

# Moment Matters: Mean and Variance Causal Graph Discovery from Heteroscedastic Observational Data



Yoichi Chikahara<sup>1</sup>.

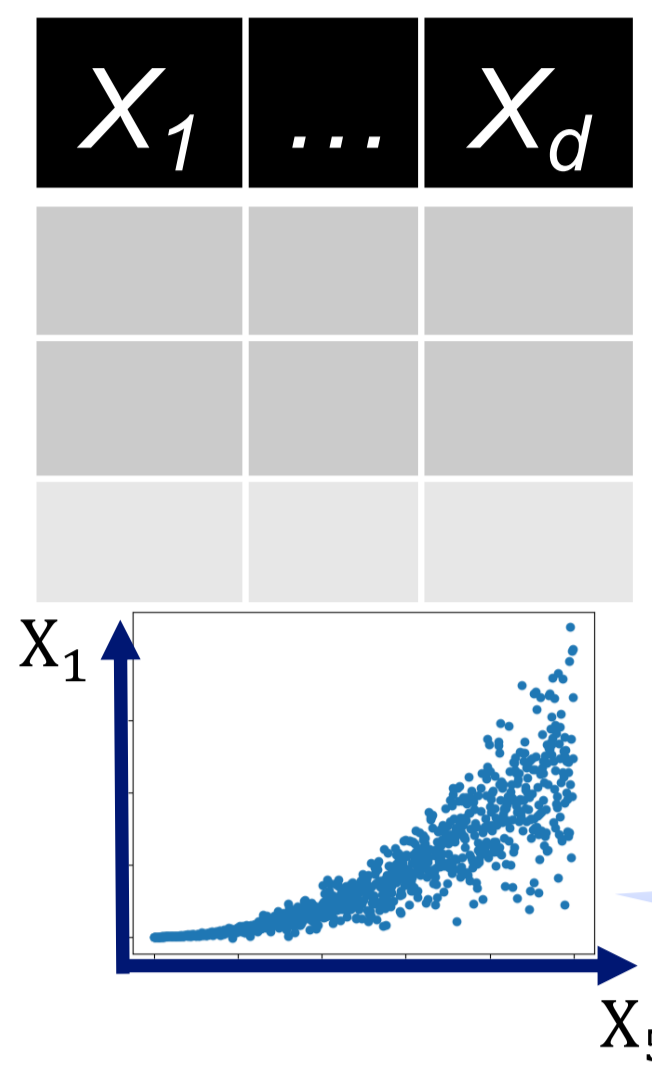
<sup>1</sup>Communication Science Laboratories, NTT, Inc.

This work was supported by JST ACT-X (JPMJAX23CF).



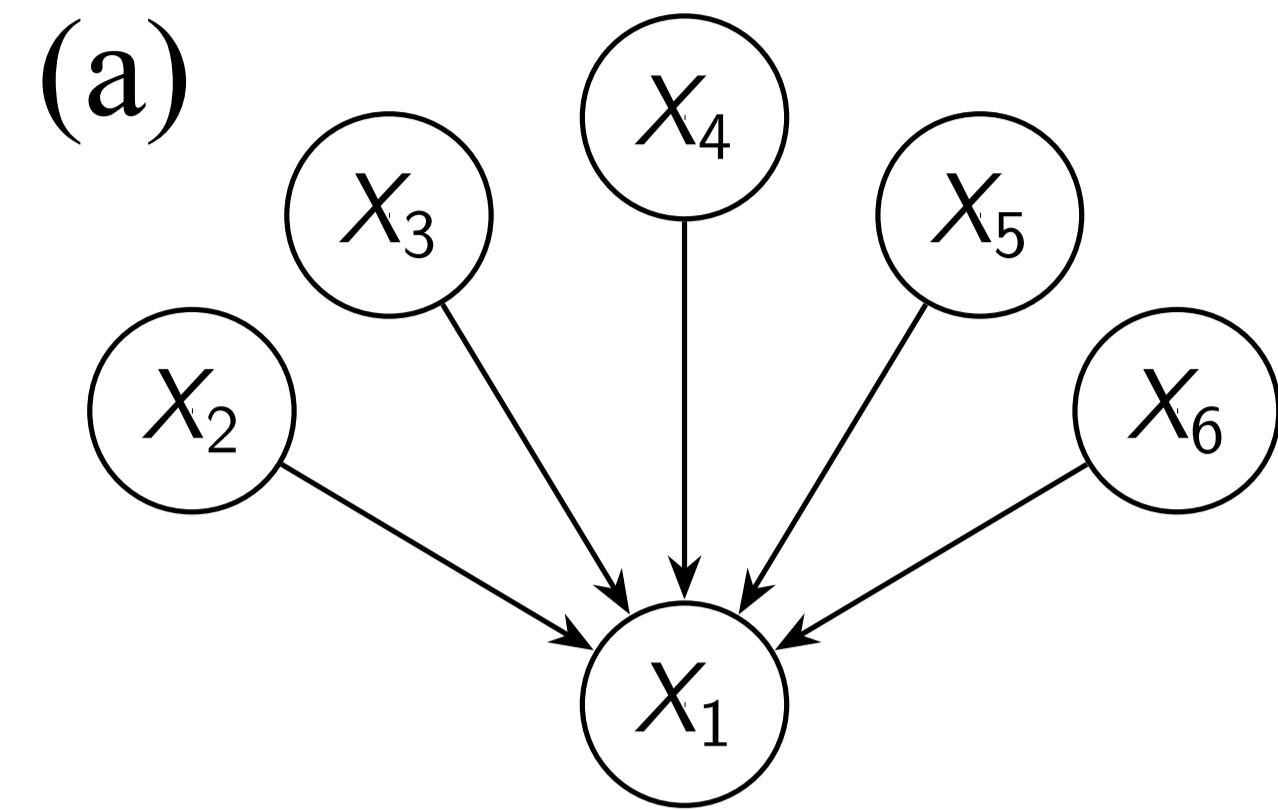
## Standard Graphs vs. Mean & Variance Graphs: Which ones are more interpretable?

Observational data  $D$

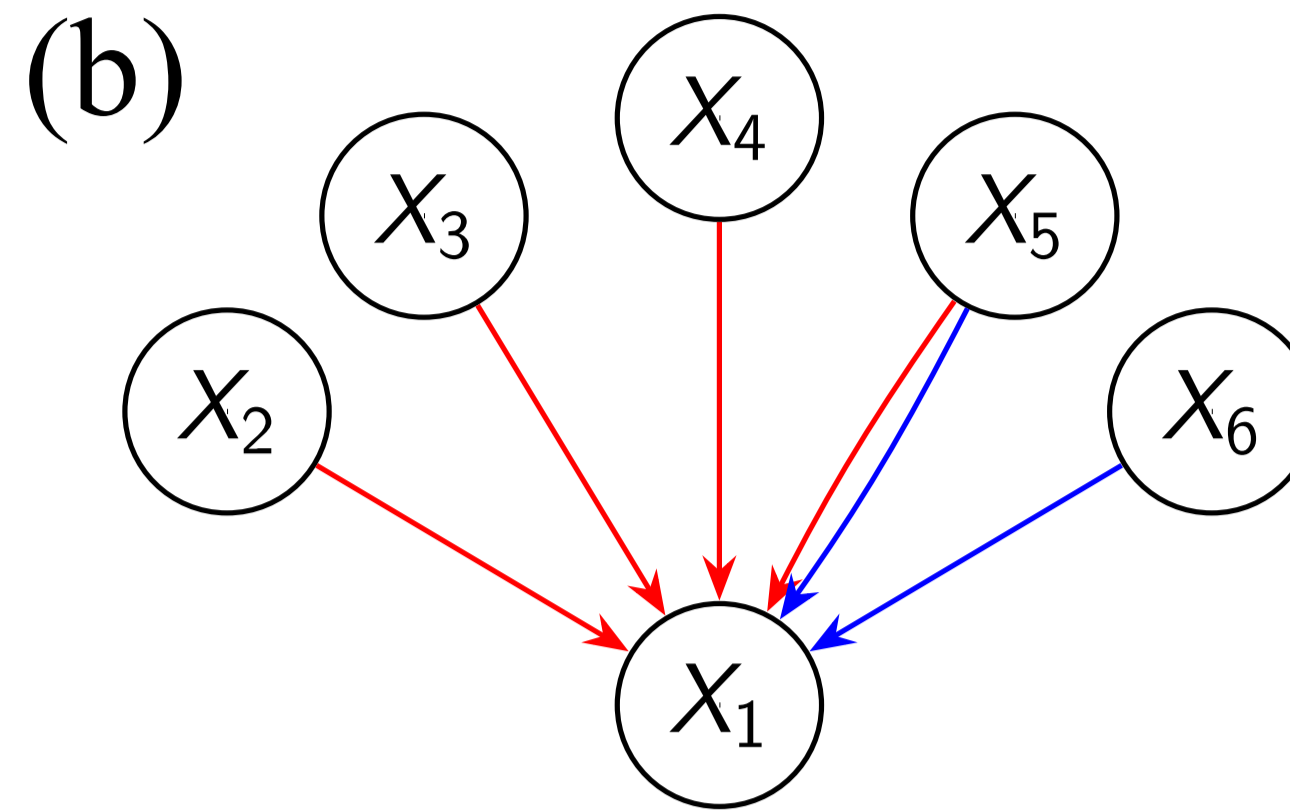


Heteroscedastic data

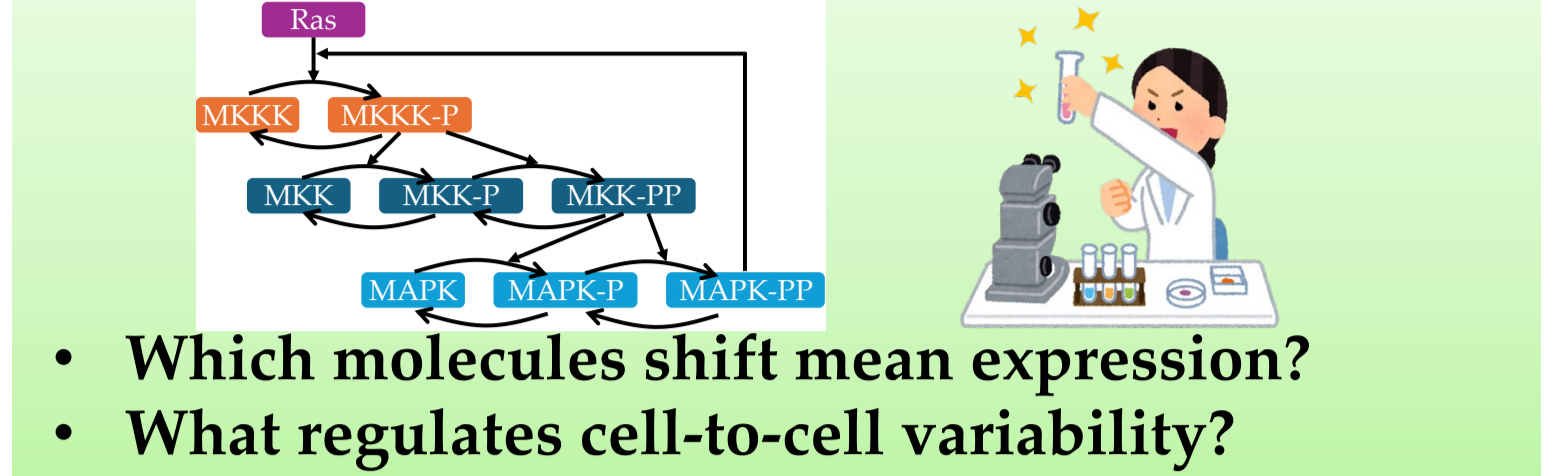
Standard causal graph



Mean/Variance causal graphs

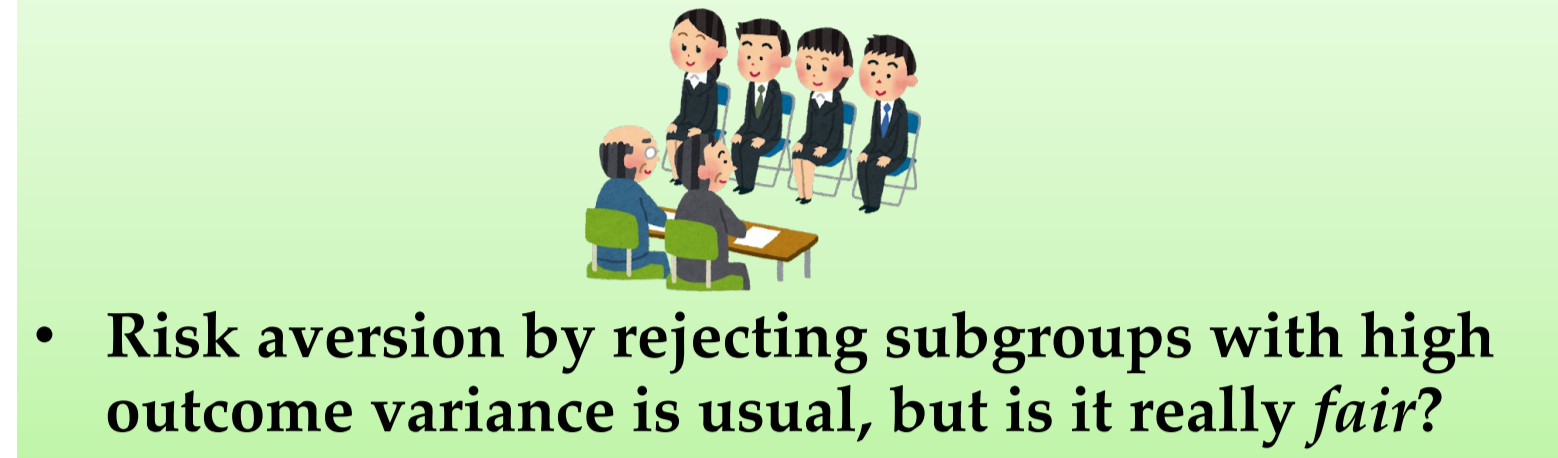


Understanding in biology & medicine:



- Which molecules shift mean expression?
- What regulates cell-to-cell variability?

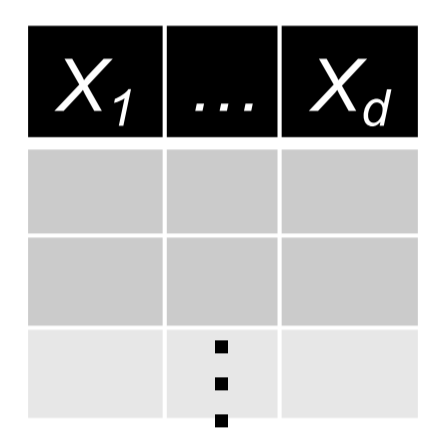
Latent discrimination detection in decision-making:



- Risk aversion by rejecting subgroups with high outcome variance is usual, but is it really fair?

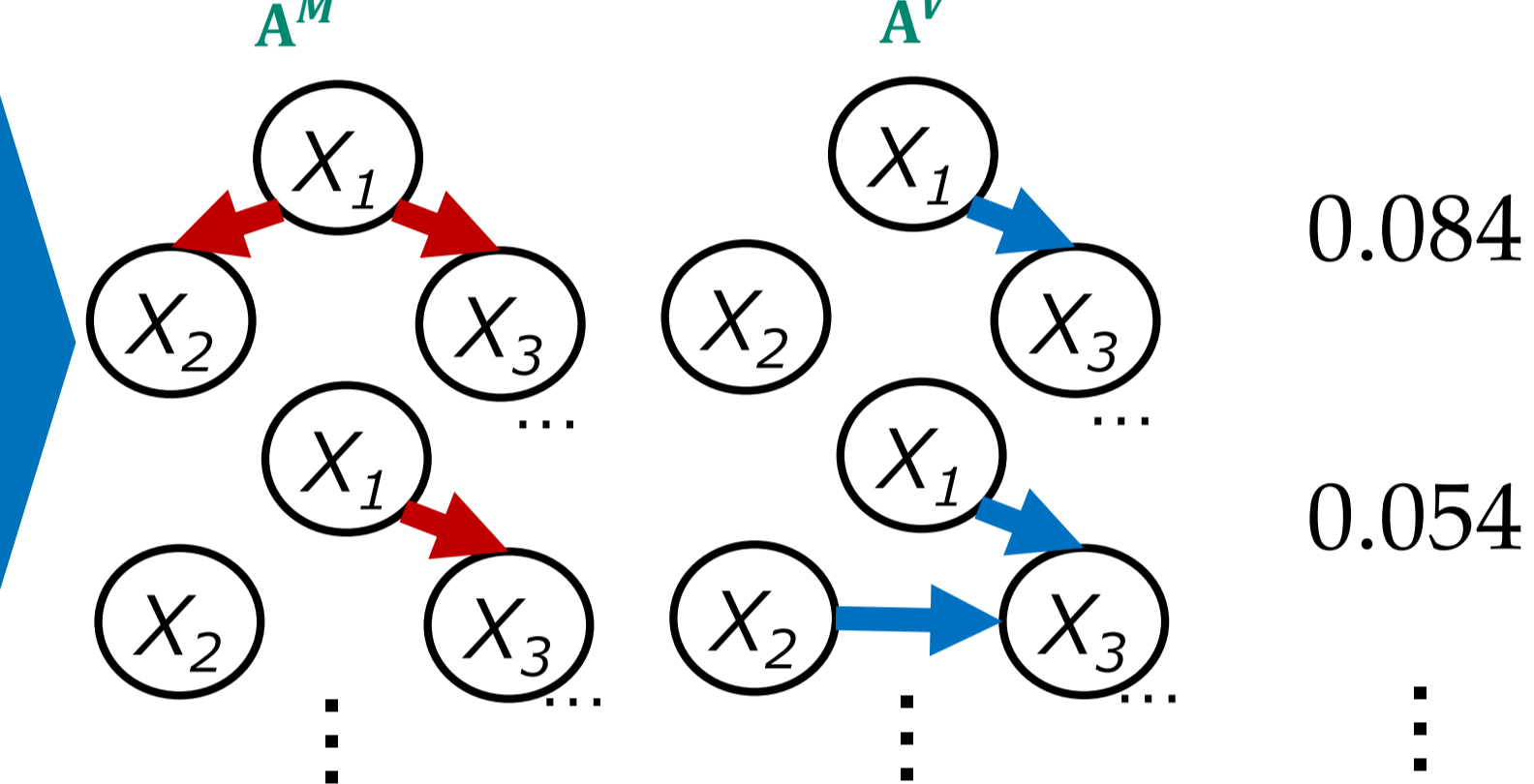
## Problem Setup: Bayesian Causal Discovery

**Input** Observational data  $D$



**Output** Posterior over causal graphs  $P(\mathbf{A}^M, \mathbf{A}^V | D)$

Mean causal graph Variance causal graph Probability



**Input** Prior knowledge about true causal graph structure

Edge of  $G^M$  Edge of  $G^V$  Order  
 $X_1 \rightarrow X_3$   $X_1 \rightarrow X_3$   $\pi(X_1) < \pi(X_2)$

- **Advantage:** Can evaluate probability on structural feature of causal graphs

$$\mathbb{E}_{P(\mathbf{A}^M, \mathbf{A}^V | D)}[f(\mathbf{A}^M, \mathbf{A}^V)]; \text{ e.g., } f(\mathbf{A}^M, \mathbf{A}^V) = \mathbb{I}(X_i \rightarrow X_j \text{ in } \mathbf{A}^M) \text{ or } \mathbb{I}(X_i \rightsquigarrow X_j \text{ in } \mathbf{A}^V) \text{ (path)}$$

## Proposed Method: Variational Approach

**Goal** Approximate the posterior as  $P(\mathbf{A}^M, \mathbf{A}^V | D) \approx P_{\Phi}(\mathbf{A}^M, \mathbf{A}^V)$

**Identifiability Conditions** Mean & Variance graphs are DAGs with shared permutation

$$\begin{aligned} \text{Adjacency matrix (DAG)} &= \text{Permutation matrix (Order)} \cdot \text{Upper-triangular matrix (Edge)} \\ \mathbf{A}^M &= \mathbf{\Pi}^T \mathbf{U}^M \mathbf{\Pi} \\ \mathbf{A}^V &= \mathbf{\Pi}^T \mathbf{U}^V \mathbf{\Pi} \end{aligned} \quad \times A_{i,j} = U_{\pi_i, \pi_j}$$

**Factorized Variational Distribution**  $P(\mathbf{A}^M, \mathbf{A}^V) = \sum_{\mathbf{U}^M, \mathbf{U}^V, \mathbf{\Pi}} P(\mathbf{U}^M) P(\mathbf{U}^V) P(\mathbf{\Pi})$

- Summation is computationally intractable  $\otimes$
- Sampling operations for binary matrices  $\mathbf{U}^M, \mathbf{U}^V, \mathbf{\Pi}$  are NOT differentiable  $\otimes$

**Differentiable sampling** Approximate gradient using differentiable function

**Gumbel-softmax sampling** [Jang+; ICLR2016]:  
 Idea: Softmax(Unnormalized log prob. + Gumbel variables)  
 For each  $i, j \in \{1, \dots, d\}$ ,  

$$U_{i,j} \in [0, 1] = \frac{e^{\log \phi_{ij} + G_{i,j,1}/\tau}}{e^{\log \phi_{ij} + G_{i,j,1}/\tau} + e^{-\log \phi_{ij} + G_{i,j,2}/\tau}}$$
 Forward: Discrete sampling by taking argmax  
 Backward: Use  $U_{i,j} \in [0, 1]$  or  $\tilde{\pi} \in [0, 1]^{d \times d}$  for gradient approximation

**Gumbel-topK + SoftSort** [Prilo+; ICML2020]:  
 Idea: SoftSort(Log prob. + Gumbel variables)  

$$\tilde{\pi} \in \mathbb{R}^d = \text{SoftSort}(\log \phi + G) \rightarrow \hat{\pi} \in [0, 1]^{d \times d}$$

**Objective Function**  $\max_{\Phi, \Theta} \mathbb{E}_{\mathbf{A}^M, \mathbf{A}^V \sim P_{\Phi}} [\log P_{\Theta}(D | \mathbf{A}^M, \mathbf{A}^V)] - \lambda \Omega_{\Phi, \Theta}$

$$\begin{aligned} \log P_{\Theta}(D | \mathbf{A}^M, \mathbf{A}^V) &= -\frac{1}{n} \sum_{i=1}^n \sum_{j=1}^d \left( \frac{(x_{i,j} - m_j(x_{i, \text{pa}^M(j)}; \theta_j^M))^2}{2(v_j(x_{i, \text{pa}^V(j)}; \theta_j^V))^2} + \log v_j(x_{i, \text{pa}^V(j)}; \theta_j^V) \right) \\ \Omega_{\Phi, \Theta} &= \lambda_{\Phi} \text{KL}(P_{\Phi}(\mathbf{U}^M, \mathbf{U}^V) || P(\mathbf{U}^M, \mathbf{U}^V)) + \lambda_{\Theta} \sum_{j=1}^d \|\theta_j^M\|_2^2 + \lambda_{\Theta^V} \sum_{j=1}^d \|\theta_j^V\|_2^2 \end{aligned}$$

**Algorithm 1** Parameter Learning for Mean and Variance Causal Graphs

**Require:** Observational dataset  $D = \{x_1, \dots, x_n\}$ ; parameters  $\Phi, \Theta$ ; node-ordering constraints  $\mathcal{O}$

- 1: Initialize all parameters  $\Phi = \{\phi^M, \phi^V, \psi\}$  and  $\Theta = \{\theta_j^M, \theta_j^V\}_{j=1}^d$
- 2: **while** Stopping criterion not met **do**
- 3: **while** Updates for parameters  $\Theta \setminus \{\theta_j^M\}_{j=1}^d$  and  $\Theta \setminus \{\theta_j^V\}_{j=1}^d$  not converged **do**
- 4: ① Sample mean and variance DAGs  $\mathbf{A}^M, \mathbf{A}^V \sim P_{\Phi}(\mathbf{A}^M, \mathbf{A}^V)$
- 5: ② Compute ELBO-based objective function in Eq. (8)
- 6: ③ Update parameters  $\Theta \setminus \{\theta_j^M\}_{j=1}^d$  and  $\Theta \setminus \{\theta_j^V\}_{j=1}^d$  using scaled gradient
- 7: **if** node-ordering constraints  $\mathcal{O}$  are provided **then**
- 8: Project updated parameters  $\psi$  onto the feasible set by Eq. (11)  $\mathcal{I}(\psi) = \{\psi_i + c_{i,j} \leq \psi_j \mid (i, j) \in \mathcal{S}\}$
- 9: **end if**
- 10: **end while**
- 11: **while** Updates for variance-related parameters  $\{\theta_j^V\}_{j=1}^d$  and  $\{\phi^V\}$  not converged **do**
- 12: ① Sample mean and variance DAGs  $\mathbf{A}^M, \mathbf{A}^V \sim P_{\Phi}(\mathbf{A}^M, \mathbf{A}^V)$
- 13: ② Compute ELBO-based objective function in Eq. (8)
- 14: ③ Update variance-related parameters  $\{\theta_j^V\}_{j=1}^d$  and  $\{\phi^V\}$  using standard gradient
- 15: **end while**
- 16: **end while**
- 17: **return** Learned parameters  $\Phi$  and  $\Theta$

**Order knowledge**  
 $\pi(X_1) < \pi(X_2)$

## Mean-Variance HNM, Mean & Variance Graphs

**Assume** each  $X_j$ 's values are determined by the mean-variance HNM:

$$X_j = m_j(\mathbf{X}_{\text{pa}^M(j)}) + v_j(\mathbf{X}_{\text{pa}^V(j)}) E_j \quad \text{for } j = 1, \dots, d,$$

mean function      variance function      zero-mean noise

$$\mathbb{E}[X_j | \mathbf{X}_{\text{pa}(j)}] = m_j(\mathbf{X}_{\text{pa}^M(j)}) \quad \mathbb{V}[X_j | \mathbf{X}_{\text{pa}(j)}] = \mathbb{V}[E_j] (v_j(\mathbf{X}_{\text{pa}^V(j)}))^2.$$

**Definition: Mean & Variance Causal Graphs**

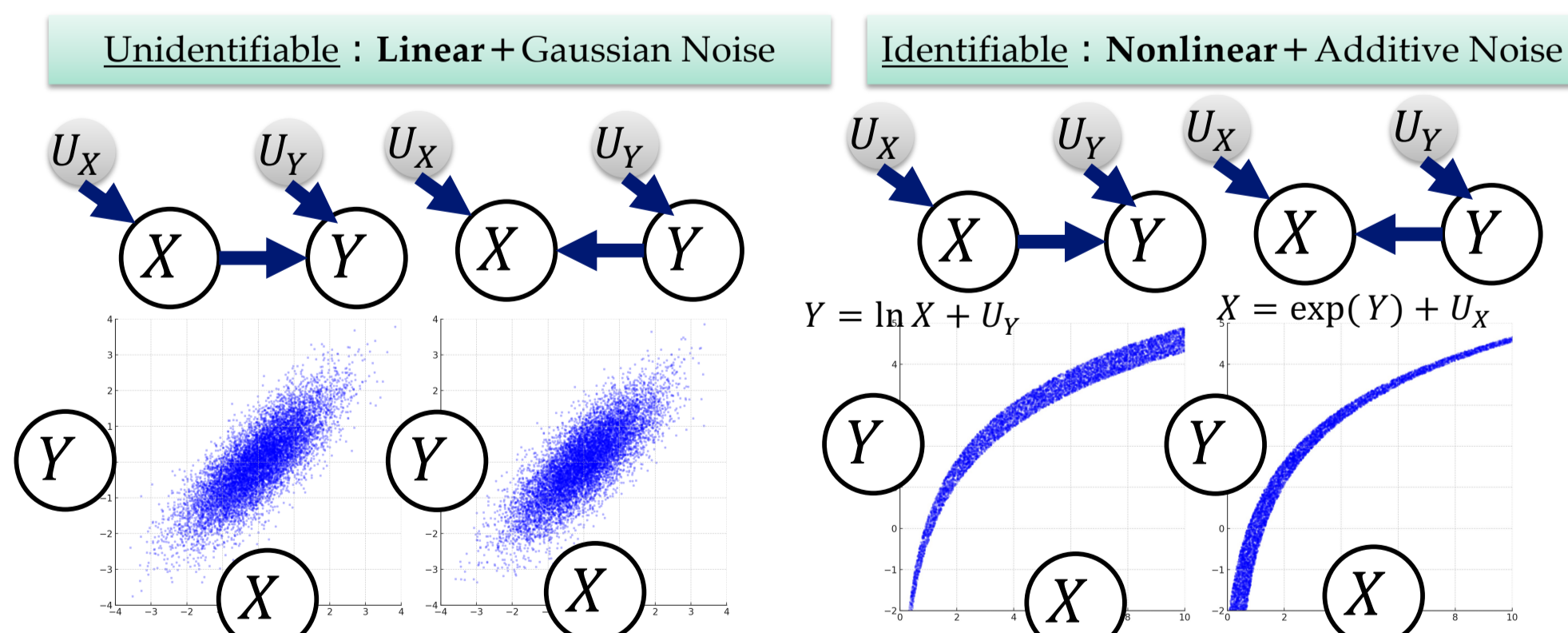
Mean causal graph  $G^M$  has edge  $X_i \rightarrow X_j$  if and only if  $X_i \in \mathbf{X}_{\text{pa}^M(j)}$ ,

variance causal graph  $G^V$  contains edge  $X_i \rightarrow X_j$  if and only if  $X_i \in \mathbf{X}_{\text{pa}^V(j)}$ .

**Adjacency Matrices:**  $\mathbf{A}^M, \mathbf{A}^V \in \{0, 1\}^{d \times d}$

## Identifiability Conditions

**Unidentifiable = Different models can yield the same observational data distribution**



**Mean & variance causal graphs are identifiable (= uniquely determined) from observational data under the following sufficient conditions:**

**Theorem 3.5:** Under Assumptions 3.1, 3.2, 3.3, and 3.4, mean and variance causal graphs  $G^M$  and  $G^V$  are identifiable from observational distribution  $P(\mathbf{X})$  if for  $j = 1, \dots, d$ , (A)  $m_j$  is a nonlinear function, (B)  $v_j$  is a piecewise function, but not a constant function, and (C)  $E_j$  is a Gaussian noise.

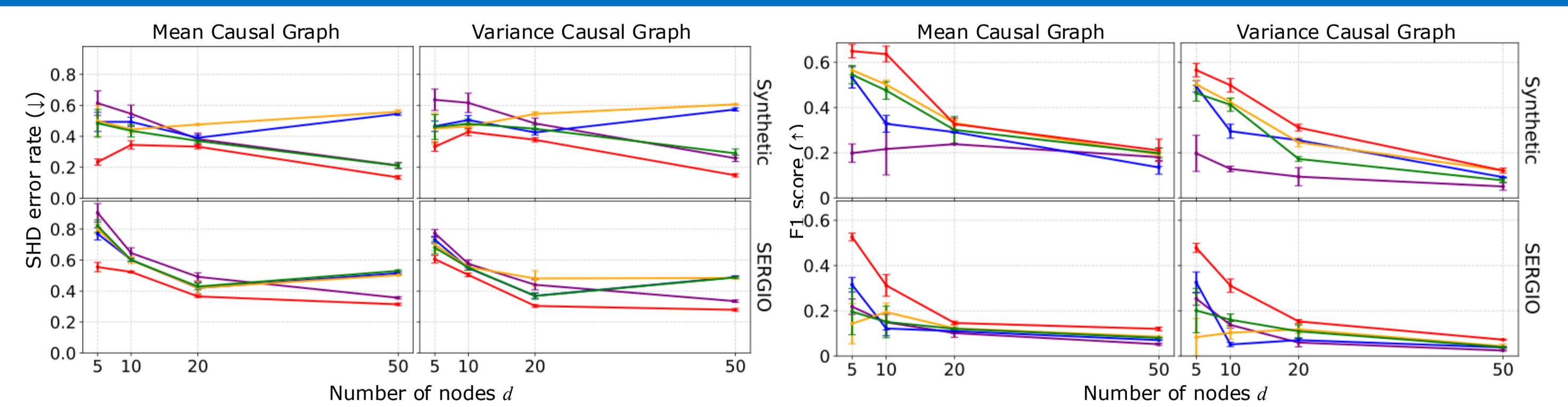
**Assumption 3.1** (Causal sufficiency): Exogenous noises satisfy  $E_i \perp E_j$  for any  $i, j \in \{1, \dots, d\}$ .

**Assumption 3.3:** Mean and variance graphs  $G^M$  and  $G^V$  are directed acyclic graphs (DAGs).

**Assumption 3.2** (Causal minimality):  $P(\mathbf{X})$  satisfies the causal minimality condition with respect to moment-agnostic causal graph  $G$ .

**Assumption 3.4** (Shared permutation condition): There exists an identical permutation (a.k.a., topological ordering) of mean and variance causal graphs  $G^M$  and  $G^V$ .

## Experimental Results



	$n = 100$		$n = 200$		$n = 853$	
	SHD ( $\downarrow$ )	F1 ( $\uparrow$ )	SHD ( $\downarrow$ )	F1 ( $\uparrow$ )	SHD ( $\downarrow$ )	F1 ( $\uparrow$ )
MC3	22.6 $\pm$ 1.1	0.19 $\pm$ 0.04	20.3 $\pm$ 1.3	0.20 $\pm$ 0.07	18.9 $\pm$ 0.7	0.14 $\pm$ 0.05
DDS	17.6 $\pm$ 0.7	0.20 $\pm$ 0.03	16.6 $\pm$ 0.6	0.20 $\pm$ 0.02	15.0 $\pm$ 0.7	0.28 $\pm$ 0.02
ICDH	20.0 $\pm$ 1.0	<b>0.21 <math>\pm</math> 0.01</b>	17.5 $\pm$ 0.5	0.22 $\pm$ 0.01	14.5 $\pm$ 0.5	0.27 $\pm$ 0.03
HOST	16.1 $\pm$ 0.5	0.19 $\pm$ 0.01	15.0 $\pm$ 0.3	<b>0.33 <math>\pm</math> 0.02</b>	<b>13.5 <math>\pm</math> 0.8</b>	<b>0.38 <math>\pm</math> 0.03</b>
PROPOSED	<b>16.0 <math>\pm</math> 0.3</b>	0.20 $\pm$ 0.02	<b>14.9 <math>\pm</math> 0.4</b>	0.31 $\pm$ 0.02	13.7 $\pm$ 0.5	0.36 $\pm$ 0.02
PROPOSED +25%	14.9 $\pm$ 0.4	0.34 $\pm$ 0.03	14.8 $\pm$ 0.5	0.35 $\pm$ 0.02	13.4 $\pm$ 0.4	0.37 $\pm$ 0.02
PROPOSED +50%	<b>13.2 <math>\pm</math> 0.5</b>	<b>0.36 <math>\pm</math> 0.02</b>	<b>13.2 <math>\pm</math> 0.3</b>	<b>0.36 <math>\pm</math> 0.02</b>	<b>13.1 <math>\pm</math> 0.2</b>	<b>0.45 <math>\pm</math> 0.03</b>

**Table 4: Estimated posterior probability of MEK  $\rightarrow$  ERK in inferred variance causal graph structure**

	$n = 100$	$n = 200$	$n = 853$
PROPOSED	0.585 $\pm$ 0.009	0.592 $\pm$ 0.008	0.620 $\pm$ 0.019