109(5):612–634, 2021.
- 511 Weizhou Shen, Siyue Wu, Yunyi Yang, and Xiaojun Quan. Directed acyclic graph network for
512 conversational emotion recognition. In *Proceedings of the 59th Annual Meeting of the Asso-*
513 *ciation for Computational Linguistics and the 11th International Joint Conference on Natural*
514 *Language Processing (Volume 1: Long Papers)*, pp. 1551–1560, Online, August 2021. Asso-
515 ciation for Computational Linguistics. doi: 10.18653/v1/2021.acl-long.123. URL <https://aclanthology.org/2021.acl-long.123>.
- 517 Tao Shi and Shao-Lun Huang. Multiemo: An attention-based correlation-aware multimodal fusion
518 framework for emotion recognition in conversations. In *Proceedings of the 61st Annual Meeting of*
519 *the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 14752–14766, 2023.

- 520 Riccardo Silini and Cristina Masoller. Fast and effective pseudo transfer entropy for bivariate
521 data-driven causal inference. *Scientific reports*, 11(1):8423, 2021.
- 522 Dipika Singhanian, Rahul Rahaman, and Angela Yao. Coarse to fine multi-resolution temporal
523 convolutional network. *arXiv preprint arXiv:2105.10859*, 2021.
- 524 Stephen M Smith, Karla L Miller, Gholamreza Salimi-Khorshidi, Matthew Webster, Christian F
525 Beckmann, Thomas E Nichols, Joseph D Ramsey, and Mark W Woolrich. Network modelling
526 methods for fmri. *Neuroimage*, 54(2):875–891, 2011.
- 527 Sebastian Stein and Stephen J McKenna. Combining embedded accelerometers with computer vision
528 for recognizing food preparation activities. In *Proceedings of the 2013 ACM international joint
529 conference on Pervasive and ubiquitous computing*, pp. 729–738, 2013.
- 530 Jiahui Wang, Zhengyou Wang, Shanna Zhuang, Yaqian Hao, and Hui Wang. Cross-enhancement
531 transformer for action segmentation. *Multimedia Tools and Applications*, pp. 1–14, 2023.
- 532 Yuhao Wang, Vlado Menkovski, Hao Wang, Xin Du, and Mykola Pechenizkiy. Causal discovery
533 from incomplete data: a deep learning approach. *arXiv preprint arXiv:2001.05343*, 2020.
- 534 Ying-Xin Wu, Xiang Wang, An Zhang, Xiangnan He, and Tat-Seng Chua. Discovering invariant
535 rationales for graph neural networks. *arXiv preprint arXiv:2201.12872*, 2022.
- 536 Yue Yu, Jie Chen, Tian Gao, and Mo Yu. Dag-gnn: Dag structure learning with graph neural networks.
537 In *International Conference on Machine Learning*, pp. 7154–7163. PMLR, 2019.
- 538 Zhongqi Yue, Qianru Sun, Xian-Sheng Hua, and Hanwang Zhang. Transporting causal mechanisms
539 for unsupervised domain adaptation. In *Proceedings of the IEEE/CVF International Conference
540 on Computer Vision*, pp. 8599–8608, 2021.
- 541 Sayyed Zahiri and Jinho D. Choi. Emotion Detection on TV Show Transcripts with Sequence-based
542 Convolutional Neural Networks. In *Proceedings of the AAAI Workshop on Affective Content
543 Analysis*, AFFCON’18, pp. 44–51, New Orleans, LA, 2018. URL <https://sites.google.com/view/affcon18>.
- 545 Duzhen Zhang, Feilong Chen, and Xiuyi Chen. Dualgats: Dual graph attention networks for emotion
546 recognition in conversations. In *Proceedings of the 61st Annual Meeting of the Association for
547 Computational Linguistics (Volume 1: Long Papers)*, pp. 7395–7408, 2023a.
- 548 Wenbo Zhang, Tong Wu, Yunlong Wang, Yong Cai, and Hengrui Cai. Towards trustworthy explana-
549 tion: On causal rationalization. *arXiv preprint arXiv:2306.14115*, 2023b.
- 550 Zhuosheng Zhang, Aston Zhang, Mu Li, and Alex Smola. Automatic chain of thought prompting in
551 large language models. *arXiv preprint arXiv:2210.03493*, 2022.
- 552 Weixiang Zhao, Yanyan Zhao, Zhuojun Li, and Bing Qin. Knowledge-bridged causal interaction
553 network for causal emotion entailment. *arXiv preprint arXiv:2212.02995*, 2022a.
- 554 Weixiang Zhao, Yanyan Zhao, and Xin Lu. Cauain: Causal aware interaction network for emotion
555 recognition in conversations. In *Proceedings of the Thirty-First International Joint Conference on
556 Artificial Intelligence, IJCAI*, pp. 4524–4530, 2022b.
- 557 Weixiang Zhao, Yanyan Zhao, Zhuojun Li, and Bing Qin. Knowledge-bridged causal interaction
558 network for causal emotion entailment. In *Proceedings of the AAAI Conference on Artificial
559 Intelligence*, volume 37, pp. 14020–14028, 2023.
- 560 Xun Zheng, Bryon Aragam, Pradeep K Ravikumar, and Eric P Xing. Dags with no tears: Continuous
561 optimization for structure learning. *Advances in neural information processing systems*, 31, 2018.
- 562 Wanqi Zhou, Shujian Yu, and Badong Chen. Causality detection with matrix-based transfer entropy.
563 *Information Sciences*, 613:357–375, 2022.

564 A SUPPLEMENTARY EXPLANATIONS FOR CAUSAL DATA

565 A.1 PRELIMINARIES

566 The joint distribution of samples within the same structure can be represented by any factorization:

$$P(x_1, x_2, \dots, x_s) = \prod_{s=1}^S P(x_s | X_{others}) \quad (4)$$

567 and it can always be consistent with the probability distribution of a certain graph. For instance, in
 568 statistical models, for any variable $x_j \in X_{others}$, there exists an undirected edge between x_j and
 569 x_i , representing the correlation between x_j and x_i . Despite the Markov property, the direction can't
 570 be directly identified by the conditional probability for an undirected edge. However, the causal
 571 model can identify the causal direction between two related variables x_j and x_i by intervention (for
 572 example, $x_j \rightarrow x_i$, $x_i \rightarrow x_j$, $x_j \rightarrow L \rightarrow x_i$, $x_j \leftarrow L \rightarrow x_i$ and so on).

573 **Definition 4** (Intervention). *Interventions are typically represented by the do operator, with the*
 574 *objective of setting the probability of an observed variable equaling a particular state to 1.*

$$do(x_t = t) := P(x_t = t) = 1 \quad (5)$$

575 where t is one of the state which probably exist in the original distribution of variable x_t .

576 When the interventions are integrated with factorization (Equation 4), the factorization satisfies causal
 577 factorization, i.e., the factorization can be converted into its corresponding Structural Causal Model
 578 (SCM).

579 **Definition 5** (Structural Causal Model). *An SCM is a 4-tuple $\langle X, \mathcal{F}, U, \mathbb{P} \rangle$, where X is the entire set*
 580 *of observed variables $X = \{x_i\}_{i=1}^S$. Structural equations $\mathcal{F} = \{f_i\}_{i=1}^S$ are functions that determine*
 581 *causal representation \hat{X} with $\hat{x}_i = f_i(Pa(x_i), u_i)$, where $Pa(x_i) \subseteq X$ represents the parent set of*
 582 *X , $u_i \in U$ represents the i.i.d. noise term. $\mathbb{P}(X)$ is a distribution set over X .*

583 Exact transformation (Rubenstein et al., 2017) or causal abstraction (Beckers & Halpern, 2019) is
 584 a method of judging causal consistency based on interventions and SCMs. The τ transformation
 585 becomes vital for making two causal models equivalent.

586 **Definition 6** (τ -transformation). *Let I_L to be a set of interventions on micro model $SCM_M =$*
 587 *$\langle X_M, \mathcal{F}_M, U_M, \mathbb{P}_M \rangle$. Similarly, let I_N be interventions on macro model $SCM_N =$*
 588 *$\langle X_N, \mathcal{F}_N, U_N, \mathbb{P}_N \rangle$. Let τ be a partial transformation function $\tau : \mathbb{P}_M(X_M) \rightarrow \mathbb{P}_N(X_N)$. Let*
 589 *$\omega : I_M \rightarrow I_N$ be*

$$\tau(\mathbb{P}_M(X_M)) \rightarrow \mathbb{P}_N(X_N) = \omega(I_M) \rightarrow I_N \quad (6)$$

590

591 A.2 EXAMPLES AND CHARACTERISTICS OF CAUSAL DATA

592 **Example 2** (Definite Data). *Arrhythmia Dataset (Guvendir et al., 1997) is a case record dataset from*
 593 *patients with arrhythmias, including 452 samples, and each sample consists of 279 single-value*
 594 *variables (e.g., age, weight, heart rate, etc.). All samples contribute a common causal graph with*
 595 *279 nodes, where the edge value indicates some causal relationship, such as the causal strength of*
 596 *how age affects heart rate.*

597 **Example 3** (Semi-definite Data (Multi-skeleton and Single-value)). *The Netsim dataset (Smith et al.,*
 598 *2011) is a simulated fMRI dataset. Because different activities in brain regions over time imply*
 599 *different categories, a set of records of one patient corresponds to one causal structure. This dataset*
 600 *includes 50 structures and each structure consists of 15 single-value variables that measure the signal*
 601 *strength of 15 brain regions.*

602 **Example 4** (Semi-definite Data (Single-skeleton and Multi-value)). *CMNIST-75sp (Fan et al., 2022)*
 603 *is a graph classification dataset with controllable bias degrees. In this dataset, all researchers*
 604 *concentrate on one causal graph including 4 variables: causal variable C , background variable B ,*
 605 *observed graph G and label Y . C is a part of the MNIST image including multi value of a group of*
 606 *pixels.*

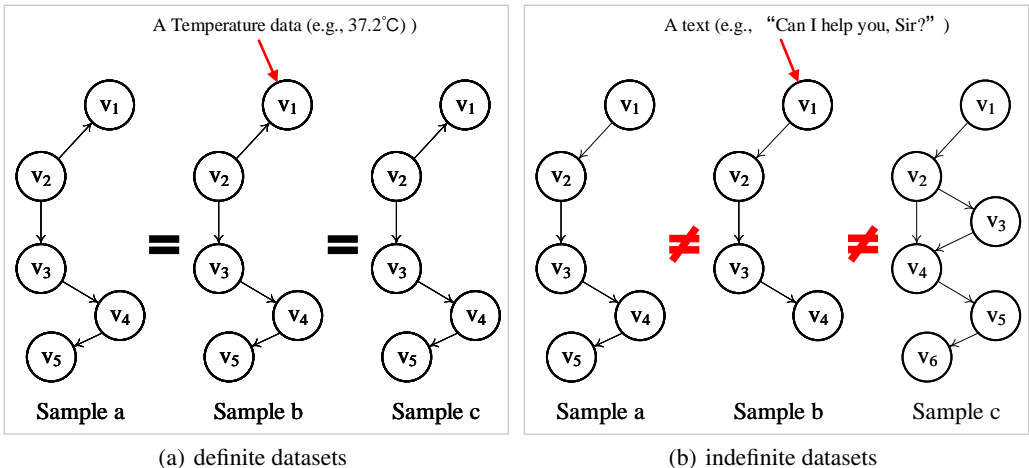


Figure 5: Differences between definite datasets (where $M = 1$ and $D = 1$) and indefinite datasets (where $M > 1$ and $D > 1$). In definite datasets, each sample corresponds to an identical causal structure, implying a single-skeleton trait as the entire dataset involves only a single causal structure. In contrast, indefinite datasets do not possess one causal structure for all samples. For instance, there might be varying numbers of causal variables in samples a, b, and c; the relationship $v_2 \rightarrow v_4$ may be absent in sample b but present in sample c. Furthermore, the causal variables in definite and indefinite datasets also differ. For example, in definite datasets, the causal variable v_1 might represent body temperature, with a causal representation of 37.3 in sample a, 37.1 in sample b and 36.8 in sample c, while v_2 might symbolize blood pressure, with a causal representation of 118, 127, and 135 separately; they are both single-value data. However, in indefinite data, within a dialogue dataset, the causal variable v_1 might be an utterance (“Can I help you, Sir?” in sample a, “Nice to meet you !” in sample b, and “What’s the matter with you” in sample c) with its causal representation being a 768-dimension word embedding in deep model, and v_2 might be a responding utterance to v_1 . In a video dataset, v_1 might denote a segment representing a particular action or event, with its causal representation as the corresponding optical flow, and v_2 might be another segment triggered by v_1 .

607 **Example 5 (Indefinite Data).** *IEM Dataset (Busso et al., 2008)* is a conversation record dataset
 608 with each sample including a dialogue between two speakers. All 100 samples are assigned into 26
 609 graphs (i.e., 26 skeletons) based on the speaker identifies and turns and each sample consists of 5-24
 610 variables where each variable is an utterance represented by embeddings $\in \mathbb{R}^{1 \times 768}$ or $\mathbb{R}^{1 \times 1024}$ in
 611 prevalent pretrained language models.

612 We aim to illustrate the relationships among three data paradigms through Examples 2,3,4,5 and
 613 Figure 5, focusing particularly on the number of skeletons (single or multi-skeleton) and the dimension
 614 of causal representations (single or multi-value).

615 **single or multi-skeleton:** Compared to single-skeleton data ($M = 1$), multi-skeleton data ($M > 1$)
 616 lacks discrimination about which samples belong to the same causal structure. Therefore, it requires
 617 algorithms capable of distinguishing between different causal structures or clustering similar samples.
 618 Simultaneously, multi-skeleton data often have trouble in low sample utilization since samples from
 619 other skeletons contribute nothing when identifying a specific causal structure. Consequently, the
 620 pathways focusing on single-skeleton and multi-skeleton data are different.

621 **single or multi-value:** multi-value data ($D > 1$) often facilitate the quantification by deep repre-
 622 sentation, such as text \rightarrow embeddings, image \rightarrow matrices, audio \rightarrow spectrum map, and video \rightarrow
 623 optical flow, as exemplified in our Figure 5. Compared to single-value data ($D = 1$), it involves
 624 more complex environments. The statistical advantages of single-value data are more significant,
 625 such as computing independence between two single-value variables. On the contrary, determining
 626 such “independence” among multi-value representations is challenging, often approximated through
 627 algorithms like cosine similarity. In Structural Causal Models (SCMs), one can assume that the
 628 noise of single-value data follows a specific distribution, but in multi-value data, the noise items

629 are multi-value and interdependent among dimensions, causing many traditional causal discovery
630 methods to make no efforts with multi-value data.

631 A.3 THE DISTINCTIONS AMONG THREE DATA PARADIGMS

632 Specifically, we employ the theory illustrated in Schölkopf et al. (2021) to explicate why the skeleton
633 (**M**) and variable dimension (**D**) are pivotal in capturing differences in causal discovery algorithms.
634 According to the assumption in Schölkopf et al. (2021), the domain of causal variables \mathcal{X} is projected
635 onto the domain of causal representations $\hat{\mathcal{X}}$ via the encoder p_φ and decoder q_θ , showcasing the
636 causal mechanism in structural equations:

$$\hat{x}_i = f_i(Pa_i, U_i) \quad (7)$$

637 where Pa_i represent the parent node set of x_i . For instance, $p_\varphi : U = (1 - A)X$ and $q_\theta : \hat{X} =$
638 $(1 - A)^{-1}U$. Without prior knowledge, there exist two pathways to recover the causal model: 1)
639 Given a fixed causal structure and known causal representation, the causal strength can be estimated
640 by the statistical strength observable in the samples. 2) If encoder and decoder are feasible, optimal
641 solutions of the causal model can be achieved by minimizing the reconstruction loss $p_\varphi \circ f \circ q_\theta$. Here
642 we would like to delimit the solvability of this process for different combinations via $M=1, M>1,$
643 $D=1,$ and $D>1$.

644 **For a single-skeleton model (M=1):** When the causal structure is fixed, causal strengths f can be
645 calculated. If the causal representation is single-value ($D=1$), the causal structure can be determined
646 without the encoder p_φ or decoder q_θ . The reconstruction loss in this case is f . However, for
647 multi-value data ($D > 1$), in the reconstruction loss function $p_\varphi \circ f \circ q_\theta$, f represents the being
648 determined part.

649 **For a multi-skeleton model (M > 1):** The multi-skeleton data induce uncertainty in causal structures,
650 unclear of which samples correspond to the same causal structure and therefore making causal
651 strengths f unsolvable directly. However, under single-value ($D=1$) condition without generated
652 representation, the precision of clustering is guaranteed. We can approach by first clustering the
653 samples, and then separate the problem to several tasks of definite data problem-solving ($M=1,$
654 $D=1$). In this regard, reconstruction loss amounts to $\{f_m\}_{m=1}^M$, representing the set containing
655 each sub-task’s f_m . Reconstruction loss can be regarded as a multi-task optimization problem,
656 $\alpha_1 f_1 + \alpha_2 f_2 + \dots + \alpha_M f_M$, where α_m is the weights of the sample quantity per structure. The
657 worst-case scenario arises with multi-value data ($D>1$), only able to attain an approximate encoder
658 $\tilde{p}_\varphi = p_\varphi \circ f_m$, which results in a final reconstruction loss of $\tilde{p}_\varphi \circ q_\theta$. Causal strength f_m comprises
659 an unassigned part.

660 In summary, for definite data ($M = 1, D = 1$), it suffices to identify the causal strength between
661 any two causal variables under a certain causal structure. Semi-definite data addresses the problem
662 of discriminating multi-skeleton structures and encoding multi-value variables separately. As for
663 indefinite data, in the absence of additional assumptions, causal discovery in such datasets presents
664 an ill-posed problem, given it requires both variable encoding and resolving structure discernibility.

665 B PROOF OF THEOREM 1

666 Existing work (Hu & Tian, 2022) proposed that two causal models satisfying causal abstraction are
667 consistent. Accordingly, we employ causal abstraction as a mediator to prove Theorem 1. In other
668 words, it suffices to verify that the concepts of equivalent distribution sets and equivalent strength
669 sets are conditionally transformable.

670 First, we would like to prove that equivalent distribution sets \rightarrow equivalent strength sets.

671 According to the causal factorization, $P(x_s | X_{others})$ satisfies $x_s = f_s(Pa(x_s, u_s))$. When we
672 assume two causal model within two variables: $\mathcal{U} : x_i \rightarrow x_j$ and $\mathcal{V} : y_i \rightarrow y_j$, we would like to
673 adopt SCM to represent the equivalent distribution $P_{x_j}(do(x_i)) = P_{y_j}(\omega(do(x_i)))$:

$$f_{i,j}(u_{x_i}) + u_{x_j} = g_{i,j}(u_{y_i}) + u_{y_j} \quad (8)$$

674 Because of $P_{x_j}(do(x_j)) = P_{y_j}(\omega(do(x_j)))$ and $P_{x_i}(do(x_i)) = P_{y_i}(\omega(do(x_i)))$, Equation 8 can be
675 written as:

$$f_{i,j} = g_{i,j} \quad (9)$$

676 According to the causal partial order, $P_{x_j}(do(x_j)) \preceq P_{x_j}(do(x_j, x_k)), (x_j \preceq x_k)$, hence:

$$f_{j,k} = g_{j,k} \quad (10)$$

677 When we convert any causal factorization into a chain of ancestral relationships through additive
678 noise formulas, it is possible to find a corresponding $g_{i,j}$ that equals $f_{i,j}$ for any step in the causal
679 chain. Finally, we can infer that if the distribution sets $P_X = P_Y$, then the strength sets $F_X = F_Y$.
680 Conversely, it can also be proven that if the strength sets $F_X = F_Y$, then the distribution sets
681 $P_X = P_Y$.

682 C DETAILS OF TWO EXAMPLES FOR IMPLEMENTATION

683 C.1 EXAMPLES IN SUPERVISED SPECIALIZED MODEL (SSM)

684 C.1.1 PREVALENT PROBABILISTIC MODEL

685 Many variational models for causal discovery, including linear SEM variational model (Yu et al.,
686 2019), autoregressive (Wang et al., 2020) and recently substitute of noise (Chen et al., 2023b), can
687 be encapsulated by a probabilistic framework:

688 1. Construct a Linear Structural Equation Model (SEM) to displace SCM. Specifically, let A
689 $\in \mathbb{R}^{N \times N}$ be the adjacency matrix, and N stands for the number of variables. $X \in \mathbb{R}^{N \times D}$
690 is a sample of N variables.

$$X = AX + E \quad (11)$$

691 where $E \in \mathbb{R}^{N \times D}$ is the matrix of independent noise ϵ_{x_n} , A represents the causal strength
692 from all variables to one observed variable.

693 2. Build a pair of Autoregression SEMs:

$$E = (I - A)X \quad (12)$$

$$X = (I - A)^{-1}E \quad (13)$$

695 Equation 13 describes a general form as a decoder of a generation model that takes noise E
696 as input and returns X as results and Equation 12 describes the corresponding encoder.

697 3. Considering a specification of noise (E) distribution sampling $\{X_s\}_{s=1}^S$ in definite data,
698 Equation 13 can be written by a maximization of log-evidence:

$$\frac{1}{S} \sum_{s=1}^S \log p(X_s) = \frac{1}{S} \sum_{s=1}^S \log \int p(X_s|E)p(E)dE \quad (14)$$

699 Continuing the theory of variational Bayes, we regard E as the latent variable in variational
700 autoencoder (VAE) (Kingma & Welling, 2022) and use variational posterior $q(E|X)$ to
701 approximate the intractable posterior $p(E|X)$, thus the evidence lower bound (ELBO) reads:

$$\mathcal{L}_{ELBO}^s = -KL(q(E|X_s)||p(E)) + E_{q(E|X_s)}[\log p(X_s|E)] \quad (15)$$

702 C.1.2 OUR SUPERVISED IMPLEMENTATION

703 Taking an example from the variational probabilistic framework mentioned in Appendix C.1.1, which
704 has become a popular choice, we simplify consider Equation 12 as the encoder p_φ and Equation 13
705 as the decoder q_θ . The model \mathcal{U} is generated by encoder: $\mathcal{G}_{s,m} = h(X_{s,m}, \varphi)$ and the model \mathcal{V}
706 is generated by encoder and decoder: $\hat{X}_{s,m} = h(X_{s,m}, \varphi, \theta)$. The adjacency matrix A represents
707 the influence between the observed variables. For example, $A_{i,j} \in [0, 1]$ describes the strength of

708 how the variable x_j influences the variable x_i . The ‘augment’ measure corresponding to $do_g(i, j)$ is
 709 defined as follows:

$$A_{m,n} = \begin{cases} 0, & (m = i \text{ or } m = j) \\ A_{m,n}, & \text{else} \end{cases} \quad (16)$$

710 The i -th and j -th row represent the influence of all parent nodes on variables i and j , respectively. After
 711 the ‘augment’ measures, the corresponding i -th and j -th rows in the adjacency matrix $W = (1 - A)^{-1}$
 712 of the decoder only contain the two non-zero terms, u_i and u_j . For the consistency check, we adopt an
 713 easily computable similarity matrix $Sim \in \mathbb{R}^{N \times N}$, where $Sim_{i,j}$ represents the similarity between
 714 variables x_i and x_j . The similarity matrices of models \mathcal{U} and \mathcal{V} are respectively generated from
 715 distinct resources: model \mathcal{U} and its similarity matrix are computed from causal representation, while
 716 model \mathcal{V} and its similarity matrix are derived from the causal structure.

717 The similarity matrix of \mathcal{U} is computed via cosine similarity:

$$Sim_{m,n}^r = \text{cossim}(\hat{x}_m^{do_g(i,j)}, \hat{x}_n^{do_g(i,j)}) * \text{Mask} \quad (17)$$

718 where $\hat{x}_m^{do_g(i,j)}$ and $\hat{x}_n^{do_g(i,j)}$ are the causal representation under the view $do_g(i, j)$, Mask stands for
 719 a lower triangular matrix. $\text{Mask}_{i,j} = 0$ when $j > i$ and $\text{Mask}_{i,j} = 1$ when $j \leq i$.

720 The similarity matrix of \mathcal{V} , its can be obtained directly from the causal strength matrix:

$$Sim_{m,n}^s = \begin{cases} 0, & (n \leq i, m = i \text{ or } n \leq j, m = j) \\ P(x_m^{do_g(i,j)} | x_n^{do_g(i,j)}), & \text{else} \end{cases} \quad (18)$$

721 Finally, the Consistency check module needs to measure these two similarity matrix within specifically
 722 MSE loss function we adopted:

$$\text{Loss} = \text{MSE}(Sim_{m,n}^r, Sim_{m,n}^s) \quad (19)$$

723 C.2 EXAMPLES FOR LARGE LANGUAGE MODEL (LLM)

724 Large language models (LLMs) demonstrate superior performance across a variety of text tasks, particu-
 725 larly showing “natural” level in human-machine conversations based on instructions. Considering
 726 dialogue as a typical data type of indefinite data, we have implemented an example where LLMs can
 727 identify causal relationships between utterances through prompt instructions.

728 Specifically, we treat each dialogue as a sample, where an utterance is regarded as a causal variable.
 729 That is, for a dialogue $D = \{Utt_1, Utt_2, \dots, Utt_N\}$, where Utt_i represents i -th utterance, N is the
 730 number of causal variables. We adopt an iterative prompt instruction wherein the LLM’s predictions
 731 are gradually corrected through feedback instructions. The causal structure predicted by the LLMs
 732 falls under model \mathcal{V} , while the causal representation computed through utterance representation
 733 belongs to model \mathcal{U} . Our iterative prompt instruction is as follows: firstly, the LLM is instructed to
 734 answer the complete the causal structure (model \mathcal{V}); secondly, the utterances’ causal representation
 735 (model \mathcal{U}) is obtained either via pre-trained models or LLMs; finally, the differences in similarity
 736 matrices of model \mathcal{V} and model \mathcal{U} are fed back, and the LLM is instructed to execute the prior steps
 737 under the acknowledgement of this difference until no difference exists between the two similarity
 738 matrices. What distinguishes our LLM implementation to the SSM implementation is that model \mathcal{U}
 739 is unlearnable and fixed in LLM, whereas the SSM permits both \mathcal{V} and model \mathcal{U} to be learnable.

740 Eventually, this leads to the discovery of the correct causal relationship. The steps are as follows:

741 **Step 1 (Prediction Causal Relationship):** Calculating the causal relationship between any two
 742 utterances within a given dialogue, the prompt instruction employs input accompanied by an example.
 743 The specific text is as follows:

“You are assuming the role of a researcher capable of distinguishing between causation and correlation, charged with the task of recognizing the causal relationships among individual utterances within a given dialogue. We prescribe that the judgment of causation between two utterances is based on whether the former is the intended target of the latter’s response. Whereas, correlation is gauged on whether the two share similar topics or vocabulary. The following is an example:

Example:

Dialogue:

- ‘1. Hazel drank too much champagne at the party.
2. Oh my goodness! That sounds like quite an eventful party.
3. Well, drinking too much alcohol can have many negative effects on the body.
4. Oh no, I can imagine Hazel waking up with a massive headache tomorrow.’

Question 1: Is there a causal relationship from utterance 1 to utterance 2?

Answer 1: Yes.

Question 2: Is there a causal relationship from utterance 1 to utterance 3?

Answer 2: Yes.

Question 3: Is there a causal relationship from utterance 1 to utterance 4?

Answer 3: Yes.

Question 4: Is there a causal relationship from utterance 2 to utterance 3?

Answer 4: No.

Question 5: Is there a causal relationship from utterance 2 to utterance 4?

Answer 5: No.

Question 6: Is there a causal relationship from utterance 3 to utterance 4?

Answer 6: Yes.

Given the above example, with its associated questions and answers, consider the following dialogue:

Dialogue:

- ‘1. Charlotte has no idea how to avoid massive estate taxes.
2. Estate taxes are a topic of concern for many people in various countries.
3. So, does anyone else have any knowledge or ideas on how to reduce estate taxes?
4. Oh, that reminds me of a story about my uncle.’

Question 1: Is there a causal relationship from utterance 1 to utterance 2?

Question 2: Is there a causal relationship from utterance 1 to utterance 3?

Question 3: Is there a causal relationship from utterance 1 to utterance 4?

Question 4: Is there a causal relationship from utterance 2 to utterance 3?

Question 5: Is there a causal relationship from utterance 2 to utterance 4?

Question 6: Is there a causal relationship from utterance 3 to utterance 4?”

744 **Step 2 (Calculating Sim^s):** Based on the results from the Step 1, the similarity matrix Sim^s of the
 745 causal structures under different views is computed according to Equation 18. This similarity matrix
 746 includes only two binary values: 0 and 1. Note that intervention here, according to Definition 2,
 747 involves the direct elimination of intervened utterances. Taking $do_g(1, 2)$ as an example, neither the

748 first nor the second utterance is included in the input. Consequently, the resulting similarity matrix,
749 $Sim \in \mathbb{R}^{(N-2) \times (N-2)}$. This applies to both Sim^r and Sim^s .

750 **Step 3 (Calculating Sim^r):** In calculating the similarity matrix Sim^r for causal representations,
751 we explored two distinct computational methods. 1), computation is conducted through LLMs. The
752 prompt text used is similar to that in Step 1, but substitutes the question “*Is there a causal relationship*
753 *from utterance A to utterance B?*” with “*Is there a correlation relationship between utterance A and*
754 *utterance B?*”. 2), using a pre-trained model such as RoBERTa, we calculate the deep representations
755 of two utterances and then compute their similarity using cosine similarity via Equation 17.

756 **Step 4 (Inconsistency Feedback):** We compare the each pair of similarity matrices Sim^r and Sim^s
757 obtained from differing views. If a condition occurs where $Sim_{i,j}^r = 1$ while $Sim_{i,j}^s = 0$, it can be
758 inferred that “*there is no common cause between the i-th utterance and the j-th utterance, and the i-th*
759 *utterance should not have a causal relationship to the j-th utterance.*” On the contrary, if a situation
760 arises where $Sim_{i,j}^r = 0$ and $Sim_{i,j}^s = 1$, we can assert that “*there is a common cause between the*
761 *i-th utterance and the j-th utterance, and the i-th utterance should have a causal relationship to the*
762 *j-th utterance.*” If ‘i’ refers to the first utterance, no response will be given to the clause relevant to
763 the ‘common cause.’ An example of this prompt instruction is as follows:

“*After verification, there is no common cause between the second utterance and the third utterance, and the second utterance should not have a causal relationship with the third utterance, and there is no common cause between the third utterance and the fourth utterance, and the third utterance should not have a causal relationship with the fourth utterance. Please re-answer based on these circumstances.*”

764 **Recursive Process:** The iterative algorithm is summarized in Algorithm 1. Steps 2 to 4 will
765 continuously loop. The end condition is reached once all instances of Sim^r and Sim^s across all
766 views are identical. The causal relationship output by the LLM during the final loop represents the
767 final results. Step 3 represents the ‘Augment’ module mentioned in Section 4.1 while Step 4 embodies
768 the ‘Consistency check’. The overall objective of the instruction is to enable the LLM to identify
769 causal relationships between utterances without causal labels.

770 Compared to the SSM example in Appendix C.1, the LLM-based example relies significantly on the
771 accuracy of Sim^r . Therefore, it often leads to a failure in achievement of desired causal relationships.
772 (We show the experiment results in Figure 7)

Algorithm 1 Iterative Prompt Instruction

Require: A dialogue text $D = \{Utt_1, Utt_2, \dots, Utt_N\}$, a set of matrices $Sim^s_Set = \{Sim_{i,j}^s, i, j \in N \text{ and } i \leq j\} = \emptyset$, a set of matrices $Sim^r_Set = \{Sim_{i,j}^r, i, j \in N \text{ and } i \leq j\} = \emptyset$, and *input_prompt* as shown in Step 1.

Ensure: Causal relation adjacency matrix $C \in \mathbb{R}^{N \times N}$ in where Sim^s_Set is consistent with Sim^r_Set (Both $\neq \emptyset$).

```

procedure INTERVENTION ( $I_* = \{do_{i,j}\}$ ) ( $i, j \in N$  and  $i \leq j$ )
while  $Sim^s\_Set = Sim^r\_Set \neq \emptyset$  do
  Predict  $C$  via LLM according to input_prompt.
  for each view  $do_{i,j}$  in  $I_*$  do
    Calculate  $Sim_{i,j}^s$  via Step 2.
     $Sim^s\_Set = Sim^s\_Set \cup Sim_{i,j}^s$ 
    Calculate  $Sim_{i,j}^r$  via Step 3.
     $Sim^r\_Set = Sim^r\_Set \cup Sim_{i,j}^r$ 
  end for
  Replace the input_prompt with the results about inconsistent pairs between  $Sim^s\_Set$  and  $Sim^r\_Set$  via Step 4.
end while
return  $C$ 

```

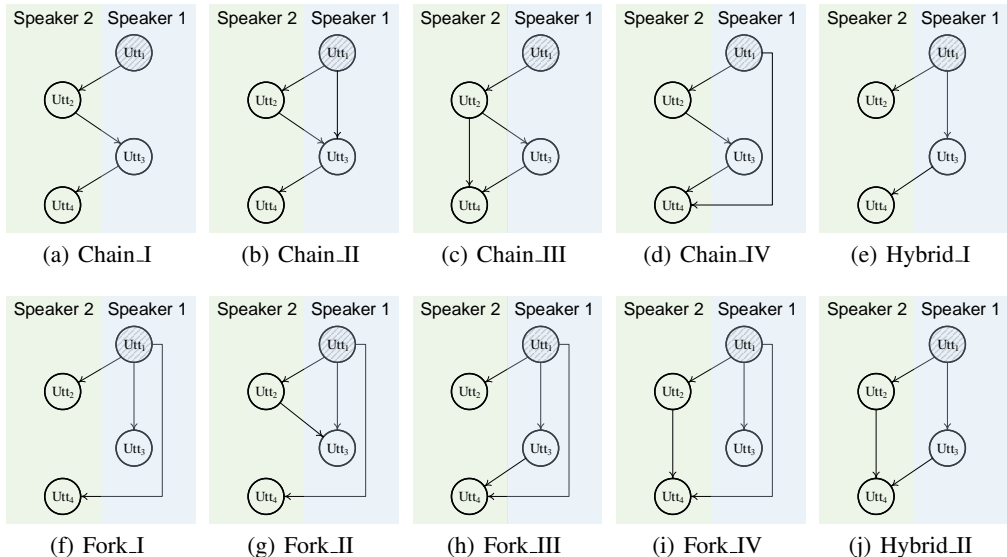


Figure 6: 10 skeletons (structures) in *Causalogue* Dataset.

773 D DETAILS ABOUT *Causalogue* DATASET

774 D.1 ATTRIBUTES

775 **Causal Variables:** We treat each dialogue as a sample, comprised of 4 utterances, which we define
 776 as 4 causal variables. Further, the first and third utterances originate from the same speaker, defined
 777 as *speaker1*. Similarly, the second and fourth utterances are from another individual, defined as
 778 *speaker2*.

779 **Skeletons:** We have designed 10 types of causal skeletons (structures) in the dataset as shown in
 780 Figure 6, listed as follows:

781 **Chain_I:** This is the most basic chain structure, serving as the prototype for Chain_II-IV models. It
 782 contains three causal relationships: $Utt_1 \rightarrow Utt_2$, $Utt_2 \rightarrow Utt_3$, and $Utt_3 \rightarrow Utt_4$, representing
 783 dialogues where two speakers interact sequentially.

784 **Chain_II:** Building upon Chain_I, this model includes an additional causal relationship from $Utt_1 \rightarrow$
 785 Utt_3 , indicating that Utt_3 considers not just the effect from Utt_2 but also from Utt_1 .

786 **Chain_III:** Building upon Chain_I, this model introduces a causal link $Utt_2 \rightarrow Utt_4$, suggesting Utt_4
 787 takes into account the effect of both Utt_3 and Utt_2 .

788 **Chain_IV:** Building upon Chain_I, this model creates an additional causal connection $Utt_1 \rightarrow Utt_4$,
 789 indicating Utt_4 considers the effects from both Utt_3 and Utt_1 —the two utterances by the one
 790 speaker.

791 **Fork_I:** This is the most basic fork structure, serving as the prototype for Fork_II-IV models. It
 792 includes three causal relationships: $Utt_1 \rightarrow Utt_2$, $Utt_1 \rightarrow Utt_3$, and $Utt_1 \rightarrow Utt_4$, representing
 793 situations where two speakers alternately respond to Utt_1 with different independent replies.

794 **Fork_II:** Building upon Fork_I, this model adds a causal relationship $Utt_2 \rightarrow Utt_3$, representing
 795 Utt_3 's response to not only Utt_1 but also Utt_2 .

796 **Fork_III:** Building upon Fork_I, this model introduces a causal link $Utt_3 \rightarrow Utt_4$, signifying that
 797 Utt_4 is not merely an independent response to Utt_1 but a combined reply to both Utt_1 and Utt_3 .

798 **Fork_IV:** Building upon Fork_I, this model incorporates a new causal relationship $Utt_2 \rightarrow Utt_4$,
 799 indicating that Utt_4 responds to both Utt_2 and Utt_1 .

800 Hybrid.I: A combination of the chain and fork structures, where the chain structure runs $Utt_1 \rightarrow$
 801 $Utt_3 \rightarrow Utt_4$, and the fork structure is $Utt_2 \leftarrow Utt_1 \rightarrow Utt_3$.

802 Hybrid.II: On the basis of Hybrid.I, this introduces an additional chain structure, $Utt_1 \rightarrow Utt_2 \rightarrow$
 803 Utt_4 . This model also results in a collider structure with $Utt_2 \rightarrow Utt_4 \leftarrow Utt_3$.

804 **Sample:** We consider a dialogue as a sample, with each sample comprising 4 utterances representing
 805 4 causal variables. Each sample corresponds to one of the 10 causal skeletons outlined above,
 806 annotating whether a causal relationship exists between any two utterances. Due to Hypothesis 1, our
 807 labels only consider forward-causal relationships. An example of a Chain.III sample is shown as
 808 follows:

“causal_type”: “Chain.III”,
 “clause”: { “1”: “Your bill is 19.”, “2”: “Before I pay the bill, I have to express my dissatisfaction with the service I received tonight.”, “3”: “I’m so sorry to hear that but I don’t know what happened.”, “4”: “Specifically, It’s understandable to feel frustrated when something unexpected happens like spilling red wine on your clothes.” },
 “dia_id”: 1,
 “label”: { “1”: “0,0,0,0”, “2”: “1,0,0,0”, “3”: “0,1,0,0”, “4”: “0,1,1,0” }

809 In the given example, the Utt_4 serves as a response to the Utt_3 , while simultaneously attach to the
 810 speaker’s Utt_2 —thereby rendering both the Utt_2 and Utt_3 as causes to the Utt_4 . Indeed, during the
 811 generation process of the Utt_4 , we made sure to inform GPT-4 of the existence of Utt_2 and Utt_3 .

812 D.2 CREATION PROCESS

813 We utilized the API interface of GPT-4² to defined the following variables: “role”, which has
 814 three types - “system”, “user”, and “assistant”. Here, “system” represents the background or a
 815 prior settings, while “user” and “assistant” are defined as speakers with two different identities.
 816 Additionally, the first utterance is pre-set. Hence, creating a dialogue requires a given combination:
 817 a fixed *first_utterance*, a specified *system* information, and a setting which previous utterances
 818 are considered. We have a total of 149 *first_utterance* options, and there are as many as 278,867
 819 combinations of *first_utterance* and *system* settings (our final samples only number in the 1638, to
 820 preserve the diversity and distinctiveness of our dialogues). What follows is an example of generating
 821 the third utterance in the skeleton of ChainII:

{ “role”: “system”, “content”: “You are Peter, you have promised to go to a Chinese Opera with your daughter, so you want to have dinner with your friends in next Sunday.” }
 { “role”: “assistant”, “content”: “Yes. Sunday sounds fine. What time?” (pre-set Utt.1) }
 { “role”: “user”, “content”: Utt.2 }

822 Upon creation, the samples are initially auto-annotated based on their designed labels, and then
 823 manually verified to ensure their validity. Our manual verification employed two annotators, who
 824 demonstrated proficient English understanding and communication skills, possessing sufficient
 825 knowledge about causality. The annotation consistency between these two annotators was tested
 826 through 833 samples, achieving a kappa coefficient of 0.92.

827 During the annotation process, if a sample was labelled differently by the two annotators, that sample
 828 was considered to possess an ambiguous causal relationship and thus was excluded from the final
 829 dataset. Only samples that were consistently labelled by both annotators were ultimately accepted.

830 Furthermore, to guarantee the freedom of manual annotation, we allowed the annotators to label
 831 structures that fell outside the predefined 10 causal structures. Specifically, we only requested
 832 annotators to judge whether any two utterances (satisfying Hypothesis 1) have a causal relationship,
 833 allowing them some discretion, which inevitably produced samples not belonging to the 10 causal
 834 structures. We classified these as the “Other” category.

²<https://platform.openai.com/docs/models/gpt-4>

835 The accuracy of labels was significantly improved after the manual annotation process. However,
 836 considering that the unverified samples might be utilized for other research areas, such as the ability
 837 of LLMs to focus on context, we have released two versions of the datasets, as demonstrated in
 838 Table 1. “Small” signifies samples that have been manually checked as correctly labeled, while
 839 “large” refers to all samples generated by GPT-4 without manual verification. We do not recommend
 840 considering the “large” version when undertaking causality-related work. Likewise, we have not
 841 taken it into our experiments.

842 E DETAILS ABOUT DATASETS, METRICS, BASELINES, AND IMPLEMENTATION

843 E.1 DATASETS

844 The *Causalogue* dataset has already been discussed in
 845 Appendix D. Hence, this section primarily focuses on
 846 the remaining real-world datasets. Their data splits
 847 and specific N -folds validation setups for SSM are
 848 exhibited in Table 5. Among them, RECCON, DD,
 849 MELD, EmoryNLP and IEM are text datasets, and
 850 GTEA, 50salads, and Breakfast are video datasets.
 851 As for LLM, we only randomly select 400 samples
 852 from *Causalogue* and RECCON datasets, respec-
 853 tively. Their overviews and prevalent metrics are
 854 detailed below.

Table 5: Statistics on Datasets

| Dataset | Train | Valid | Test | Folds |
|------------|-------|-------|------|-------|
| Causalogue | 1338 | 100 | 200 | 10 |
| RECCON | 833 | 47 | 225 | 10 |
| DD | 11118 | 1000 | 1000 | 5 |
| MELD | 1038 | 114 | 280 | 5 |
| EmoryNLP | 713 | 99 | 85 | 5 |
| IEM | 100 | 20 | 31 | 5 |
| GTEA | 19 | 2 | 7 | 10 |
| 50salads | 36 | 4 | 10 | 10 |
| Breakfast | 1314 | 146 | 252 | 10 |

855 E.1.1 EMOTION-CAUSE 856 PAIR EXTRACTION (ECPE) TASK

857 **RECCON** (Poria et al., 2021): The first dataset for emotion cause recognition of conversation
 858 including RECCON-DD and RECCON-IE (emulating an out-of-distribution generalization test).
 859 RECCON-DD includes 5380 labeled ECPs and 5 cause spans (*no-context*, *inter-personal*, *self-*
 860 *contagion*, *hybrid*, and *latent*).

861 E.1.2 EMOTION RECOGNITION IN CONVERSATION (ERC) TASK

862 **DD** (Li et al., 2017): A Human-written dialogs dataset with 7 emotion labels (*neutral*, *happiness*,
 863 *surprise*, *sadness*, *anger*, *disgust*, and *fear*). We follow Chen et al. (2023b) to regard utterance turns
 864 as the speaker turns.

865 **MELD** (Poria et al., 2019): A multimodal ERC dataset with 7 emotion labels as the same as DD.

866 **EmoryNLP** (Zahiri & Choi, 2018): A TV show scripts dataset with 7 emotion labels (*neutral*, *sad*,
 867 *mad*, *scared*, *powerful*, *peaceful*, *joyful*).

868 **IEM** (Busso et al., 2008): A multimodal ERC dataset with 9 emotion labels (*neutral*, *happy*, *sad*,
 869 *angry*, *frustrated*, *excited*, *surprised*, *disappointed*, and *fear*). However, models in ERC field are
 870 often evaluated on samples with the first six emotions due to the too few samples of the latter three
 871 emotions. 20 dialogues for validation set is following (Chen et al., 2023b).

872 E.1.3 TEMPORAL ACTION SEGMENTATION (TAS) TASK

873 **GTEA** (Fathi et al., 2011) Georgia Tech Egocentric Activities is comprised of 28 videos captured
 874 from a first-person perspective. It documents 7 different daily activities performed by 4 test actors,
 875 therefore, the dataset is partitioned into four 4 based on the actors. Each video contains approximately
 876 20 fine-grained instances, with each video divided by action segments as labels.

877 **50salads** (Stein & McKenna, 2013) A cooking dataset includes 50 videos highlighting the complete
 878 process of salad preparation undertaken by 25 people, with each video housing between 9,000 to
 879 18,000 RGB frames and containing 17 action class labels. Each video, named after the complete
 880 process of salad making by an individual, is segregated into 5 groups.

881 **Breakfast** A cooking action dataset consists of 10 cooking activities performed by 52 different actors
 882 at various kitchen locations. It encompasses 1,989 videos and offers over 77 hours of content. Each
 883 video is characterized by a sub-cooking activity accomplished by an actor; the complete preparation
 884 process comprises 20-30 such action segments. As the largest among the mentioned datasets, it is
 885 divided into 4 groups.

886 E.2 EVALUATION METRICS

887 E.2.1 CAUSAL DISCOVERY (CD) TASK

888 The CD task was evaluated on the *Causalogue* dataset, a brand new dataset released by us. In our
 889 experiments, we endeavored to assess three outcomes: the accuracy of causal graphs, the accuracy of
 890 causal representations, and the consistency between causal graphs and representations. Consequently,
 891 we employed AUROC and Hamming Distance (HD) to measure causal graphs, AUROC and F1
 892 scores for causal representation evaluation, and MSE and 1–AUROC for measuring the distance of
 893 inconsistencies. These metrics are common and well-accepted. Simultaneously, for each outcome,
 894 we ensured two different metrics to comprehensively evaluate the performance.

895 E.2.2 EMOTION-CAUSE PAIR EXTRACTION (ECPE) TASK

896 We continue to employ the F1 score as the evaluation metric, as initially proposed in Poria et al.
 897 (2021). This metric is broadly accepted and utilized in current research works (Li et al., 2023a;b).

898 E.2.3 EMOTION RECOGNITION IN CONVERSATION (ERC) TASK

899 Similarly to ECPE task, We continue to employ the F1 score as the evaluation metric, as initially
 900 proposed in Shen et al. (2021). This metric is broadly accepted and utilized in current research
 901 works (Chen et al., 2023b; Zhang et al., 2023a; Shi & Huang, 2023).

902 E.2.4 TEMPORAL ACTION SEGMENTATION (TAS) TASK

903 Commonly used metrics include frame-level accuracy (Acc), segmental edit distance (Edit), and
 904 segmental F1 scores with different overlapping threshold k (F1@ k) ($k = \{10, 25, 50\}$). Moreover, to
 905 evaluate the causal consistency of the segmentation results, we proposed an additional causal edit
 906 distance (C-Dis) to measure the dissimilarity between the adjacency matrix and the ground truth. For
 907 the final segmentation results, we constructed causal adjacency matrices $\hat{C} \in \mathbb{R}^{T \times T}$ and ground truth
 908 matrices $C \in \mathbb{R}^{T \times T}$, based on the constraints in consistent mapping condition and calculated the
 909 dissimilarity between them.

$$C - Dis := num(\hat{C}_{i,j} \neq C_{i,j}) \text{ for } i, j = 1, 2, \dots, T \quad (20)$$

910 A lower causal edit distance indicates that the causal relationship at the frame-level has less dissimi-
 911 larity with the ground truth, demonstrating stronger learning ability with causal representation in the
 912 model, and hence a higher level of causal consistency in the segmentation results.

913 E.3 BASELINES

914 E.3.1 BASELINES ON CD TASK

915 SSM

916 **ACD**: leverages shared dynamics to learn to infer causal relationships from multi-skeleton time-series
 917 data via a single, amortized model.

918 **DAG-GNN**: leverages SCM to construct a gnn-based variational model adopting independent noise
 919 E as latent variable.

920 **ACCD**: discover causal relationships in multi-value data via designing a common skeleton and
 921 generating a substitute for independent noise.

922 **biCD**:proposes a dynamic variational inference model leveraging the causal strength instead of
 923 independent noise as the latent variable to construct ELBO for indefinite data.

924 **DisC**: designs a new method for intervention in deep models, combining causal patterns with different
925 shortcuts to achieve the goal of intervention in causal nodes.

926 **DIR**:distinguishes between positive and negative samples after intervention by designing a dynamic
927 loss function, Similar to the DisC thereby effectively intervening in the causal pattern.

928 Since DisC and DIR do not have complete causal discovery models, we incorporate their intervention
929 modules into DAG-GNN (namely, DAG-DisC and DAG-DIR) to demonstrate their intervention
930 strategies for indefinite data.

931 LLM

932 **Zero-shot and Zero-shot-CoT**: proposes a new prompt paradigm like “Let’s think step by step”
933 which is task-agnostic and does not need input-output demonstrations.

934 **Auto-CoT**: proposes an auto prompt method which could cluster the samples first and then select an
935 example for prompt text.

936 E.3.2 BASELINES ON ECPE TASK

937 **EDKA-GM**:introduces an experiencer identification task and present a document-level heterogeneous
938 graph network for capturing global experiencer information to enrich experiencer-based cross-clause
939 association.

940 **seF**:includes two main components: core clause selector and emotion-cause pairs extractor to jointly
941 extract emotion-cause pairs.

942 E.3.3 BASELINES ON ERC TASK

943 **DAG-ERC**: proposes a gnn&rnn-based model to learn the relationship of different speakers and
944 sequential information.

945 **DualGAT**:introduces Dual Graph Attention networks to concurrently consider the complementary
946 aspects of discourse structure and speaker-aware context.

947 **MultiEMO**:proposes a novel attention-based correlation-aware multimodal fusion framework effec-
948 tively integrating multimodal cues by capturing cross-modal mapping relationships across textual,
949 audio and visual modalities.

950 E.3.4 BASELINES ON TAS TASK

951 **MS-TCN**: This is the first method to introduce a multi-stage action segmentation framework based
952 on Temporal Convolutional Networks (TCN). Each stage inputs the initial prediction output from the
953 preceding one for further modification and adjustment.

954 **MS-TCN++**: On the foundation of MS-TCN, this method introduces a dual dilated layer, implement-
955 ing parameter sharing and optimizing segmentation performance.

956 **ASRF**: This method proposes an improved technique based on MS-TCN, composed of a long-term
957 feature extractor and two branches: the Action Segmentation Branch (ASB) and the Boundary
958 Regression Branch (BRB).

959 **CETNet**: Leveraging Transformer, this method connects every layer of convolutional feature mapping
960 in the encoder with a group of features generated through self-attention in the decoder.

961 **C2F**: Utilizing TCN, this method puts forward a novel temporal encoder-decoder to tackle the
962 sequence fragment issue. Its decoder conforms to a coarse-to-fine structure with multi-timescale
963 implicit integration.

964 E.4 IMPLEMENTATION DETAILS

965 E.4.1 THE MODEL ON CD TASK

966 In our Experiments, we utilized RoBERTa-base (768) as our pre-trained model for generating word
 967 embeddings in the SSMs. Throughout the training process, a learning rate of 1e-5 was set, with the
 968 batch size and epochs set to 16 and 50, respectively. The dimension of the hidden layers within the
 969 network was also set to 768. The entire training procedure was conducted on a NVIDIA GEFORCE
 970 RTX 3090 graphics processing unit.

971 E.4.2 THE MODEL ON ECPE TASK AND ERC TASK

972 In the word embedding, we adopt the affect-based pre-trained features³ proposed by Shen et al. (2021)
 973 for all baselines and models.

974 In the hyper-parameters, we follow the setting of Chen et al. (2023b) in the ERC task. Moreover, in
 975 the ECPE, the learning rate is set to 3e-5, batch size is set to 32, and epoch is set to 60. Further in our
 976 approach, hidden size of GNN is set to 300, and dropout rate is 0.3. All experiments were conducted
 977 on a NVIDIA GEFORCE RTX 3090 for both training and testing.

978 E.4.3 THE MODEL ON TAS TASK

979 We employed the features extracted from I3D (Carreira & Zisserman, 2017) as the input for our
 980 model. To avoid random bias, we applied our augmented approach across different backbones while
 981 retaining the seed setup from their original studies, ensuring that the specific training epochs are
 982 consistent with the backbones. All experiments were conducted on a NVIDIA GEFORCE RTX 3090
 983 for both training and testing. We set the learning rate to 0.0005, established a weight decay of 0.001,
 984 and utilized Adam as the optimizer. To enhance the training efficiency and avert degenerate matrices
 985 during whitening, we set the batch size of the frame images for a video segment to 512. We followed
 986 the recommendation in Ermolov et al. (2021) to further subdivide the batches during the whitening
 987 process, setting the sub-batch size to 128.

988 F RESULTS ON LLMs

989 In our LLMs experiments, we devised a simple task:
 990 predicting the existence of a causal relationship between
 991 any two utterances (yes or no) to gauge the
 992 LLM’s capability for causal reasoning. Table 6 illustrates
 993 some interesting conclusions - existing prompting
 994 methods is difficult to yield effective outcomes
 995 for this task. For instance, the “step by step” thinking
 996 guided by the CoT approach tends to make LLMs
 997 involve many correlation-based responses. The cluster
 998 approach of Auto-CoT also contributes meaning-
 999 lessly when the samples are too similar. Conversely,
 1000 our iterative prompts instruction enable LLMs to
 1001 uncover causal inconsistencies in its previous answers,
 1002 thereby allowing for self-correction. This self-
 1003 supervised idea appears to impose the LLMs with a capability of “reasoning”.

Table 6: The F1 score of causal relationship recognition of prompt Models. Ours¹_{LLM} represents that calculating Sim^r via LLM while Ours¹_{LLM} represents it via RoBERTa pre-trained model.

| Model | Causalogue | RECCON |
|----------------------------------|------------|--------|
| Zero-Shot | 0.61 | 0.52 |
| Zero-Shot-CoT | 0.58 | 0.51 |
| Auto-CoT | 0.62 | 0.51 |
| Ours ¹ _{LLM} | 0.74 | 0.66 |
| Ours ² _{LLM} | 0.72 | 0.69 |

1004 However, the accuracy of our proposed iterative prompt is substantially dependent on the precision of
 1005 Sim^r . An incorrect Sim^r can lead the correct results to be modified. Furthermore, different from the
 1006 SSMs that could make Sim^r and Sim^s be trained together, the LLM does not provide a learnable
 1007 module to refine Sim^r , only leading to the Sim^s close to the fixed Sim^r . Figure 7 (a, b, c) represents
 1008 the Sim^r calculations obtained through three different measures, where ‘c’ derives from labels,
 1009 providing 100% accurate. The red dashed line in the lineplot denotes the F1 of Sim^r . Our method

³https://drive.google.com/file/d/1R5K_2P1Z3p3RFQ1Ycgmo3TgxvYBzptQG/view?usp=sharing

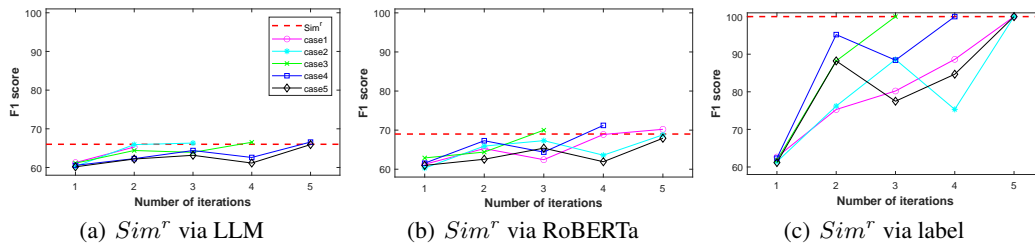


Figure 7: The upper bound of Sim^r and F1 scores of 5 cases in different Sim^r calculating measures

1010 was applied across 5 cases, using each of 3 measures respectively. The results consistently converge
 1011 around the Sim^r which indicates that Sim^r is the upper bound of identifying causal relationships.
 1012 Finally, in order to demonstrate the specific question and answer process of our iterative prompt, we
 1013 provide a complete case as follows.

***user**:* You are assuming the role of a researcher capable of distinguishing between causality and correlation, charged with the task of assessing the causal relationships among individual utterances within a given dialogue. We prescribe that the judgment of causality between two sentences is based on whether the former is the intended target of the latter’s response. Whereas, correlation is gauged on whether the two share similar topics or vocabulary. The following is an example:

Example:

Dialogue:

Hazel drank too much champagne at the party. Oh my goodness! That sounds like quite an eventful party. Well, drinking too much alcohol can have many negative effects on the body. Oh no, I can imagine Hazel waking up with a massive headache tomorrow.

Question 1: Is there a causal relationship from utterance 1 to utterance 2? Answer 1: Yes.

Question 2: Is there a causal relationship from utterance 1 to utterance 3? Answer 2: Yes.

Question 3: Is there a causal relationship from utterance 1 to utterance 4? Answer 3: Yes.

Question 4: Is there a causal relationship from utterance 2 to utterance 3? Answer 4: No.

Question 5: Is there a causal relationship from utterance 2 to utterance 4? Answer 5: No.

Question 6: Is there a causal relationship from utterance 3 to utterance 4? Answer 6: Yes.

Given the above example, with its associated questions and answers, consider the following dialogue:

Dialogue:

Charlotte has no idea how to avoid massive estate taxes. Estate taxes are a topic of concern for many people in various countries. So, does anyone else have any knowledge or ideas on how to reduce estate taxes? Oh, that reminds me of a story about my uncle.

Question 1: Is there a causal relationship from utterance 1 to utterance 2?

Question 2: Is there a causal relationship from utterance 1 to utterance 3?

Question 3: Is there a causal relationship from utterance 1 to utterance 4?

Question 4: Is there a causal relationship from utterance 2 to utterance 3?

Question 5: Is there a causal relationship from utterance 2 to utterance 4?

Question 6: Is there a causal relationship from utterance 3 to utterance 4?

****chatgpt****: Answer 1: Yes. Utterance 2 is a response to the topic brought up in utterance 1, specifically the issue of estate taxes.

Answer 2: Yes. Utterance 3 directly follows up on the issue presented in utterance 1, asking if others have knowledge about the same topic.

Answer 3: No. While utterance 4 may be about a related issue, there is no explicit causal relationship based on the given information.

Answer 4: Yes. Utterance 3 is a continuation of the discussion about estate taxes, which was the topic of utterance 2.

Answer 5: No. Utterance 4 does not indicate a response or direct continuation of the topic presented in utterance 2.

Answer 6: Yes. The story mentioned in utterance 4 could be a response to the question asked in utterance 3.

****user****: After verification, the first utterance should have a causal relationship with the fourth utterance, and there is no common cause between the third utterance and the fourth utterance, and the third utterance should not have a causal relationship with the fourth utterance, and there is no common cause between the second utterance and the third utterance, and the second utterance should not have a causal relationship with the third utterance. Please re-answer based on these circumstances.

****chatgpt****: With these new circumstances:

Answer 1: Yes. Utterance 2 is a response to the topic brought up in utterance 1, specifically the issue of estate taxes.

Answer 2: No. Even though utterance 3 follows utterance 1, we are told that there's no causal relationship between them.

Answer 3: Yes. Based on your new information, utterance 4 now has a causal relationship with utterance 1.

Answer 4: No. Taking into account the revised situation, there's no causal link between utterance 2 and utterance 3.

Answer 5: No. Utterance 4 does not indicate a response or direct continuation of the topic presented in utterance 2.

Answer 6: No. According to the new context, there's no causal relationship between utterance 3 and utterance 4.

****user****: After verification, the first utterance should have a causal relationship with the fourth utterance. Please re-answer based on these circumstances.

1014

****chatgpt****: Given the new specification:

Answer 1: Yes. Utterance 2 is a response to the topic brought up in utterance 1, specifically the issue of estate taxes.

Answer 2: Yes. Utterance 3 directly follows up on the issue presented in utterance 1, asking if others have knowledge about the same topic.

Answer 3: Yes. Based on your new information, utterance 4 now has a causal relationship with utterance 1 (Charlotte’s lack of knowledge about estate taxes somehow triggers the story about the speaker’s uncle).

Answer 4: Yes. Utterance 3 is a continuation of the discussion about estate taxes, which was the topic of utterance 2.

Answer 5: No. Utterance 4 does not indicate a response or direct continuation of the topic presented in utterance 2.

Answer 6: Yes. Under normal circumstances, the story mentioned in utterance 4 can be perceived as a response to the question asked in utterance 3. However, if we strictly adhere to the new context provided(“no common cause between the third utterance and the fourth utterance”), then the answer would be: No. There is no causal relationship between utterance 3 and utterance 4.

1015

***user**:* After verification, there is no common cause between the second utterance and the third utterance, and the second utterance should not have a causal relationship with the third utterance, and there is no common cause between the third utterance and the fourth utterance, and the third utterance should not have a causal relationship with the fourth utterance. Please re-answer based on these circumstances.

***chatgpt**:* Given the additional constraints:

Answer 1: Yes. Utterance 2 is a response to the topic brought up in utterance 1, specifically the issue of estate taxes.

Answer 2: Yes. Despite the absence of causal links to utterances 2 and 4, utterance 3 directly follows up on the issue presented in utterance 1, asking if others have knowledge about the same topic.

Answer 3: Yes. Based on your new information, utterance 4 now has a causal relationship with utterance 1 (Charlotte’s lack of knowledge about estate taxes reminded the speaker of a story about their uncle).

Answer 4: No. Based on your new information, there is no causal relationship between utterance 2 and utterance 3.

Answer 5: No. Utterance 4 does not indicate a response or direct continuation of the topic presented in utterance 2.

Answer 6: No. Based on your updated context, there is no causal relationship between utterance 3 and utterance 4.

1016 G RESULTS ON ECPE, ERC, AND TAS TASKS

1017 We assessed the performance of Ours_{SSM} on three downstream tasks: ECPE, ERC, and TAS. These
 1018 tasks not only represent typical instances of indefinite data (text and video), but they also embody
 1019 causally related tasks. For instance, in ECPE and ERC tasks, mastering the causal relationships
 1020 between utterances is vital, while in the TAS task, recognizing and effectuating transformation from
 1021 intra-frame relations to intra-segment relations is crucial.

1022 From the outcomes presented in Tables 7 and 8, Ours_{SSM} exhibits a remarkable improvement when
 1023 dealing with these high-level causal models. The underlying reason for this enhancement is that,
 1024 under the conditions of ensured causal consistency, an increase in the accuracy of the causal model
 1025 promotes enhancements in both the causal structure (C-Dis in Table 8) and causal representation
 1026 (Edit in Table 8), surpassing other methods, hence improving the final results.

Table 7: Results of SSMs on ECPE and ERC tasks. 95% confidence interval shown. All evaluation metrics used in the Table were F1 scores (Appendix E.2). The backbone of Ours_{SSM} are biCD and DualGAT, respectively.

| ECPE | | | ERC | | | | |
|---------------------|-------------------------|--|---------------------|-------------------------|-------------------------|-------------------------|-------------------------|
| Model | RECCON | | Model | MELD | EmoryNLP | DD | IEM |
| ACCD | 73.17 \pm 1.1 | | ACCD | 63.81 \pm 0.11 | 39.54 \pm 0.12 | 59.53 \pm 0.01 | 69.17 \pm 0.15 |
| biCD | 74.14 \pm 0.74 | | biCD | 63.22 \pm 0.17 | 38.21 \pm 0.11 | 59.64 \pm 0.07 | 67.15 \pm 0.09 |
| EDKA-GM | 72.14 \pm 0.93 | | DAG-ERC | 63.65 \pm 0.05 | 39.02 \pm 0.13 | 59.33 \pm 0.01 | 68.03 \pm 0.15 |
| seF | 74.55 \pm 0.98 | | DualGAT | 66.72 \pm 0.12 | 40.88 \pm 0.15 | 61.80 \pm 0.02 | 67.74 \pm 0.21 |
| - | - | | MultiEMO | 61.23 \pm 1.26 | 37.14 \pm 0.11 | 57.46 \pm 0.01 | 64.41 \pm 0.16 |
| Ours _{SSM} | 76.89 \pm 1.21 | | Ours _{SSM} | 67.79 \pm 0.18 | 40.95 \pm 0.08 | 62.57 \pm 0.01 | 69.81 \pm 0.26 |

Table 8: Results of SSMs on TAS task. All evaluation metrics used in the Table were introduced in Appendix E.2.4. The backbone of Ours_{SSM} is CETnet.

| Model | GTEA | | | | | | 50salads | | | | Breakfast | | | | | | | |
|---------------------|-----------------|-------------|-------------|-------------|-------------|------------|-----------------|-------------|-------------|-------------|-------------|------------|-----------------|-------------|-------------|-------------|-------------|-------------|
| | F1@{10, 25, 50} | | Edit | Acc | C-Dis | | F1@{10, 25, 50} | | Edit | Acc | C-Dis | | F1@{10, 25, 50} | Edit | Acc | C-Dis | | |
| MSTCN++ | 82.3 | 83.6 | 71.9 | 79.8 | 77.6 | 8.4 | 79.4 | 77.3 | 69.3 | 71.6 | 82.8 | 3.3 | - | - | - | - | | |
| ASRF | 85.5 | 83.8 | 73.6 | 76.9 | 74.7 | 9.0 | 80.3 | 77.4 | 67.4 | 74.2 | 77.6 | 4.9 | 69.1 | 63.4 | 50.8 | 66.6 | 63.0 | 55.8 |
| CETnet | 90.5 | 89.6 | 78.9 | 85.7 | 79.4 | 7.1 | 87.6 | 87.3 | 80.9 | 82.8 | 87.3 | 2.6 | 72.5 | 68.7 | 57 | 72.8 | 74.2 | 38.1 |
| C2F | 88 | 86.6 | 78.3 | 81.6 | 80.6 | 7.4 | 83.5 | 81.5 | 71.8 | 75.7 | 86.9 | 2.8 | 71.6 | 68.0 | 57.1 | 68.1 | 74.6 | 49.8 |
| Ours _{SSM} | 91.4 | 90.2 | 80.5 | 87.2 | 79.7 | 6.9 | 88.9 | 87.6 | 81.4 | 83.1 | 88.9 | 2.5 | 78.7 | 74.9 | 63.4 | 78.3 | 75.6 | 35.4 |

1027 To better illustrate the role of the causal model in these downstream tasks, we demonstrate two
 1028 visualizations in Figures 8 and 9. Figure 8 displays a visualization of the adjacency matrix for the
 1029 ECPE task, which can be equated with a causal graph, showing how the model assigns weights to
 1030 the context when learning utterance relationships. Figure 8 demonstrates that the superiority of the
 1031 causal method over non-causal ones lies in turning the adjacency matrix into a DAG, thus avoiding
 1032 the factual error of treating earlier utterances as outcomes of latter ones. However, due to unknown
 1033 causal labels, there is not a sufficiently strong constraint for causal graph, which often leaves the
 1034 model uncertain about which edges in the DAG should exist. Our model mitigates this issue by using
 1035 causal consistency constraints, enabling the model to identify the correct edges through contrastive
 1036 learning of the causal representation and structure.

1037 Figure 9 shows the TAS task’s visualization results, illustrating that causal consistency between
 1038 frames and segments can significantly reduce the exsitings of trivial segments. This harmonizes
 1039 with our intuition on indefinite data: all frames within a segment share a similar causal relationship.
 1040 Such a causal model goes beyond the scope of this paper, and we thoroughly discuss these extended
 1041 contributions in the discussion section.

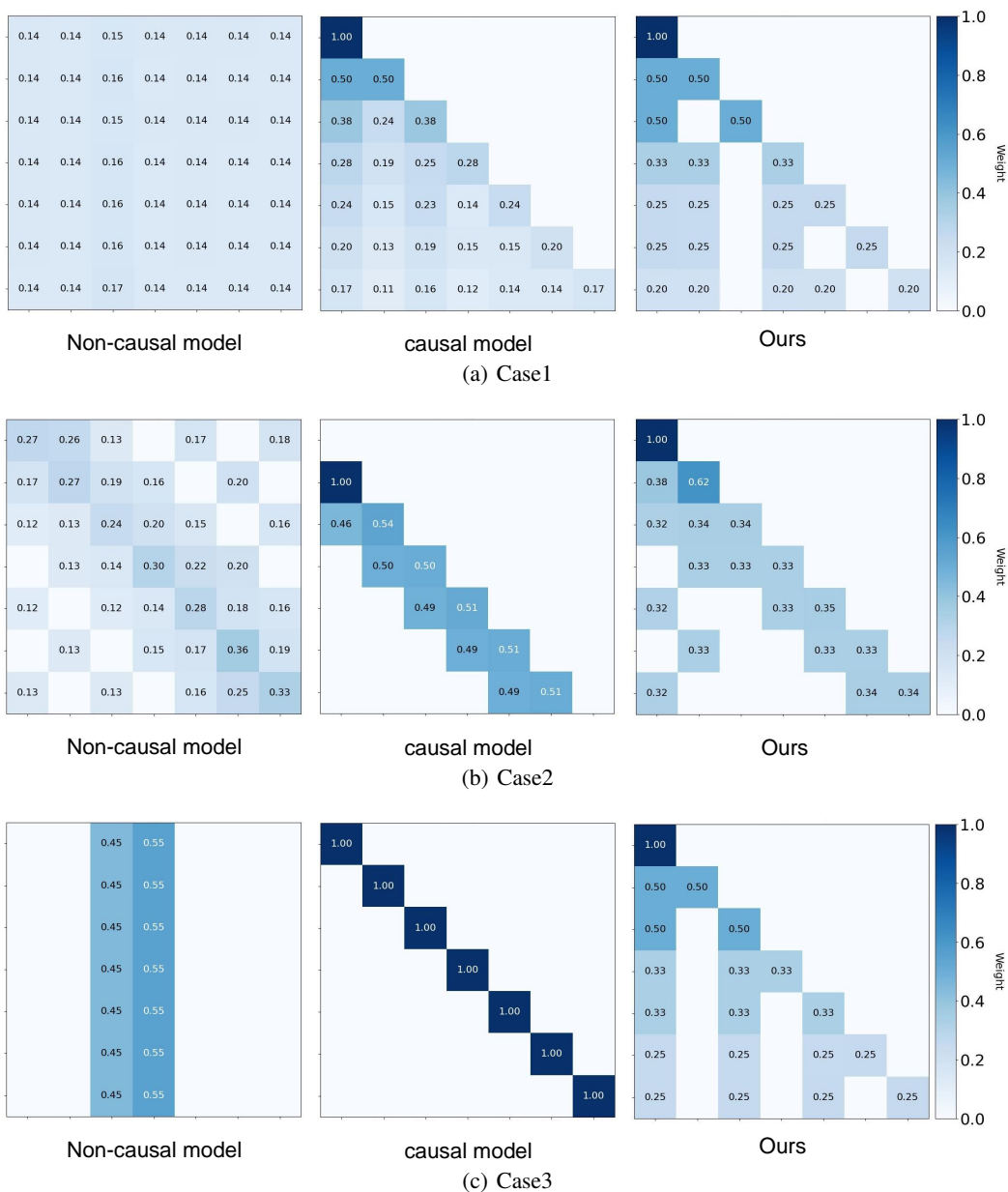
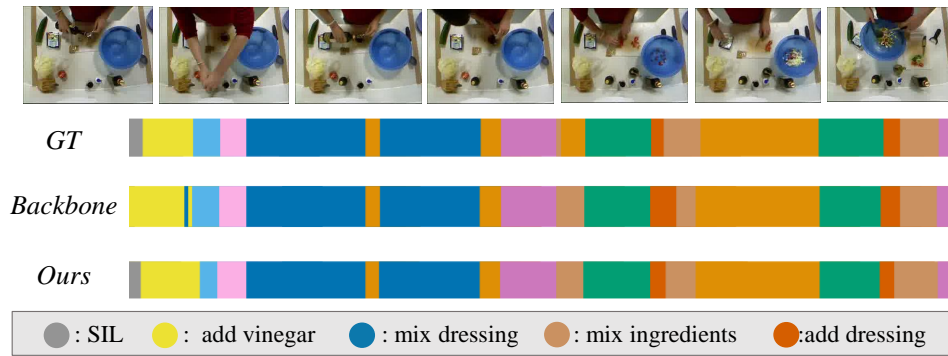
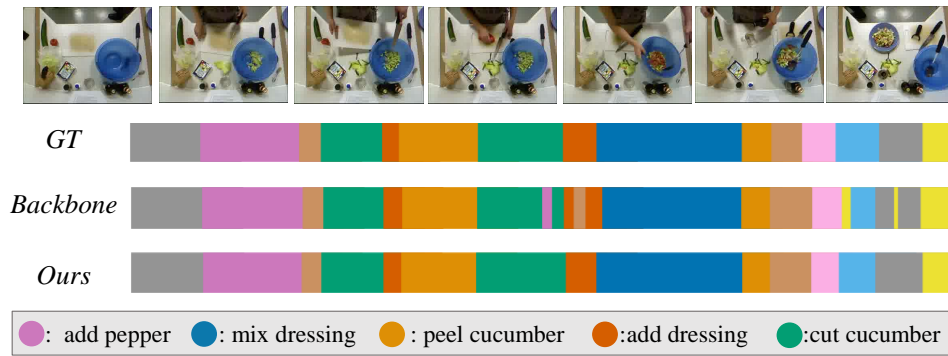


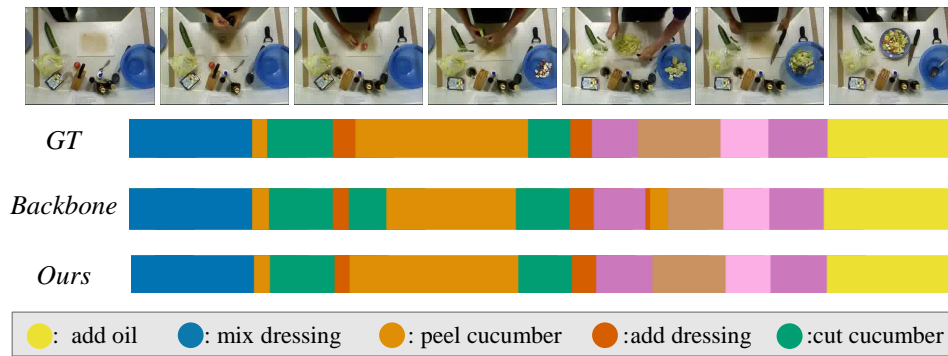
Figure 8: Visualization of adjacent matrices of 3 cases on ECPE task. Non-causal model is EDKA-GM, and we choose ACCD as causal model. The adjacent matrix is $N \times N$ representing the relationship between any two utterances.



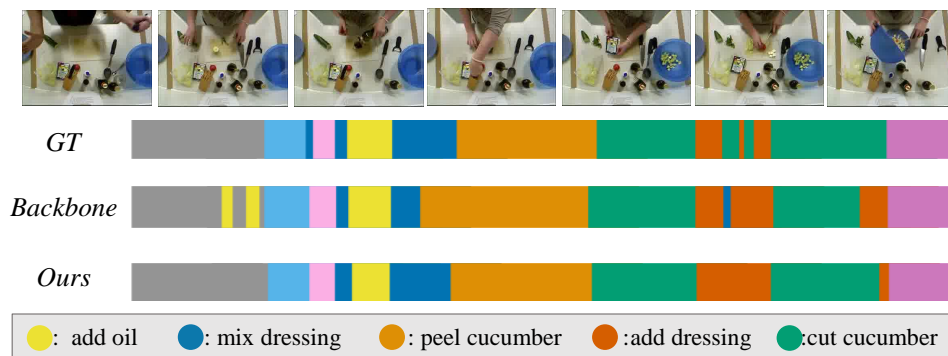
(a) Case1



(b) Case2



(c) Case3



(d) Case4

Figure 9: Visualization of results of 4 cases on 50salads dataset. GT represents the Ground Truth, Backbone we choose is the MSTCN++.