



浙江大學
ZHEJIANG UNIVERSITY

ICML | 2024

Forty-first International
Conference on Machine Learning

Learning Causal Relations from Subsampled Time Series with Two Time-Slices

Anpeng Wu¹, Haoxuan Li², Kun Kuang^{1*}, Keli Zhang³, Fei Wu^{145*}

1 Department of Computer Science and Technology, Zhejiang University, Hangzhou, China

2 Center for Data Science, Peking University, Peking, China

3 Huawei Noah's Ark Lab, Huanwei, Shenzhen, China

4 Shanghai Institute for Advanced Study, Zhejiang University, Shanghai, China

5 Shanghai AI Laboratory, Shanghai, China

 anpwu@zju.edu.cn

*Corresponding authors



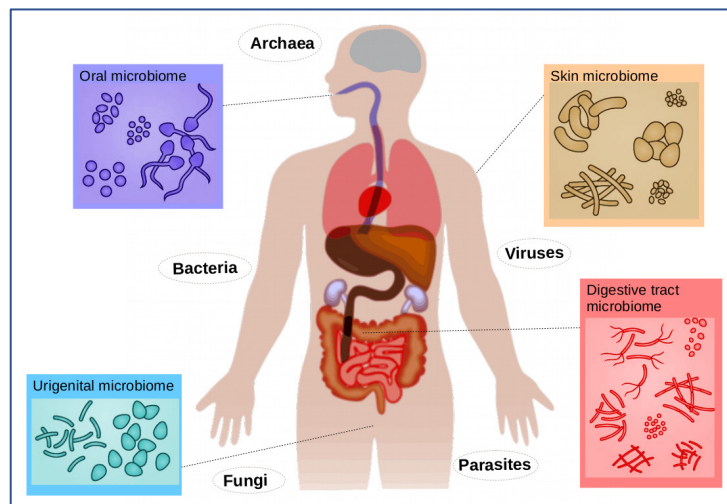
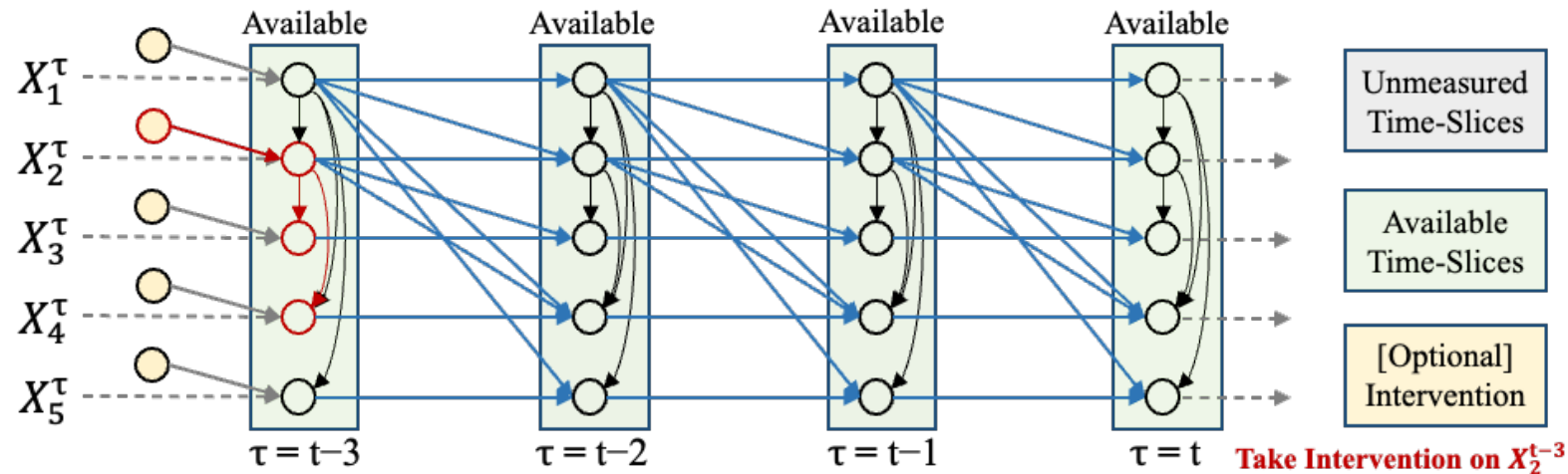


CLIMATE

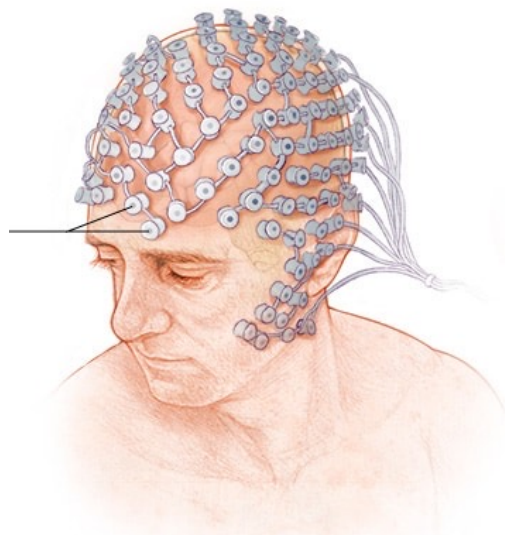
Time-Series Data is ubiquitous in real-applications.



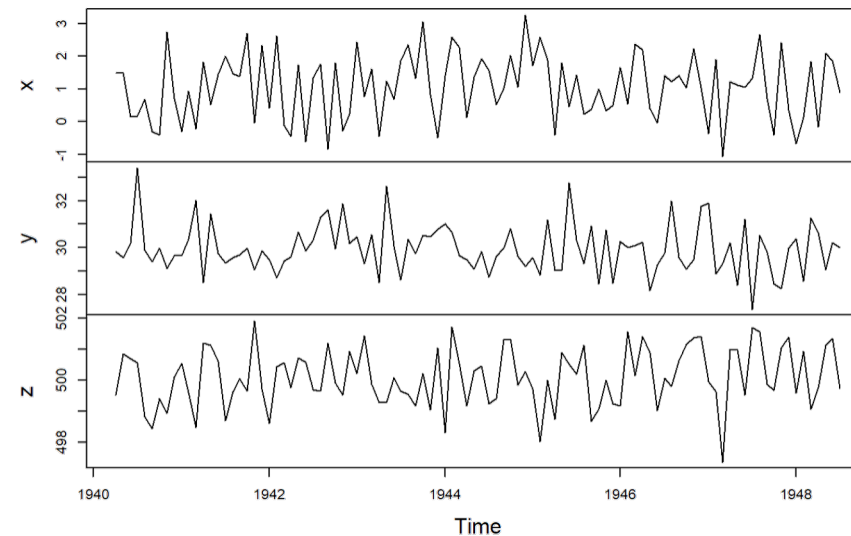
Climate



Human Microbiome



EEG Electrodes



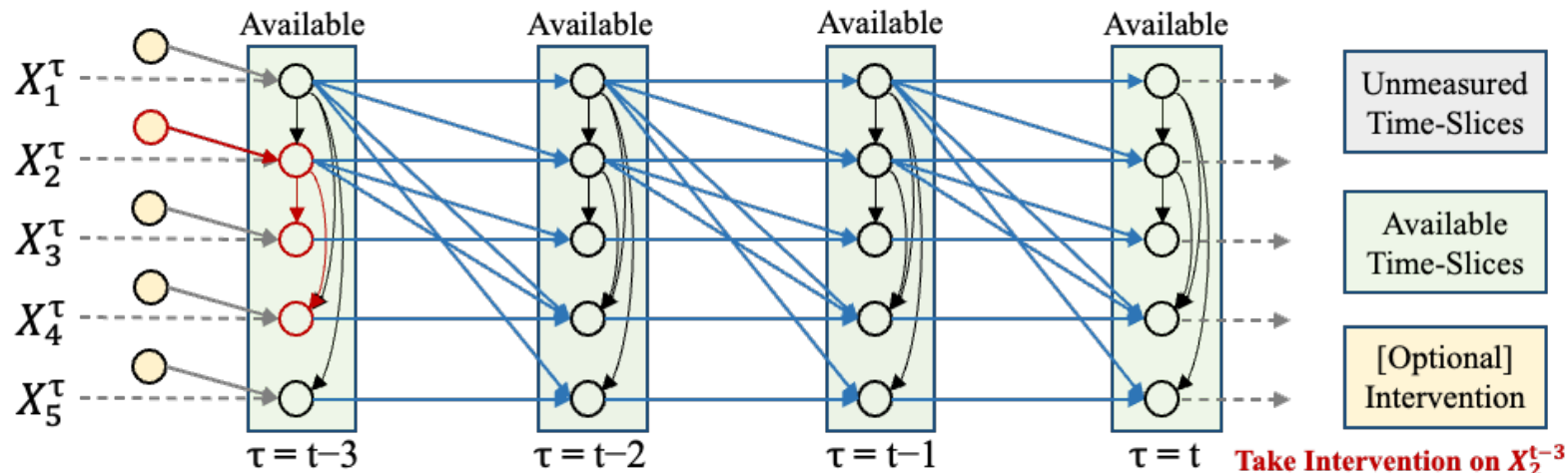
Panel: Multivariable Time Series



Time-Series Data is ubiquitous in real-applications.



Climate



Gold Method: Granger causality, which builds an autoregressive to model causal relations.

Advantages: Granger causality does not require specific assumptions about functional forms or distributions, making it applicable to various types of time series data.

Disadvantages: Granger causality requires

- **Correct Time Lags:** Granger causality relies on lagged terms to model causal relationships.
- **Numerous Time Slices:** Sufficient time slices are needed for accurate inference, and inadequate data length or large time intervals may lead to unreliable results.
- **Without Instantaneous Effects:** Granger causality assumes causal relationships based on lagged terms and cannot capture instantaneous causality, such as simultaneous changes between variables at the same time point.

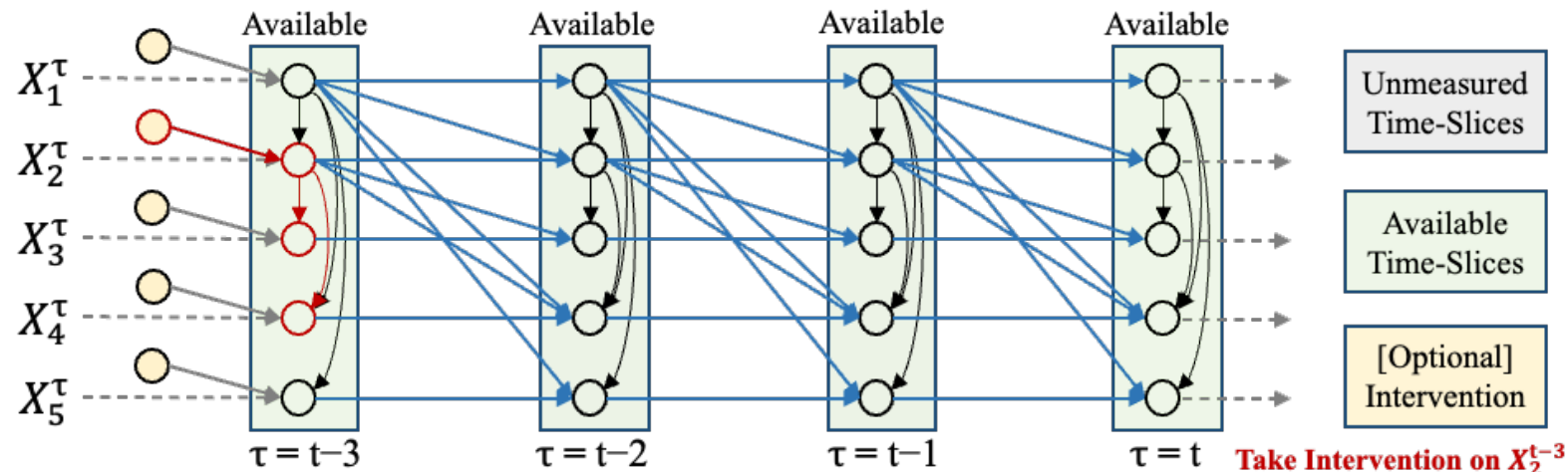


CLIMATE

Time-Series Data is ubiquitous in real-applications.



Climate



Stational Methods.

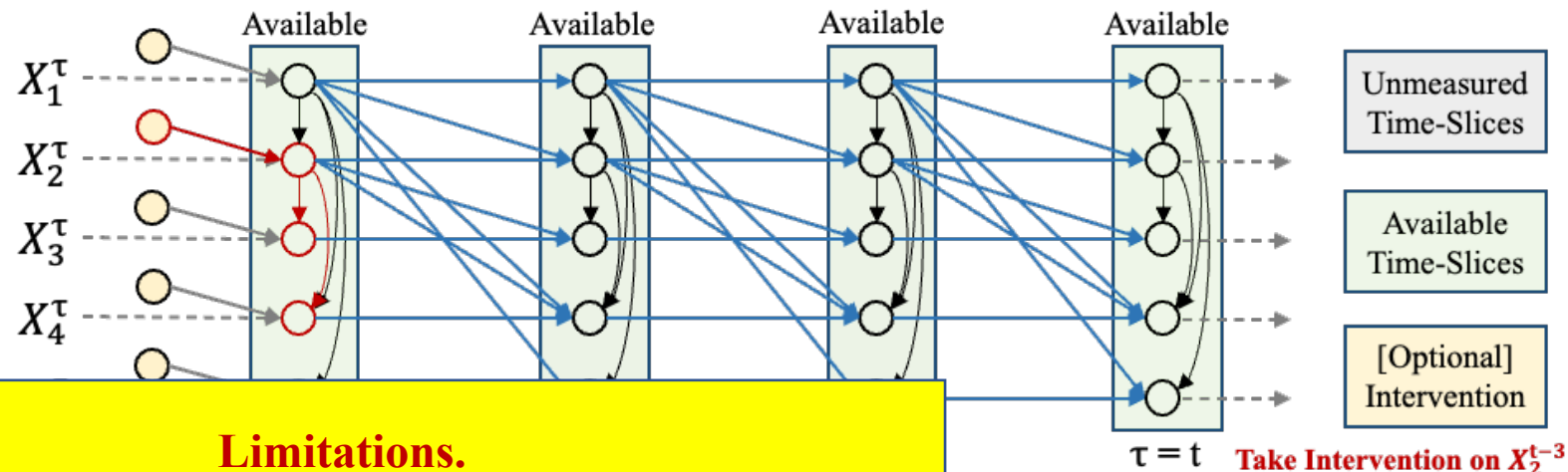
- **Constraint-based methods:**
 - PC, FCI, SGS and ICPs.
- **Score-based methods:**
 - GES and GIES.
- **Continuous-optimization methods:**
 - GraNDAG, GOLEM, NOTEARS, and ReScore.
- **Hybrid methods:**
 - GSP and IGSP.

Temporal Methods.

- **Granger causality and Autoregressive:**
 - PWGC, MVGC, TCDF.
- **Variants of Stational Methods:**
 - PCMCI, oCSE, ANLTSM, tsFCI, SVAE-FCI, VarLiNGAM, DYNOTEARS.
- **Others:**
 - CD-NOD, VarLiNGAM, TiMINo.



Time-Series Data is ubiquitous in real-applications.



Limitations.

- Some works rely on Linear Assumption;
- Some works rely on No Instantaneous Effects;
- All works rely on Modeling Causal Structures at the System Timescale.
- All works rely on Causal Sufficiency Assumption;
- All works rely on Numerous Time Slices Data;

Stational Methods.

- **Constraint-based**
 - PC, FCI, SGS and
- **Score-based methods**
 - GES and GIES.
- **Continuous-optimization methods:**
 - GraNDAG, GOLEM, NOTEARS, and ReScore.
- **Hybrid methods:**
 - GSP and IGSP.

and Autoregressive:

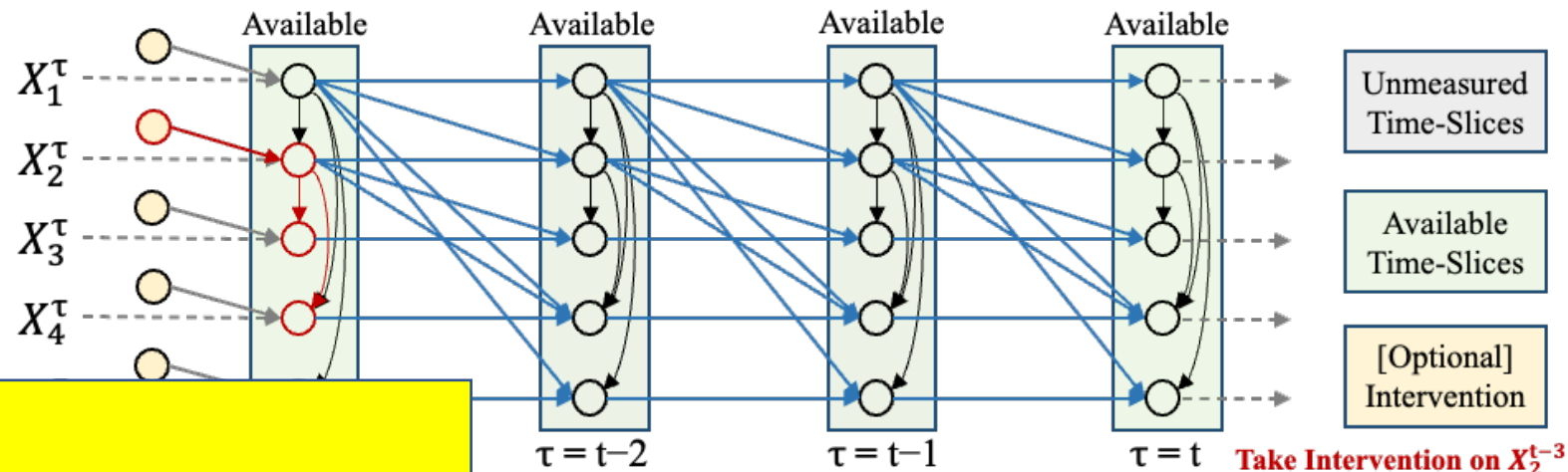
TCDF.

al Methods:

- PCMCI, oCSE, ANLTSM, tsFCI, SVAE-FCI, VarLiNGAM, DYNOTEARS.
- **Others:**
 - CD-NOD, VarLiNGAM, TiMINo.



Subsampled Time-Series Data is ubiquitous in real-applications.



Limitations.

- Some works rely on Linear Assumption;
- Some works rely on No Instantaneous Effects;
- All works rely on Modeling Causal Structures at the System Timescale.
- All works rely on Causal Sufficiency Assumption;
- All works rely on Numerous Time Slices Data;

- **Continuous-optimization methods:**
 - GraNDAG, GOLEM, NOTEARS, and ReScore.
- **Hybrid methods:**
 - GSP and IGSP.

Subsampled Time Series.

Measurements are sparse and sampled at a coarser timescale than the causal timescale of the underlying system

Temporal M
Granger
AC
Variants

- PCMCI, oCSE, ANLTSM, tsFCI, SVAE-FCI, VarLiNGAM, DYNOTEARS.

- **Others:**
 - CD-NOD, VarLiNGAM, TiMINo.

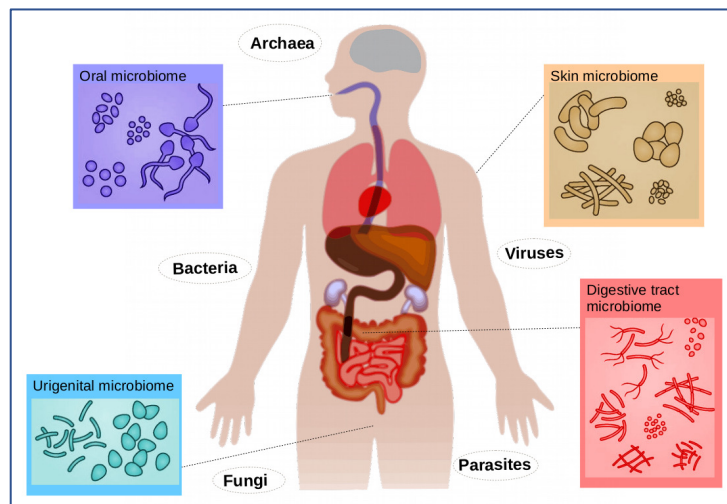
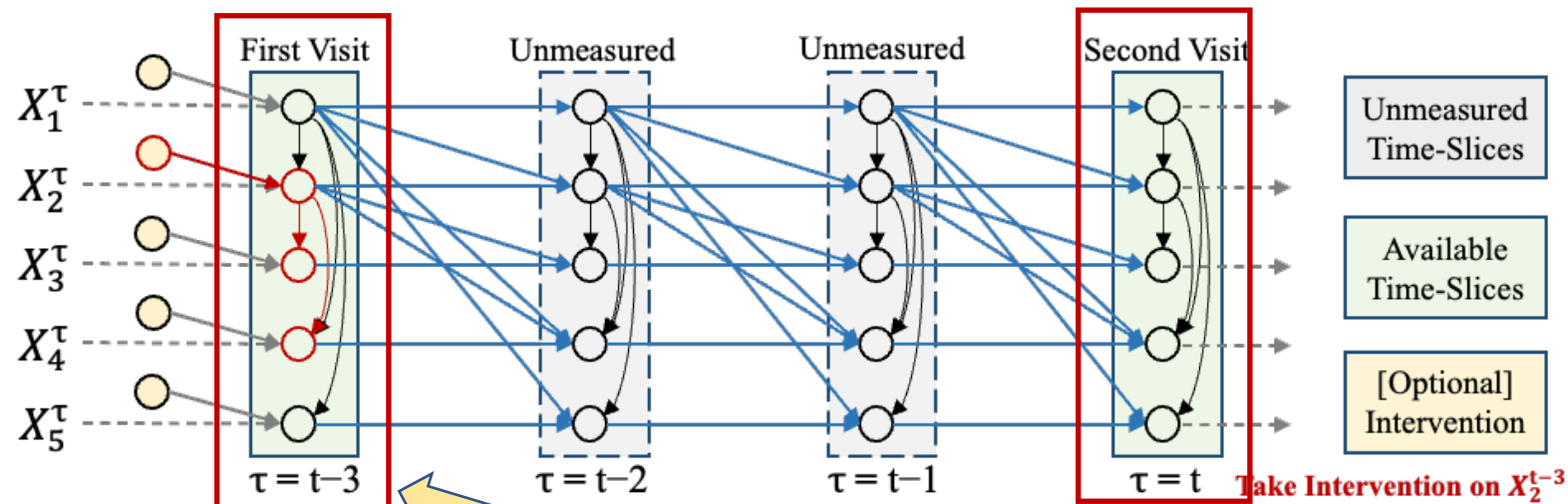


CLIMATE

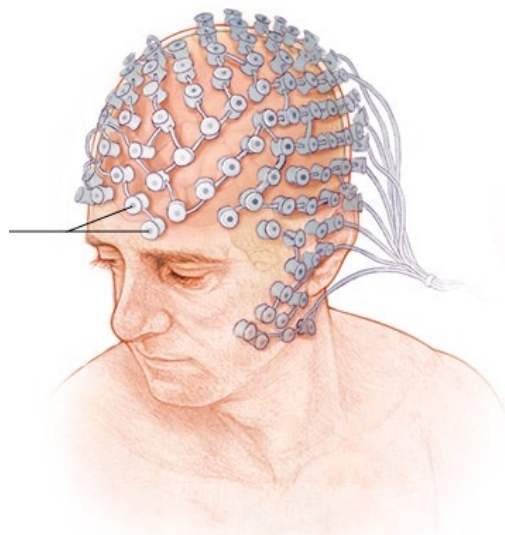
Take Medical Scenarios as an Example.



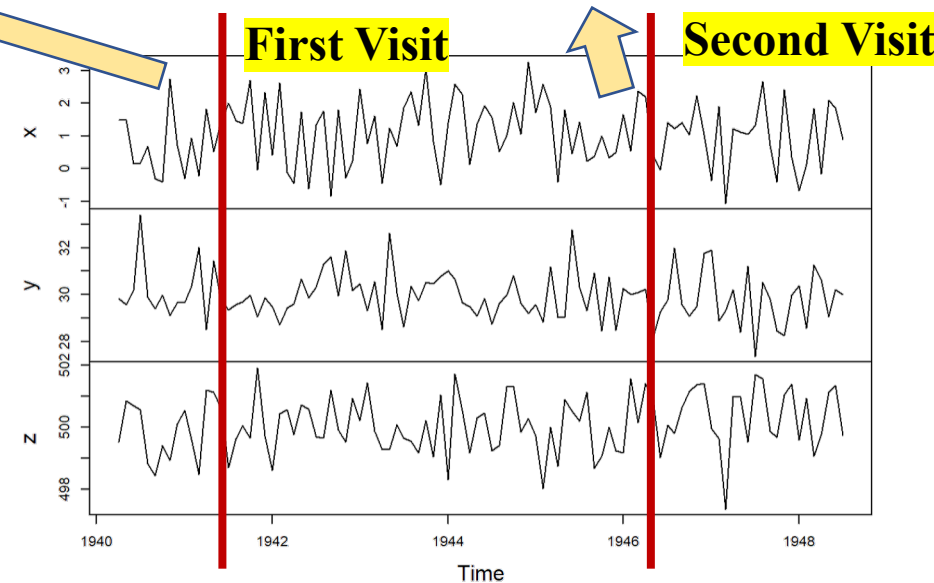
Climate



Human Microbiome



EEG Electrodes



Panel: Multivariable Time Series

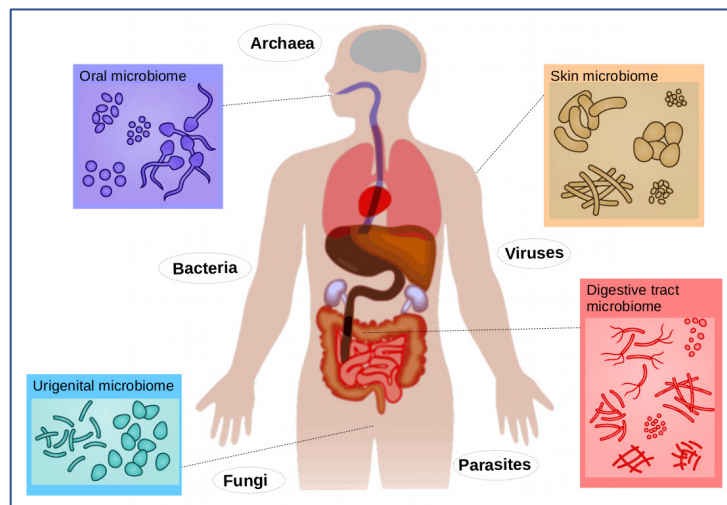
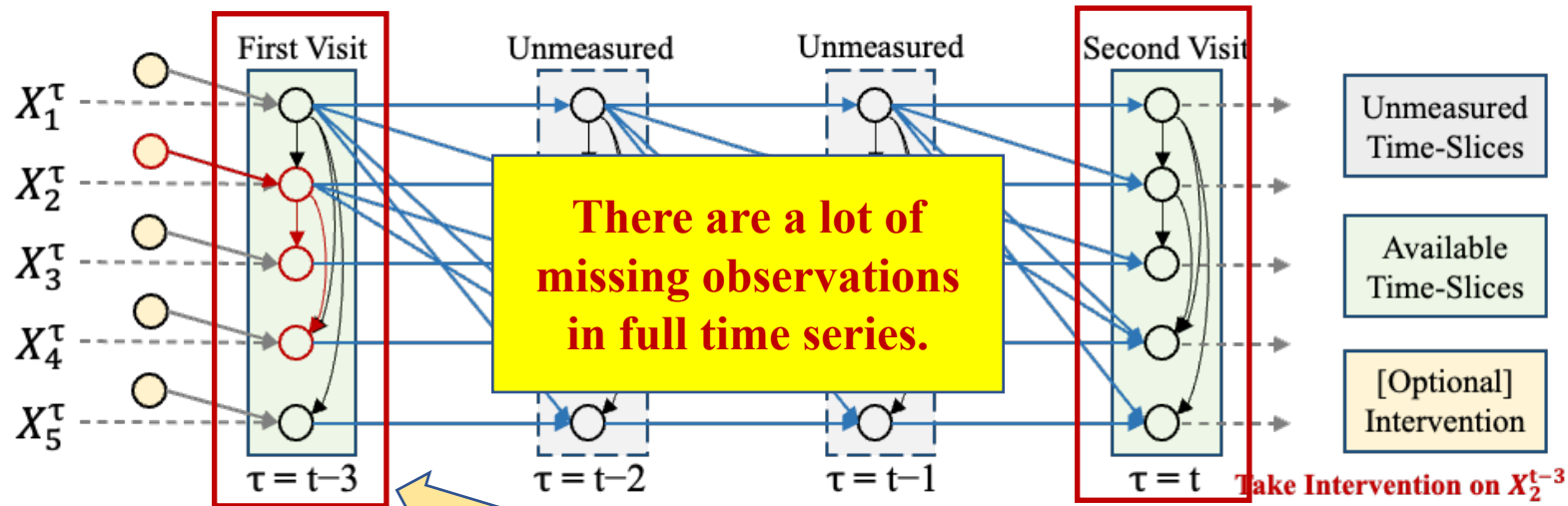


CLIMATE

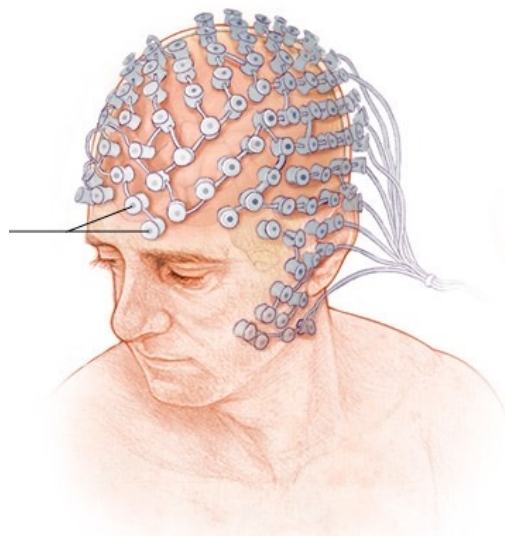
Take Medical Scenarios as an Example.



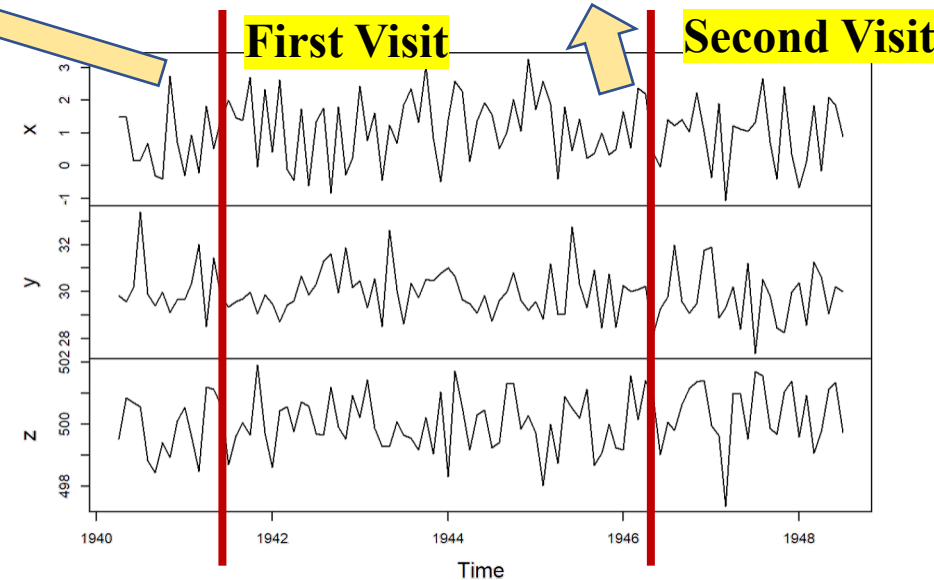
Climate



Human Microbiome



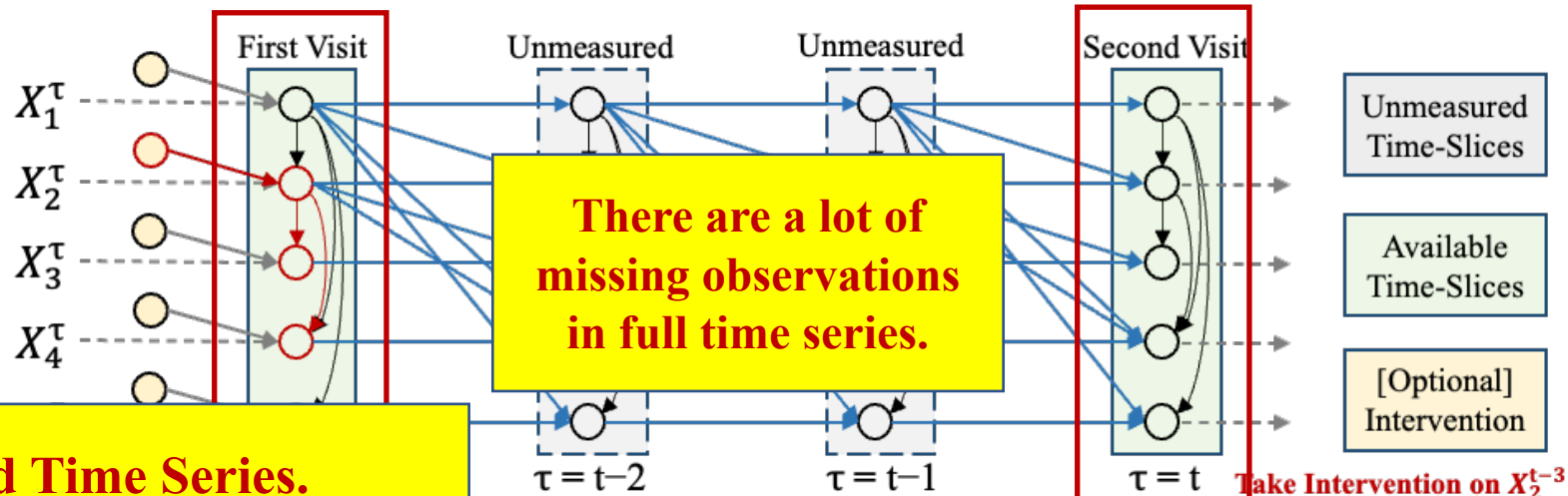
EEG Electrodes



Panel: Multivariable Time Series

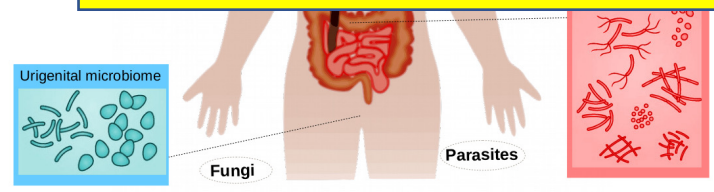


Take Medical Scenarios as an Example.



Challenges in Subsampled Time Series.

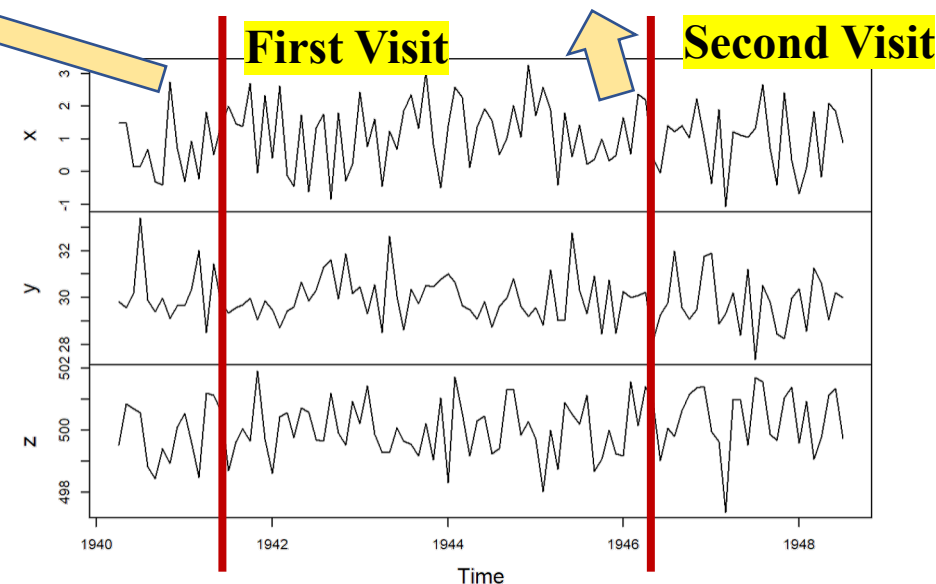
- Some works rely on Linear Assumption;
- Some works rely on No Instantaneous Effects;
- All works rely on Modeling Causal Structures at the System Timescale.
- All works rely on Causal Sufficiency Assumption;
- All works rely on Numerous Time Slices Data;



Human Microbiome



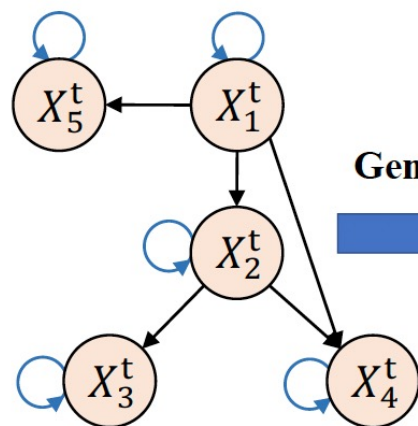
EEG Electrodes



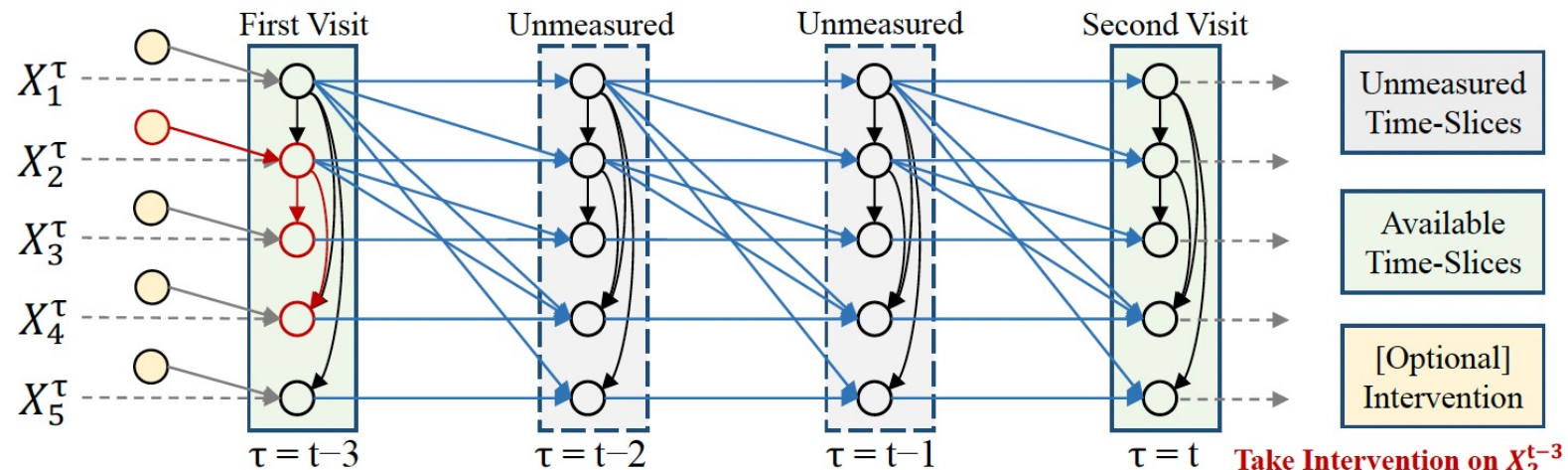
Panel: Multivariable Time Series



Abstract Figure.

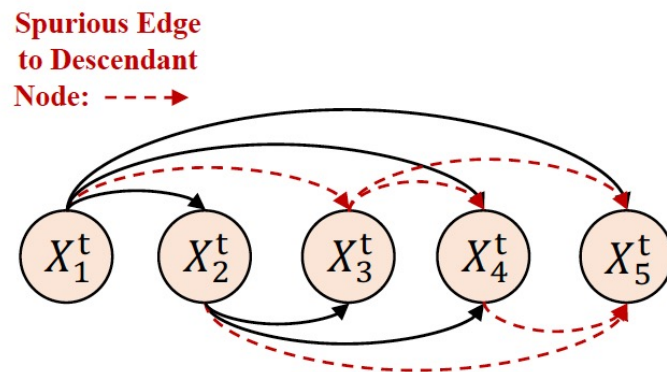
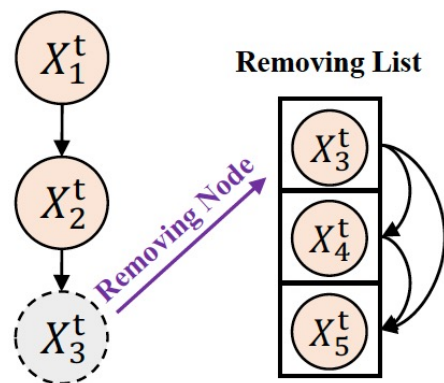


Generate

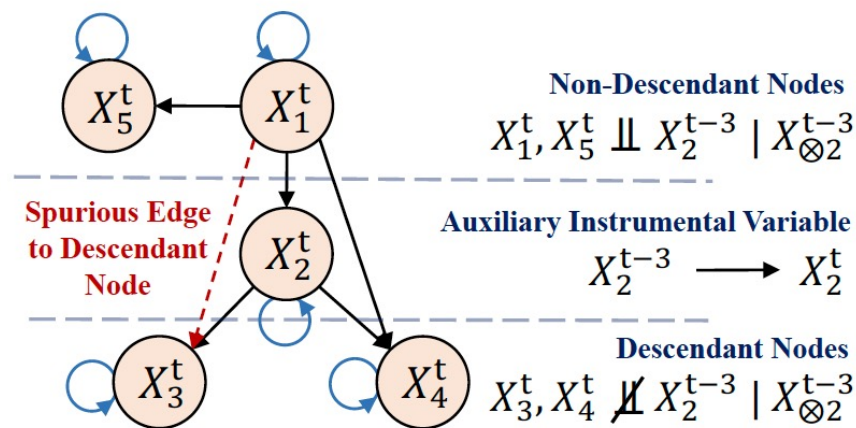


(a) Summary Causal Graph

(b) Graphical Models for Subsampled Time-Series (Interventions or Two Time-Slices)



(I) Sequentially Identifying and Removing the Leaves (II) Complete Topological Ordering
(c) Identifying Leaf Nodes for Complete Topological Ordering (SCORE)



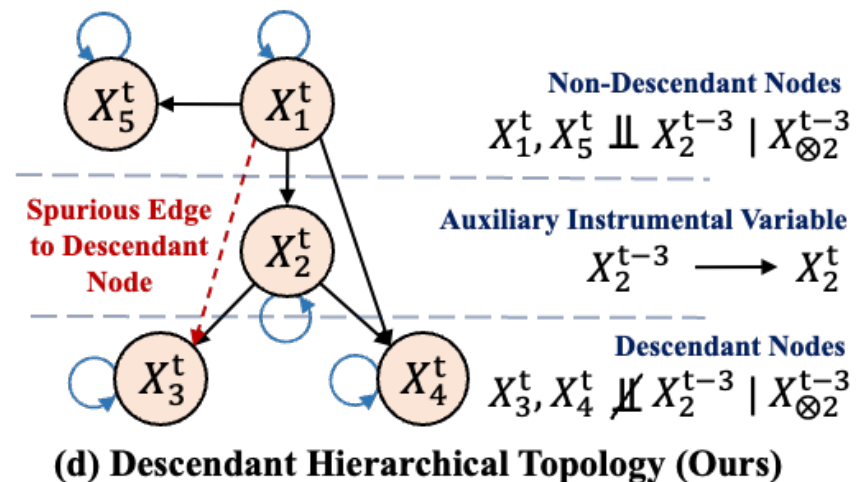
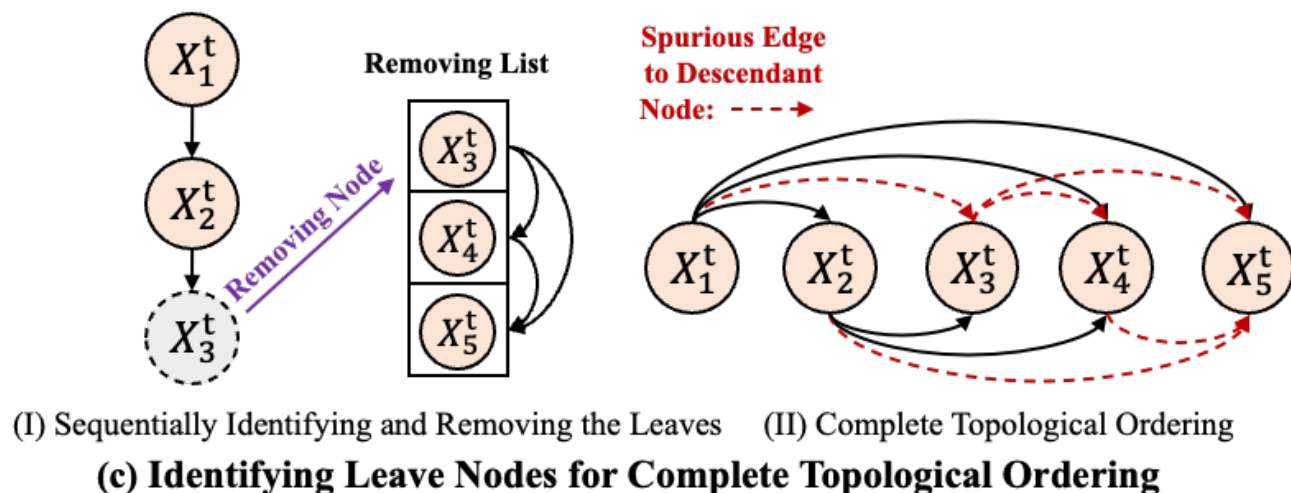
(d) Descendant Hierarchical Topology (Ours)



Descendant Hierarchical Topology

Definition 4.1 (Complete Topological Ordering). The complete topological ordering $(\pi(\mathbf{X}) = (X_{\pi_1}, X_{\pi_2}, \dots, X_{\pi_d}), \pi_i$ is the reordered index of node) is a sorting of all nodes in a DAG such that for any pair of nodes X_{π_i} and X_{π_j} , if there exists a directed edge from X_{π_i} to X_{π_j} , then $i > j$.

Definition 4.2 (Hierarchical Topological Ordering). In the hierarchical topological ordering e.g., $\Pi(\mathbf{X}) = (\{X_{\pi_1}\}_{L_1}, \{X_{\pi_2}, X_{\pi_3}\}_{L_2}, \dots)$, each layer is denoted by L_i and the located layer of X_j are represented as l_j . If there is a directed edge from X_{π_i} to X_{π_j} , then $l_{\pi_i} > l_{\pi_j}$.



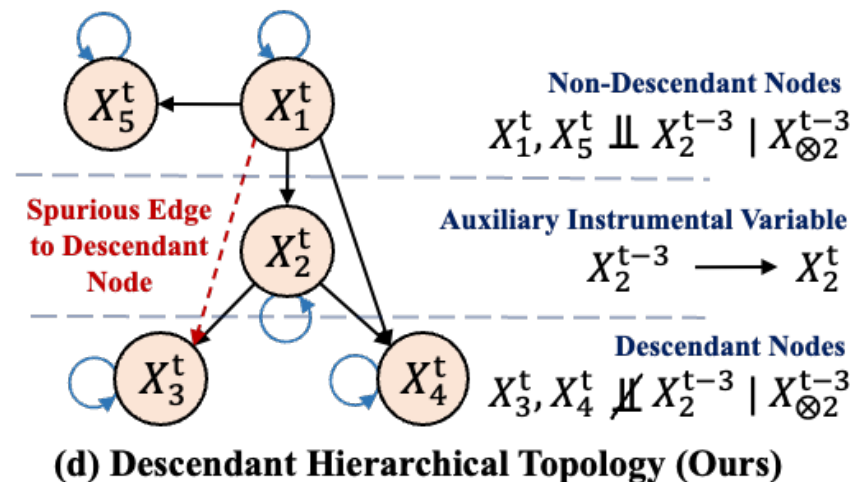
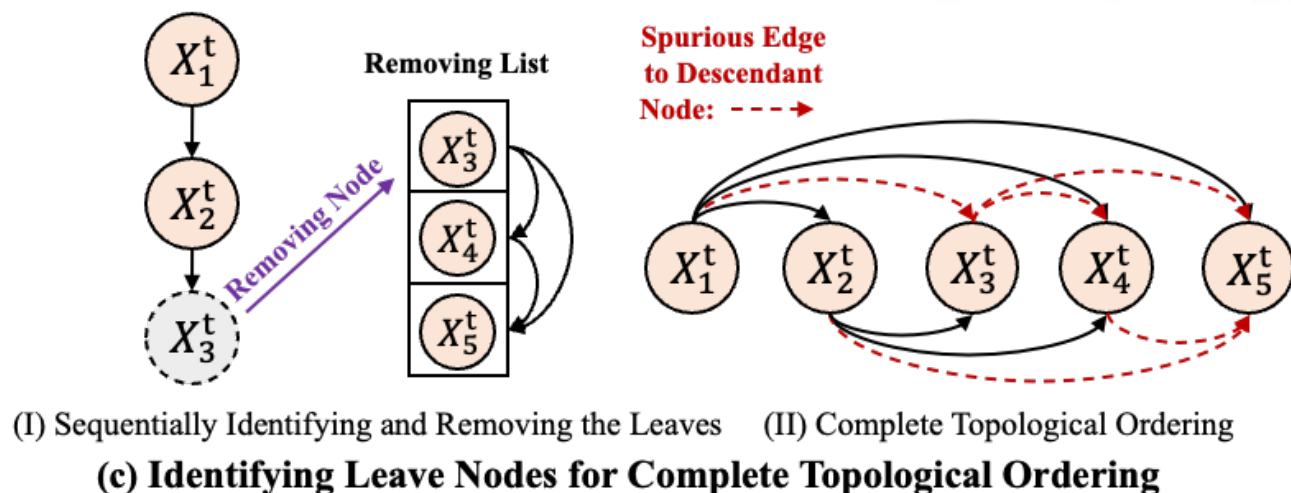


Descendant Hierarchical Topology

Definition 4.1 (Complete Topological Ordering). The complete topological ordering ($\pi(\mathbf{X}) = (X_{\pi_1}, X_{\pi_2}, \dots, X_{\pi_d})$, π_i is the reordered index of node) is a sorting of all nodes in a DAG such that for any pair of nodes X_{π_i} and X_{π_j} , if there exists a directed edge from X_{π_i} to X_{π_j} , then $i > j$.

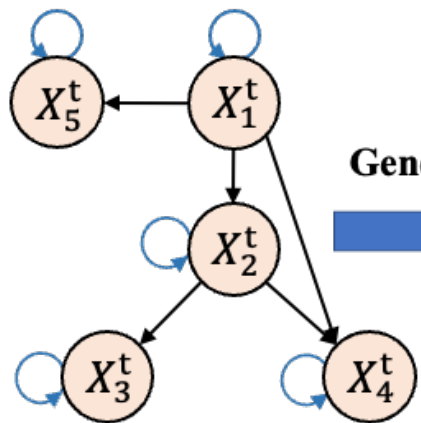
Definition 4.2 (Hierarchical Topological Ordering). In the hierarchical topological ordering e.g., $\Pi(\mathbf{X}) = (\{X_{\pi_1}\}_{L_1}, \{X_{\pi_2}, X_{\pi_3}\}_{L_2}, \dots)$, each layer is denoted by L_i and the located layer of X_j are represented as l_j . If there is a directed edge from X_{π_i} to X_{π_j} , then $l_{\pi_i} > l_{\pi_j}$.

Definition 4.3 (Descendant Hierarchical Topology). In the descendant hierarchical topology, each node X_i^t identifies other nodes as either non-descendant nodes or descendant nodes, and each node X_i^t establishes direct edges pointing to its descendants de_i^t , i.e., $X_i^t \rightarrow \text{de}_i^t, i \in \{1, 2, \dots, d\}$.



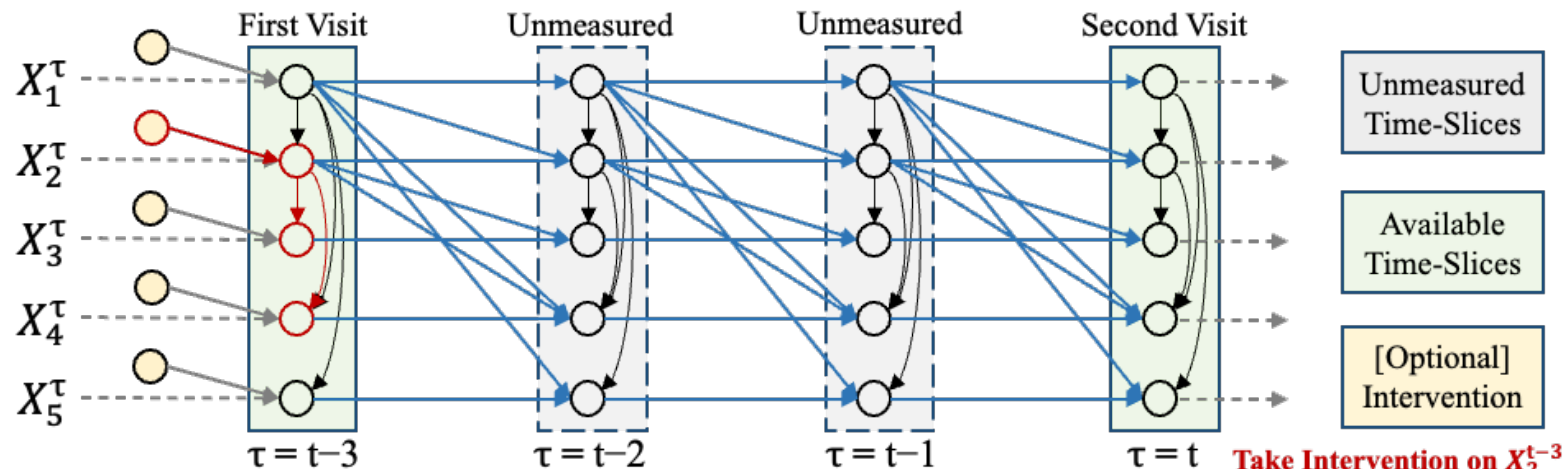


Using Conditional Instrumental Variables to Replace Interventions.



(a) Summary Causal Graph

Generate



(b) Graphical Models for Subsampled Time-Series (Interventions or Two Time-Slices)

Assumption 3.1 (Markov Property). The Markov property of time series assumes the future slice \mathbf{X}^{t+1} depends on current state \mathbf{X}^t but does not depend on history $\mathbf{X}^{1 \dots t-1}$.

Assumption 3.2 (Acyclic Summary Causal Graph, Section 5.2.1 in Assaad et al. (2022)). The summary causal graph of a time series is considered acyclic if the lagged effect of each variable solely affects its own value and its descendants, without any influence on its non-descendants.

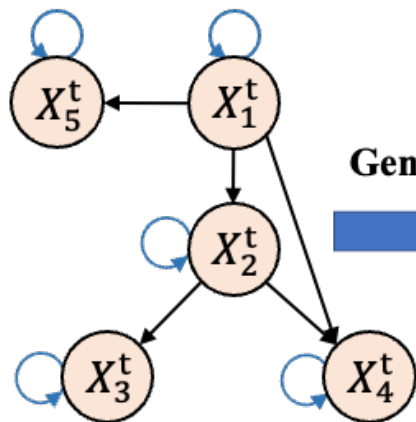
Assumption 3.3 (Consistency Throughout Time, Definition 7 in Assaad et al. (2022)). A causal graph \mathcal{G} for a multivariate time series \mathbf{X} is said to be consistent throughout time if all the causal relationships remain constant throughout time, also referred to as stationary full-time graph.

$$X_i^\tau = f_i(\mathbf{pa}_i^\tau, X_i^{\tau-1}, \mathbf{pa}_i^{\tau-1}) + \epsilon_i^\tau, \quad (1)$$

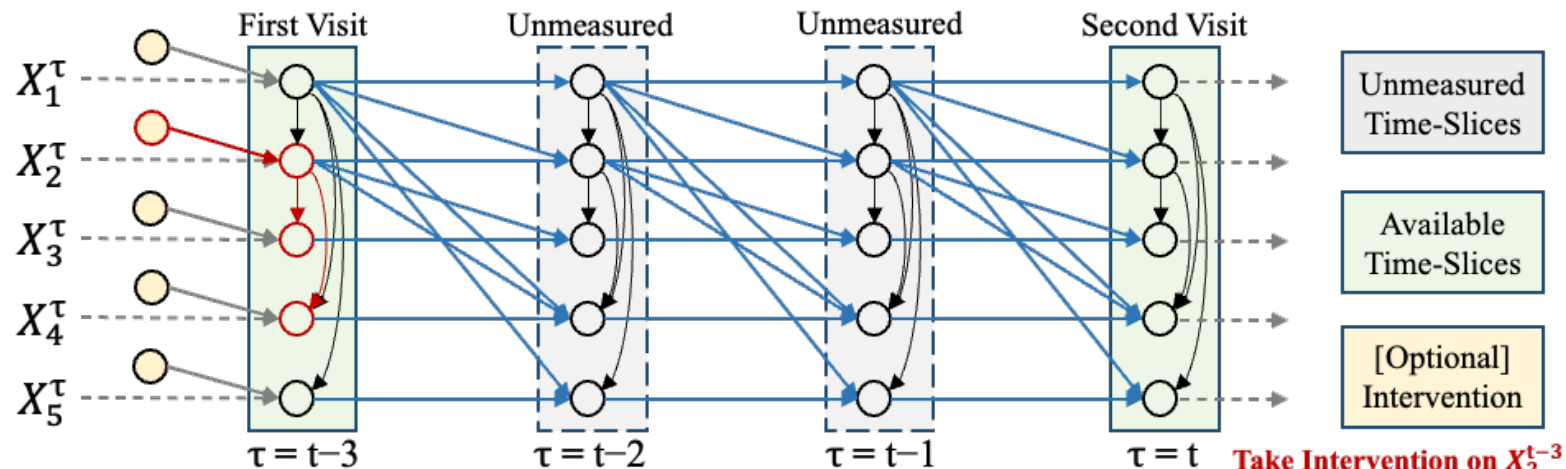
where $f_i(\cdot)$ is a twice continuously differentiable function, which embeds the instantaneous effects from its parents \mathbf{pa}_i^τ at time τ and non-zero time-lagged effects from previous variable $\mathbf{X}^{\tau-1}$; and ϵ_i^τ denotes the *Additive Noise* term at time τ . In the generation function $f_i(\cdot)$, we require that the time-lagged effect of $X_i^{\tau-1}$ on X_i^τ is non-zero.



Using Conditional Instrumental Variables to Replace Interventions.



(a) Summary Causal Graph



(b) Graphical Models for Subsampled Time-Series (Interventions or Two Time-Slices)

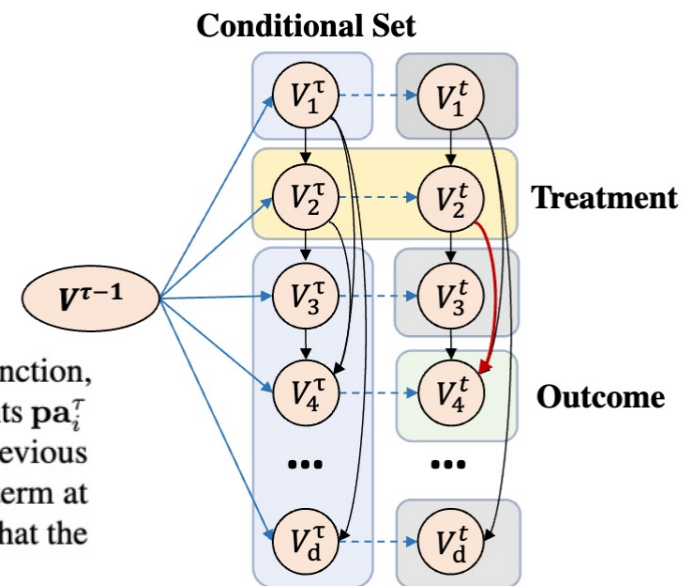
Assumption 3.1 (Markov Property). The Markov property of time series assumes the future slice \mathbf{X}^{t+1} depends on current state \mathbf{X}^t but does not depend on history $\mathbf{X}^{1 \dots t-1}$.

Assumption 3.2 (Acyclic Summary Causal Graph, Section 5.2.1 in Assaad et al. (2022)). The summary causal graph of a time series is considered acyclic if the lagged effect of each variable solely affects its own value and its descendants, without any influence on its non-descendants.

Assumption 3.3 (Consistency Throughout Time, Definition 7 in Assaad et al. (2022)). A causal graph \mathcal{G} for a multivariate time series \mathbf{X} is said to be consistent throughout time if all the causal relationships remain constant throughout time, also referred to as stationary full-time graph.

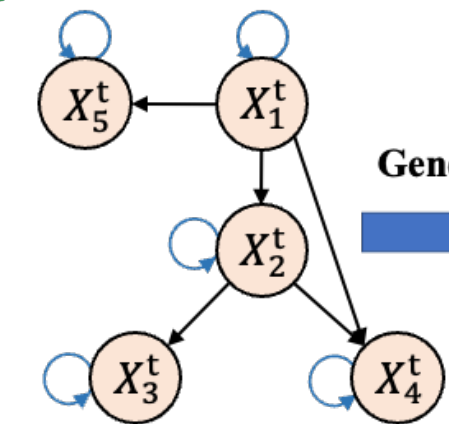
$$X_i^\tau = f_i(\mathbf{pa}_i^\tau, X_i^{\tau-1}, \mathbf{pa}_i^{\tau-1}) + \epsilon_i^\tau,$$

where $f_i(\cdot)$ is a twice continuously differentiable function, which embeds the instantaneous effects from its parents \mathbf{pa}_i^τ at time τ and non-zero time-lagged effects from previous variable $\mathbf{X}^{\tau-1}$; and ϵ_i^τ denotes the *Additive Noise* term at time τ . In the generation function $f_i(\cdot)$, we require that the time-lagged effect of $X_i^{\tau-1}$ on X_i^τ is non-zero.



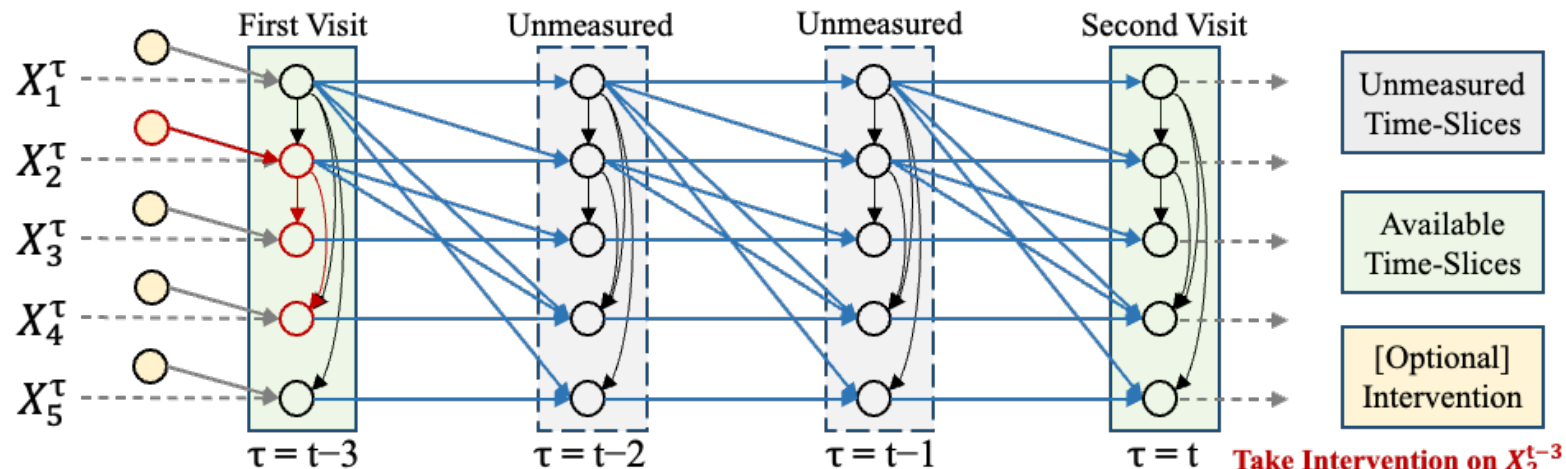


Descendant-Oriented Conditional Independence Criteria



(a) Summary Causal Graph

Generate



(b) Graphical Models for Subsampled Time-Series (Interventions or Two Time-Slices)

Theorem 4.4 (Descendant-Oriented Conditional Independence Criteria). *Given observations $\mathcal{D} = \{X^{t_a}, X^{t_b}\}_{t_a < t_b}$ satisfying Assumptions 3.1, 3.2, and 3.3, for variables $X_i^{t_a}$ and $X_i^{t_b}$, where $i \in \{1, 2, \dots, d\}$, we can conclude that $X_j^{t_b}$ is a descendant node of $X_i^{t_b}$ iff $X_i^{t_a} \not\perp\!\!\!\perp X_j^{t_b} \mid \text{an}_i^{t_a}$.*



Descendant-Oriented Conditional Independence Criteria

Theorem 4.4 (Descendant-Oriented Conditional Independence Criteria). *Given observations $\mathcal{D} = \{\mathbf{X}^{t_a}, \mathbf{X}^{t_b}\}_{t_a < t_b}$ satisfying Assumptions 3.1, 3.2, and 3.3, for variables $X_i^{t_a}$ and $X_j^{t_b}$, where $i \in \{1, 2, \dots, d\}$, we can conclude that $X_j^{t_b}$ is a descendant node of $X_i^{t_a}$ iff $X_i^{t_a} \not\perp\!\!\!\perp X_j^{t_b} \mid \mathbf{an}_i^{t_a}$.*

Corollary 4.5. *Given observations $\mathcal{D} = \{\mathbf{X}^{t_a}, \mathbf{X}^{t_b}\}_{t_a < t_b}$, for variables X_i and X_j where $i, j \in \{1, 2, \dots, d\}$, X_j is a descendant node of X_i iff $X_i^{t_a} \not\perp\!\!\!\perp X_j^{t_b} \mid \mathbf{X}_{\otimes i}^{t_a}$.*

Proof. From the the non-zero *time-lagged effect* and Assumptions 3.1, 3.2, and 3.3, we can infer that:

- (a) The effect of $X_i^{t_a}$ on $X_j^{t_b}$ is non-zero, i.e., $X_i^{t_a} \dashrightarrow X_j^{t_b}$;
- (b) Under Markov property, $\mathbf{X}^\tau \not\rightarrow X_j^{t_b}$ for $\tau < t_a < t_b$;
- (c) Under acyclic assumption, $X_i^{t_a} \not\rightarrow \mathbf{an}_i^{t_b}$ for $t_a < t_b$;
- (d) Under stationary time series, $X_i^{t_a} \dashrightarrow X_j^{t_a} \dashrightarrow X_j^{t_b}$.

Under conditions (a), (b), (c) and (d), if $X_j^{t_b} \in \mathbf{an}_i^{t_b}$, then there are only two causal paths between $X_i^{t_a}$ and $X_j^{t_b}$: $X_i^{t_a} \leftarrow \mathbf{an}_i^{t_a} \dashrightarrow X_j^{t_b}$ and $X_i^{t_a} \dashrightarrow \{X_i^{t_b}, \mathbf{de}_i^{t_b}\} \leftarrow X_j^{t_b}$. Hence, once we cut off all backdoor paths by controlling the conditional set $\mathbf{an}_i^{t_a}$, then the confounding effect between $X_i^{t_a}$ and $X_j^{t_b}$ would be eliminated, leading to $X_i^{t_a} \perp\!\!\!\perp X_j^{t_b} \mid \mathbf{an}_i^{t_a}$. Similarity, if $X_j^{t_b} \in \mathbf{sib}_i^{t_b}$, then the summary backdoor path is $X_i^{t_a} \leftarrow \mathbf{an}_i^{t_a} \dashrightarrow \mathbf{an}_j^{t_b} \dashrightarrow X_j^{t_b}$. In summary, if $X_j^{t_b}$ is a non-descendant node of $X_i^{t_b}$, then $X_i^{t_a} \perp\!\!\!\perp X_j^{t_b} \mid \mathbf{an}_i^{t_a}$. In turn, given the condition $X_i^{t_a} \not\perp\!\!\!\perp X_j^{t_b} \mid \mathbf{an}_i^{t_a}$, $X_j^{t_b}$ is a descendant node of $X_i^{t_b}$. \square



DHT-CIT Algorithm

Algorithm 1 DHT-CIT: Descendant Hierarchical Topology with Conditional Independence Test

Input: Two time-slices $\mathcal{D} = \{X^{t_a}, X^{t_b}\}_{t_a < t_b}$ with d nodes; two significance threshold $\alpha = 0.01$ and $\beta = 0.001$ for conditional independence test and pruning process; the layer index $k = 0$.

Output: One adjacency matrix of descendant hierarchical topology A^{TP} , one DAG \mathcal{G} .

Components: Conditional independence test $\text{HSIC}(\dots)$; and pruning process $\text{CAM}(\dots)$.

Stage 1 - Identifying Descendant Hierarchical Topology:

for $i = 1$ **to** d **do**

Construct the conditional set $X_{\otimes i}^{t_a}$ via an independence test $X_{\otimes i}^{t_a} = \{X_j^{t_a} \mid X_j^{t_a} \perp\!\!\!\perp X_i^{t_a}\}$

for $j = 1$ **to** d **do**

$p_{i,j} = \text{HSIC}(X_i^{t_a}, X_j^{t_b} \mid X_{\otimes i}^{t_a})$

$a_{i,j}^{TP} = \mathbb{I}(p_{i,j} \leq \alpha)$

end for

end for

We obtain $P = \{p_{i,j}\}_{d \times d}$ and $A^{TP} = \{a_{i,j}^{TP}\}_{d \times d}$

Stage 2 - Adjusting the Topological Ordering:

while The causal relationship between the unprocessed nodes is a directed cyclic graph **do**

$k := k + 1$

$X_{M_{i,k}} = \{X_i^{t_a} / X_i^{t_b}, L_{1:k-1}\}$

$X_i^{t_b} \in L_k$, if $a_{i,j}^{TP} = 0$ for all $j \in M_{i,k}$

while $L_k = \emptyset$ **do**

$p_{i^*,j^*} := 2\alpha$ and $a_{i^*,j^*}^{TP} = 0$, $(i^*, j^*) = \arg \max_{i,j} (p_{i,j} \leq \alpha)$

$X_i^{t_b} \in L_k$, if $a_{i,j}^{TP} = 0$ for all $j \in M_{i,k}$

end while

We obtain $P = \{p_{i,j}\}_{d \times d}$ and $A^{TP} = \{a_{i,j}^{TP}\}_{d \times d}$

end while

Stage 3 - Pruning Spurious Edges:

We obtain $\mathcal{G} = \text{CAM}(\mathcal{D}, A^{TP}, \beta)$

Return: A^{TP} and \mathcal{G}

Data-nodes-edges with different complex functions: Sin-10-10, Sin-20-20, Sigmoid-10-10, Poly-10-10

Table 1. The results (mean \pm std) on Sin- d - e using observational data ($\mathcal{D} = \{X^1, X^2\}$).

	Sin-10-10 Graph with Observational Data ($\mathcal{D} = \{X^1, X^2\}$)					Sin-20-20 Graph with Observational Data ($\mathcal{D} = \{X^1, X^2\}$)				
Method	SHD \downarrow	SID \downarrow	F1-Score \uparrow	Dis. \downarrow	#Prune \downarrow	SHD \downarrow	SID \downarrow	F1-Score \uparrow	Dis. \downarrow	#Prune \downarrow
PC	12.8 \pm 5.03	43.6 \pm 9.94	0.56 \pm 0.12	3.51 \pm 0.72	-	21.5 \pm 6.75	98.2 \pm 31.8	0.61 \pm 0.11	4.59 \pm 0.69	-
FCI	15.3 \pm 3.77	71.0 \pm 11.5	0.54 \pm 0.09	3.89 \pm 0.46	-	30.5 \pm 4.09	237. \pm 59.1	0.54 \pm 0.05	5.51 \pm 0.37	-
GOLEM	0.50 \pm 0.80	1.80 \pm 2.70	0.97 \pm 0.03	0.38 \pm 0.59	-	1.30 \pm 1.10	5.60 \pm 4.40	0.97 \pm 0.03	0.93 \pm 0.66	-
NOTEARS	1.20 \pm 0.60	2.30 \pm 1.20	0.94 \pm 0.02	1.02 \pm 0.30	-	2.60 \pm 1.49	6.00 \pm 3.40	0.94 \pm 0.03	1.55 \pm 0.46	-
ReScore	1.00 \pm 0.63	1.40 \pm 1.36	0.95 \pm 0.03	0.88 \pm 0.47	-	2.00 \pm 0.77	5.10 \pm 2.90	0.95 \pm 0.01	1.38 \pm 0.28	-
Granger	31.3 \pm 11.6	66.8 \pm 30.8	0.21 \pm 0.04	5.48 \pm 1.10	-	104 \pm 20.7	368 \pm 8.82	0.10 \pm 0.03	10.1 \pm 1.01	-
VarLiNGAM	35.0 \pm 0.00	69.4 \pm 3.20	0.36 \pm 0.00	5.91 \pm 0.00	-	170 \pm 0.00	339 \pm 3.20	0.19 \pm 0.00	13.0 \pm 0.00	-
CD-NOD	5.40 \pm 0.92	15.5 \pm 4.70	0.74 \pm 0.04	2.32 \pm 0.19	-	-	-	-	-	-
CAM	3.70 \pm 2.95	13.2 \pm 10.6	0.84 \pm 0.13	1.79 \pm 0.74	80.00 \pm 0.00	10.3 \pm 6.50	41.6 \pm 34.7	0.79 \pm 0.12	3.07 \pm 0.98	360.0 \pm 0.00
SCORE	5.60 \pm 3.92	21.2 \pm 16.1	0.78 \pm 0.14	2.25 \pm 0.78	35.80 \pm 0.98	7.40 \pm 2.41	31.3 \pm 21.7	0.85 \pm 0.04	2.68 \pm 0.47	172.1 \pm 0.22
DHT-CIT	1.00 \pm 1.22	3.20 \pm 3.70	0.95 \pm 0.05	0.68 \pm 0.72	13.20 \pm 4.30	1.00 \pm 1.32	3.10 \pm 4.40	0.98 \pm 0.03	0.51 \pm 0.61	30.60 \pm 7.70

* CD-NOD on Sin-20-20 takes over 5 hours and #Prune on one-stage methods is not meaningful. We don't discuss these results and represent them with '-'.

Table 2. The results (mean \pm std) on Sigmoid-10-10 & Poly-10-10 data.

	Sigmoid-10-10 data with Gaussian Noise ($\mathcal{D} = \{X^1, X^2\}$)					Poly-10-10 data with Gaussian Noise ($\mathcal{D} = \{X^1, X^2\}$)				
Method	SHD \downarrow	SID \downarrow	F1-Score \uparrow	Dis. \downarrow	#Prune \downarrow	SHD \downarrow	SID \downarrow	F1-Score \uparrow	Dis. \downarrow	#Prune \downarrow
GOLEM	4.30 \pm 2.19	18.4 \pm 7.92	0.78 \pm 0.11	2.00 \pm 0.51	-	19.00 \pm 4.00	59.4 \pm 13.6	0.20 \pm 0.12	4.33 \pm 0.45	-
NOTEARS	12.5 \pm 5.40	45.3 \pm 17.9	0.46 \pm 0.21	3.44 \pm 0.78	-	17.8 \pm 5.36	56.4 \pm 16.9	0.23 \pm 0.18	4.16 \pm 0.64	-
ReScore	12.2 \pm 4.30	45.6 \pm 14.4	0.45 \pm 0.17	3.43 \pm 0.63	-	17.7 \pm 4.73	57.3 \pm 14.1	0.22 \pm 0.15	4.16 \pm 0.56	-
CAM	3.70 \pm 3.43	10.4 \pm 7.86	0.82 \pm 0.17	1.55 \pm 1.20	80.00 \pm 0.00	8.00 \pm 4.69	19.8 \pm 7.88	0.63 \pm 0.21	2.68 \pm 0.95	80.00 \pm 0.00
SCORE	9.90 \pm 3.81	32.8 \pm 11.6	0.56 \pm 0.16	3.09 \pm 0.61	38.90 \pm 1.60	18.90 \pm 4.33	40.4 \pm 10.9	0.23 \pm 0.13	4.32 \pm 0.52	42.20 \pm 1.48
DHT-CIT	0.67 \pm 1.12	1.80 \pm 2.99	0.96 \pm 0.06	0.46 \pm 0.72	8.67 \pm 2.92	3.22 \pm 3.15	10.8 \pm 5.69	0.84 \pm 0.15	1.51 \pm 1.03	11.33 \pm 3.87

Data-nodes-edges for denser graphs: Sigmoid-10-20 , Sigmoid-10-30, Sigmoid-10-40 , Sigmoid-20-60, Sigmoid-20-100

Table 4. The experiments (mean \pm std) on Sigmoid-d-e using observations $D = \{X^{t_a}, X^{t_b}\}$ with Subsampling Rate $u = t_b - t_a$.

	Sigmoid-10-20 on $\mathcal{D} = \{X^1, X^4\}$			Sigmoid-10-30 on $\mathcal{D} = \{X^1, X^4\}$			Sigmoid-10-40 on $\mathcal{D} = \{X^1, X^4\}$		
Method	SHD \downarrow	SID \downarrow	#Prune \downarrow	SHD \downarrow	SID \downarrow	#Prune \downarrow	SHD \downarrow	SID \downarrow	#Prune \downarrow
GOLEM	10.70 \pm 2.93	63.70 \pm 11.85	-	26.90 \pm 4.83	71.40 \pm 8.59	-	35.20 \pm 3.25	67.00 \pm 11.20	-
SCORE	15.10 \pm 3.65	53.10 \pm 12.95	31.20 \pm 1.60	14.80 \pm 5.71	46.40 \pm 8.04	20.80 \pm 1.60	23.60 \pm 2.24	38.40 \pm 13.46	10.30 \pm 2.00
DHT-CIT	6.30\pm2.90	25.30\pm13.66	13.50\pm2.91	14.10\pm4.46	38.80\pm11.02	11.10\pm1.58	23.60\pm2.24	41.20\pm6.90	6.80\pm2.48
	Sigmoid-20-20 on $\mathcal{D} = \{X^1, X^4\}$			Sigmoid-20-60 on $\mathcal{D} = \{X^1, X^4\}$			Sigmoid-20-100 on $\mathcal{D} = \{X^1, X^4\}$		
Method	SHD \downarrow	SID \downarrow	#Prune \downarrow	SHD \downarrow	SID \downarrow	#Prune \downarrow	SHD \downarrow	SID \downarrow	#Prune \downarrow
GOLEM	26.0 \pm 5.60	138.0 \pm 47.15	-	60.10 \pm 5.49	322.3 \pm 23.84	-	100.0 \pm 5.32	336.4 \pm 16.19	-
SCORE	8.40 \pm 6.20	39.10 \pm 38.82	173.2 \pm 1.99	37.10 \pm 8.14	257.9 \pm 34.38	144.7 \pm 4.27	57.5 \pm 11.00	266.4 \pm 48.18	112.4 \pm 3.10
DHT-CIT	0.70\pm0.90	3.20\pm3.49	30.10\pm9.84	22.10\pm3.75	173.5\pm38.71	58.8\pm6.52	53.5\pm8.43	233.3\pm31.78	75.4\pm6.05
	Sigmoid-10-20 on $\mathcal{D} = \{X^2, X^4\}$			Sigmoid-10-20 on $\mathcal{D} = \{X^2, X^6\}$			Sigmoid-10-20 on $\mathcal{D} = \{X^2, X^{10}\}$		
Method	SHD \downarrow	SID \downarrow	#Prune \downarrow	SHD \downarrow	SID \downarrow	#Prune \downarrow	SHD \downarrow	SID \downarrow	#Prune \downarrow
GOLEM	17.40 \pm 4.96	58.40 \pm 13.81	-	21.60 \pm 4.50	67.20 \pm 11.41	-	21.80 \pm 4.87	69.70 \pm 10.66	-
SCORE	12.20 \pm 1.78	47.20 \pm 5.21	29.20 \pm 0.40	15.80 \pm 3.76	54.70 \pm 10.17	30.00 \pm 1.20	22.30 \pm 4.43	67.60 \pm 10.34	32.30 \pm 2.40
DHT-CIT	8.30\pm3.82	26.00\pm12.77	0.46\pm0.72	8.30\pm1.55	38.10\pm5.19	11.20\pm2.28	14.60\pm3.75	36.20\pm12.86	13.60\pm0.92

Data-nodes-edges for large graphs: Sin-50-100, Sin-100-100

Table 7. The experiments on Sin-50-100 & Sin-100-100 datasets.

	Sin-50-100 data with Gauss noise ($\mathcal{D} = \{X^1, X^3\}$)			
Method	SID↓	SHD↓	#Prune↓	Running Time(s)↓
SCORE	247.0 \pm 102.5	23.0 \pm 8.56	1127. \pm 2.06	1027s
DHT-CIT	2039. \pm 84.14	234. \pm 2.71	397.6 \pm 25.54	3217s
DHT-CIT (50 Intervention)	203.0 \pm 61.1	14.8 \pm 3.90	149.0 \pm 27.00	357s
DHT-CIT+SCORE	97.4 \pm 101.6	7.60 \pm 5.28	352.6 \pm 23.69	1249s
DHT-CIT+SCORE (10 Intervention)	53.2 \pm 20.29	4.80 \pm 0.98	284.0 \pm 26.58	1109s
	Sin-100-100 data with Gauss noise ($\mathcal{D} = \{X^1, X^3\}$)			
Method	SID↓	SHD↓	#Prune↓	Running Time(s)↓
SCORE	381. \pm 156.5	28.67 \pm 4.5	4850 \pm 0.2	4689s
DHT-CIT	2377 \pm 427	218. \pm 12.9	787.0 \pm 49.1	19655s
DHT-CIT (100 Intervention)	5.33 \pm 7.50	1.00 \pm 1.41	347.0 \pm 9.10	1074s
DHT-CIT+SCORE	28.67 \pm 9.53	4.67 \pm 0.47	925.0 \pm 83.3	6342s
DHT-CIT+SCORE (10 Intervention)	14.67 \pm 11.9	3.33 \pm 1.25	797.3 \pm 29.4	6108s

Real Data

Variable	Description
$PM_{2.5}(T)$	Annual county PM2.5 concentration, $\mu g/m^3$
$CMR(Y)$	Annual county cardiovascular mortality rate, deaths/100,000 person-years
$Unemploy(X_1)$	Civilian labor force unemployment rate in 2010
$Income(X_2)$	Median household income in 2009
$Female(X_3)$	Family households - female householder, no spouse present in 2010 / Family households in 2010
$Vacant(X_4)$	Vacant housing units in 2010 / Total housing units in 2010
$Owner(X_5)$	Owner-occupied housing units - percent of total occupied housing units in 2010
$Edu(X_6)$	Educational attainment - persons 25 years and over - high school graduate (includes equivalency) in 2010
$Poverty(X_7)$	Families below poverty level in 2009

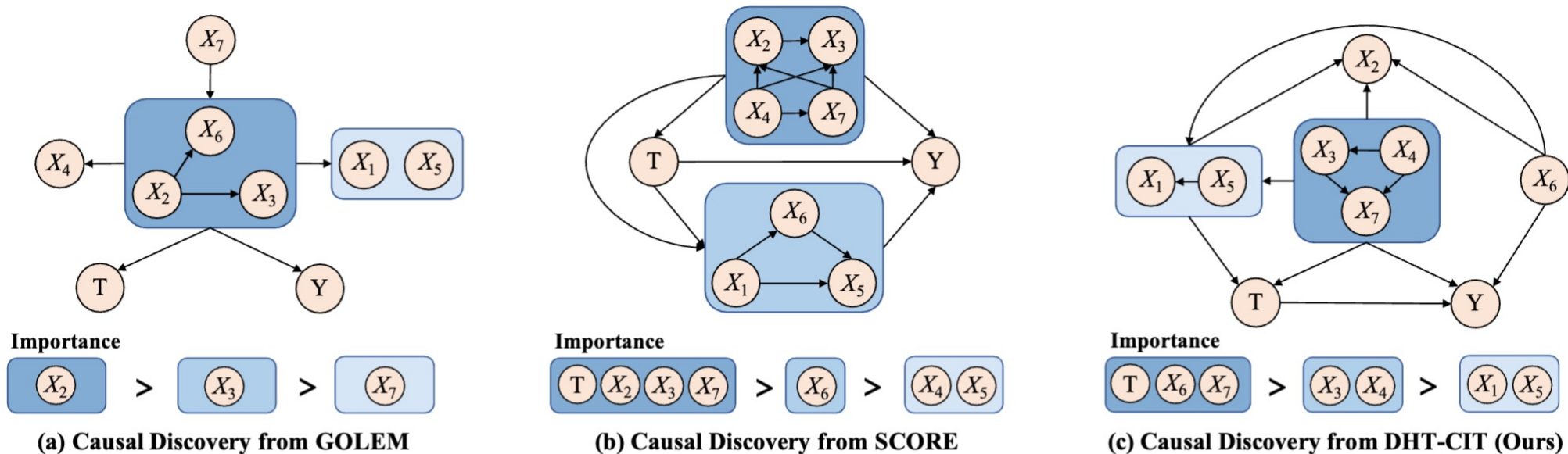


Figure 2. Causal Discovery on the PM-CMR Dataset.

Thanks

Acknowledgement

This work was supported by the National Natural Science Foundation of China (62441605, 62376243, 62037001, U20A20387, 623B2002), and the Starry Night Science Fund of Zhejiang University Shanghai Institute for Advanced Study (SN-ZJU-SIAS-0010).