

---

# Appendix for Bayesian Active Causal Discovery with Multi-Fidelity Experiments

---

Anonymous Author(s)  
Affiliation  
Address  
email

1	<b>Contents</b>	
2	<b>A Monte Carlo Approximation for <math>f(j, v, m)</math></b>	<b>3</b>
3	A.1 Derivation Process for $f(j, v, m)$ . . . . .	3
4	A.2 Sampling from $p(\phi_m D)$ . . . . .	4
5	A.3 Sampling from $p(\phi_m \phi_M, D)$ . . . . .	7
6	A.4 Calculation of $p(\mathbf{x} e, \phi_m)$ . . . . .	8
7	<b>B Bayesian Optimization for Determining <math>(j^*, v^*, m^*)</math></b>	<b>9</b>
8	<b>C Detailed Training Process of ELBO</b>	<b>10</b>
9	C.1 Derivation Process of ELBO . . . . .	10
10	C.2 Estimation of ELBO . . . . .	10
11	C.3 Gaussian Reparameterization Trick . . . . .	12
12	C.4 Gumbel-softmax Reparameterization Trick . . . . .	12
13	C.5 Optimization of ELBO . . . . .	13
14	<b>D Training Process of Constraint based ELBO</b>	<b>13</b>
15	<b>E Proof of Theory 3</b>	<b>14</b>
16	<b>F Proof of Theory 4</b>	<b>14</b>
17	<b>G Algorithm</b>	<b>15</b>
18	<b>H More Experiments</b>	<b>16</b>
19	H.1 Experimental Settings . . . . .	16
20	H.1.1 Datasets . . . . .	16
21	H.1.2 Baselines . . . . .	16
22	H.1.3 Metrics . . . . .	17

23	H.2 Details of Configurations and Computation . . . . .	17
24	H.3 Experiments on DREAM Dataset . . . . .	18
25	H.4 Experiments on More Nodes . . . . .	18
26	<b>I Potentially Negative Social Impact</b>	<b>18</b>

27 **A Monte Carlo Approximation for  $f(j, v, m)$**

28 **A.1 Derivation Process for  $f(j, v, m)$**

29 Considering that the mutual information is not directly tractable, we approximate  $f(j, v, m)$  by:

$$f(j, v, m) = - \frac{1}{\lambda_m \cdot K_1 \cdot L_1} \sum_{k_1=1}^{K_1} \sum_{l_1=1}^{L_1} \log \left[ \frac{1}{C} \sum_{c_1=1}^{C_1} p(\mathbf{x}_m^{(k_1, l_1)} | \phi_m^{(c_1)}, \mathbf{e}) \right] \\ + \frac{1}{\lambda_m \cdot K_2 \cdot L_2 \cdot C_2} \sum_{k_2=1}^{K_2} \sum_{l_2=1}^{L_2} \sum_{c_2=1}^{C_2} \log \left[ p(\mathbf{x}_m^{(k_2, l_2, c_2)} | \phi_M^{(k_2)}, \mathbf{e}) \right],$$

30 where  $\mathbf{e} = \{(j, v), m\}$  is the experiment to be designed,  $\phi_m^{(c_1)}, \phi_m^{(k_1)} \sim p(\phi_m | D)$ ,  $\mathbf{x}_m^{(k_1, l_1)} \sim$   
31  $p(\mathbf{x} | \phi_m^{(k_1)}, \mathbf{e})$ ,  $\phi_M^{(k_2)} \sim p(\phi_M | D)$ ,  $\phi_m^{(k_2, l_2)} \sim p(\phi_m | \phi_M^{(k_2)}, D)$  and  $\mathbf{x}_m^{(k_2, l_2, c_2)} \sim p(\mathbf{x} | \phi_m^{(k_2, l_2)}, \mathbf{e})$ .

32 We present the detailed approximation process as follows:

$$f(j, v, m) = \frac{1}{\lambda_m} I(\mathbf{x}; \phi_M | \mathbf{e}, D) \\ = \frac{1}{\lambda_m} [H(\mathbf{x} | \mathbf{e}, D) - H(\mathbf{x} | \phi_M, \mathbf{e}, D)] \\ = \frac{1}{\lambda_m} [-\mathbb{E}_{p(\mathbf{x} | \mathbf{e}, D)} [\log p(\mathbf{x} | \mathbf{e}, D)] + \mathbb{E}_{p(\phi_M | D)} [\mathbb{E}_{p(\mathbf{x} | \phi_M, \mathbf{e})} [\log p(\mathbf{x} | \mathbf{e}, \phi_M)]]] \\ = \frac{1}{\lambda_m} \underbrace{[-\mathbb{E}_{p(\mathbf{x} | \mathbf{e}, D)} [\log \mathbb{E}_{p(\phi_m | \mathbf{e}, D)} [p(\mathbf{x} | \mathbf{e}, \phi_m)]]]}_E \\ + \frac{1}{\lambda_m} \underbrace{[\mathbb{E}_{p(\phi_M | D)} [\mathbb{E}_{p(\mathbf{x} | \mathbf{e}, \phi_M)} [\log p(\mathbf{x} | \mathbf{e}, \phi_M)]]]}_F$$

33 For part  $E$ , we can estimate it by

$$E = - \frac{1}{\lambda_m \cdot K_1 \cdot L_1} \sum_{k_1=1}^{K_1} \sum_{l_1=1}^{L_1} \log \left[ \frac{1}{C} \sum_{c_1=1}^{C_1} p(\mathbf{x}_m^{(k_1, l_1)} | \phi_m^{(c_1)}, \mathbf{e}) \right],$$

34 where for the first expectation on  $p(\phi_m | \mathbf{e}, D)$ , we first sample  $\phi_m^{(k_1)}$  from  $\phi_m^{(k_1)} \sim p(\phi_m | \mathbf{e}, D)$  for  
35  $K_1$  times, and then for each  $\phi_m^{(k_1)}$ , we sample  $\mathbf{x}_m^{(k_1, l_1)}$  from  $\mathbf{x}_m^{(k_1, l_1)} \sim p(\mathbf{x} | \phi_m^{(k_1)}, \mathbf{e})$  for  $L_1$  times.

36 For the second expectation on  $p(\phi_m | \mathbf{e}, D)$ , we sample  $\phi_m^{(c_1)} \sim p(\phi_m | \mathbf{e}, D)$  for  $C_1$  times.

37 For part  $F$ , we have

$$F = \frac{1}{\lambda_m} \cdot [\mathbb{E}_{p(\phi_M | D)} [\mathbb{E}_{p(\mathbf{x} | \phi_M, \mathbf{e})} [\log p(\mathbf{x} | \phi_M, \mathbf{e})]]] \\ = \frac{1}{\lambda_m} \cdot [\mathbb{E}_{p(\phi_M | D)} \left[ \int p(\mathbf{x} | \phi_M, \mathbf{e}) \log p(\mathbf{x} | \phi_M, \mathbf{e}) d\mathbf{x} \right]] \\ = \frac{1}{\lambda_m} \cdot [\mathbb{E}_{p(\phi_M | D)} \left[ \int \int p(\mathbf{x} | \phi_m, \mathbf{e}) p(\phi_m | \phi_M) d\phi_m \log p(\mathbf{x} | \phi_M, \mathbf{e}) d\mathbf{x} \right]] \\ = \frac{1}{\lambda_m} \cdot [\mathbb{E}_{p(\phi_M | D, \mathbf{e})} \left[ \int \int p(\mathbf{x} | \phi_m, \mathbf{e}) p(\phi_m | \phi_M) \log p(\mathbf{x} | \phi_M, \mathbf{e}) d\mathbf{x} d\phi_m \right]] \\ = \frac{1}{\lambda_m} \cdot [\mathbb{E}_{p(\phi_M | D, \mathbf{e})} \left[ \int \mathbb{E}_{p(\mathbf{x} | \phi_m, \mathbf{e})} [p(\phi_m | \phi_M) \log p(\mathbf{x} | \phi_M, \mathbf{e})] d\phi_m \right]] \\ = \frac{1}{\lambda_m} \cdot [\mathbb{E}_{p(\phi_M | D, \mathbf{e})} [\mathbb{E}_{p(\phi_m | \phi_M)} [\mathbb{E}_{p(\mathbf{x} | \phi_m, \mathbf{e})} [\log p(\mathbf{x} | \phi_M, \mathbf{e})]]]] .$$

38 It can be estimated by

$$\frac{1}{\lambda_m \cdot K_2 \cdot L_2 \cdot C_2} \sum_{k_2=1}^{K_2} \sum_{l_2=1}^{L_2} \sum_{c_2=1}^{C_2} \log \left[ p(\mathbf{x}_m^{(k_2, l_2, c_2)} | \phi_M^{(k_2)}, \mathbf{e}) \right],$$

39 where for the expectation on  $p(\phi_M|D, e)$ , we sample  $\phi_M^{(k_2)}$  from  $\phi_M^{(k_2)} \sim p(\phi_M|e, D)$  for  $K_2$  times.  
 40 For the expectation on  $p(\phi_m|\phi_M)$ , for each  $\phi_M^{(k_2)}$ , we sample  $\phi_m^{(k_2, l_2)}$  from  $\phi_m^{(k_2, l_2)} \sim p(\phi_m|\phi_M^{(k_2)})$   
 41 for  $L_2$  times. For the expectation on  $p(\mathbf{x}|\phi_m, e)$ , for each  $\phi_M^{(k_2)}$  and  $\phi_m^{(k_2, l_2)}$ , we sample  $\mathbf{x}_m^{(k_2, l_2, c_2)}$   
 42 from  $\mathbf{x}_m^{(k_2, l_2, c_2)} \sim p(\mathbf{x}|\phi_m^{(k_2, l_2)}, e)$  for  $C_2$  times.

43 Therefore, we can conclude that  $f(j, v, m)$  can be estimated by

$$f(j, v, m) = - \frac{1}{\lambda_m \cdot K_1 \cdot L_1} \sum_{k_1=1}^{K_1} \sum_{l_1=1}^{L_1} \log \left[ \frac{1}{C} \sum_{c_1=1}^{C_1} p(\mathbf{x}_m^{(k_1, l_1)} | \phi_m^{(c_1)}, e) \right] \\ + \frac{1}{\lambda_m \cdot K_2 \cdot L_2 \cdot C_2} \sum_{k_2=1}^{K_2} \sum_{l_2=1}^{L_2} \sum_{c_2=1}^{C_2} \log \left[ p(\mathbf{x}_m^{(k_2, l_2, c_2)} | \phi_M^{(k_2)}, e) \right],$$

44 where  $\phi_m^{(c_1)}, \phi_m^{(k_1)} \sim p(\phi_m|D)$ ,  $\mathbf{x}_m^{(k_1, l_1)} \sim p(\mathbf{x}|\phi_m^{(k_1)}, e)$ ,  $\phi_M^{(k_2)} \sim p(\phi_M|D)$ ,  $\phi_m^{(k_2, l_2)} \sim$   
 45  $p(\phi_m|\phi_M^{(k_2)}, D)$  and  $\mathbf{x}_m^{(k_2, l_2, c_2)} \sim p(\mathbf{x}|\phi_m^{(k_2, l_2)}, e)$ .

46 Obviously, the above approximation of  $f(j, v, m)$  only depends on  $p(\phi_m|D)$ ,  $p(\phi_m|\phi_M, D)$  and  
 47  $p(\mathbf{x}|\phi_m, e)$ . In the next, we show how to sample from them in Section A.2, A.3 and A.4, respectively.

## 48 A.2 Sampling from $p(\phi_m|D)$

49 Basically, sampling from the posterior of “ $p(\cdot|D)$ ” is not easy. To solve this problem, as mentioned  
 50 in the main paper, we introduce a variational probability “ $q$ ” to approximate “ $p$ ”. In specific, in order  
 51 to sample from  $p(\phi_m|D)$ , we first obtain a sample  $\phi_1$  from  $\phi_1 \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ , and then get  $\phi_m$  from  
 52 the distribution  $q(\phi_m|\phi_1)$ .

53 Since

$$q(\phi_m|\phi_1) = \int_{\phi_{m-1}} \cdots \int_{\phi_2} q(\phi_m, \phi_{m-1}, \dots, \phi_2|\phi_1) d\phi_{m-1} \dots d\phi_2.$$

54 and

$$q(\phi_m, \phi_{m-1}, \dots, \phi_2|\phi_1) = \prod_{i=2}^m q(\phi_i|\phi_{i-1}) \\ q(\phi_i|\phi_{i-1}) = \mathcal{N}(\mathbf{c}_i \phi_{i-1} + \mathbf{d}_i, \sigma_i^2 \mathbf{I}),$$

55 we have  $q(\phi_m|\phi_1)$  is a Gaussian distribution, which is easy for sampling.

56 In our model,  $\mathbf{c}_i$  and  $\sigma_i^2 \mathbf{I}$  are diagonal matrices, which means that the dimensions in  $\phi_i$  are independ-  
 57 ent with each other. We denote  $\phi_i = [\phi_{i,1}, \phi_{i,2} \dots \phi_{i,d}]$ , where  $\phi_{i,j}$  is the  $j$ th element of  $\phi_i$ . Then,  
 58 we have

$$q(\phi_m, \phi_{m-1}, \dots, \phi_2|\phi_1) = \prod_{j=1}^d q(\phi_{m,j}, \phi_{m-1,j}, \dots, \phi_{2,j}|\phi_{1,j}).$$

59 So our target can be converted to calculate the probability  $q(\phi_{m,j}, \phi_{m-1,j}, \dots, \phi_{2,j}|\phi_{1,j})$  for all  
 60 dimensions  $\forall 1 \leq j \leq d$ . Let  $c_{i,j}$ ,  $d_{i,j}$  and  $\sigma_{i,j}^2$  be the  $j$ th element of  $\mathbf{c}_i$ ,  $\mathbf{d}_i$  and  $\sigma_i^2$ , respectively. We  
 61 assume that  $\sigma_{i,j} = \sqrt{c_{i-1,j}^2 + 1} \cdot \sigma_{i-1,j}$  ( $i \geq 4$ ) and  $\sigma_{3,j} = \sigma_{2,j} = e$ , where  $e$  is the hyper-parameter.  
 62 Suppose  $\mu_{i,j}$  is the mean of the Gaussian distribution for  $q(\phi_{i,j}|\phi_{i-1,j})$ , that is,

$$\mu_{i,j} = c_{i,j} \phi_{i-1,j} + d_{i,j} \quad (i \geq 2),$$

63 then, the approximated joint distribution can be represented as

$$q(\phi_{m,j}, \phi_{m-1,j}, \dots, \phi_{2,j}|\phi_{1,j}) = \prod_{i=2}^m q(\phi_{i,j}|\phi_{i-1,j}) \\ = \prod_{i=2}^m \frac{1}{\sqrt{2\pi\sigma_{i,j}}} \cdot e^{\frac{-1}{2\sigma_{i,j}^2} (\phi_{i,j} - \mu_{i,j})^2}.$$

64 Then, we integrate  $\phi_{2,j}, \phi_{3,j}, \dots, \phi_{m-1,j}$  sequentially to obtain  $q(\phi_{m,j}|\phi_{1,j})$ .

65 First of all, we integrate  $\phi_{2,j}$  for the joint distribution, where we have:

$$\begin{aligned}
& q(\phi_{m,j}, \phi_{m-1,j}, \dots, \phi_{3,j}|\phi_{1,j}) \\
&= \int q(\phi_{m,j}, \phi_{m-1,j}, \dots, \phi_{3,j}, \phi_{2,j}|\phi_{1,j}) d\phi_{2,j} \\
&= \int \prod_{i=2}^m \frac{1}{\sqrt{2\pi}\sigma_{i,j}} \cdot e^{-\frac{1}{2\sigma_{i,j}^2}(\phi_{i,j}-\mu_{i,j})^2} d\phi_{2,j} \\
&= \prod_{i=4}^m \frac{1}{\sqrt{2\pi}\sigma_{i,j}} \cdot e^{-\frac{1}{2\sigma_{i,j}^2}(\phi_{i,j}-\mu_{i,j})^2} \cdot \frac{1}{\sqrt{2\pi}\sigma_{3,j}} \cdot \frac{1}{\sqrt{2\pi}\sigma_{2,j}} \cdot \int e^{-\frac{1}{2\sigma_{3,j}^2}[\phi_{3,j}-(c_{3,j}\phi_{2,j}+d_{3,j})]^2} \\
&\quad e^{-\frac{1}{2\sigma_{2,j}^2}[\phi_{2,j}-(w_2\phi_{1,j}+d_{3,j})]^2} d\phi_{2,j}.
\end{aligned}$$

66 Denote  $\bar{c}_{2,j} = c_{2,j}$  and  $\bar{d}_{2,j} = d_{2,j}$ , and because of  $\sigma_{3,j} = \sigma_{2,j}$ , we have

$$\begin{aligned}
& q(\phi_{m,j}, \phi_{m-1,j}, \dots, \phi_{3,j}|\phi_{1,j}) \\
&= \frac{1}{\sqrt{2\pi}\sigma_{2,j}} \cdot \prod_{i=4}^m \frac{1}{\sqrt{2\pi}\sigma_{i,j}} \cdot e^{-\frac{1}{2\sigma_{i,j}^2}(\phi_{i,j}-\mu_{i,j})^2} \cdot \frac{1}{\sqrt{2\pi}\sigma_{3,j}} \int e^{-\frac{1}{2\sigma_{3,j}^2}[\phi_{3,j}-(c_{3,j}\phi_{2,j}+d_{3,j})]^2} \\
&\quad e^{-\frac{1}{2\sigma_{3,j}^2}[\phi_{2,j}-(\bar{c}_{2,j}\phi_{1,j}+\bar{d}_{2,j})]^2} d\phi_{2,j} \\
&= \frac{1}{\sqrt{2\pi}\sigma_{2,j}} \cdot \prod_{i=4}^m \frac{1}{\sqrt{2\pi}\sigma_{i,j}} \cdot e^{-\frac{1}{2\sigma_{i,j}^2}(\phi_{i,j}-\mu_{i,j})^2} \cdot \frac{1}{\sqrt{2\pi}\sigma_{3,j}} \int \cdot \\
&\quad e^{-\frac{1}{2\sigma_{3,j}^2} \left\{ \underbrace{[\phi_{3,j} - (c_{3,j}\phi_{2,j} + d_{3,j})]^2 + [\phi_{2,j} - (\bar{c}_{2,j}\phi_{1,j} + \bar{d}_{2,j})]^2}_{S_1} \right\}} d\phi_{2,j}.
\end{aligned}$$

67 For  $S_1$ , we have

$$\begin{aligned}
S_1 &= [\phi_{3,j} - (c_{3,j}\phi_{2,j} + d_{3,j})]^2 + [\phi_{2,j} - (\bar{c}_{2,j}\phi_{1,j} + \bar{d}_{2,j})]^2 \\
&= \phi_{3,j}^2 + c_{3,j}^2\phi_{2,j}^2 + d_{3,j}^2 + 2d_{3,j}c_{3,j}\phi_{2,j} - 2c_{3,j}\phi_{3,j}\phi_{2,j} - 2d_{3,j}\phi_{3,j} + \phi_{2,j}^2 \\
&\quad + \bar{c}_{2,j}^2\phi_{1,j}^2 + \bar{d}_{2,j}^2 + 2\bar{d}_{2,j}\bar{c}_{2,j}\phi_{1,j} - 2\bar{c}_{2,j}\phi_{1,j} - 2d_{3,j}\phi_{3,j} \\
&= (c_{3,j}^2 + 1) \cdot \\
&\quad \left[ \phi_{2,j}^2 + \frac{2(d_{3,j}c_{3,j} - c_{3,j}\phi_{3,j} - \bar{c}_{2,j}\phi_{1,j} - \bar{d}_{2,j})}{c_{3,j}^2 + 1} + \left( \frac{d_{3,j}c_{3,j} - c_{3,j}\phi_{3,j} - \bar{c}_{2,j}\phi_{1,j} - \bar{d}_{2,j}}{c_{3,j}^2 + 1} \right)^2 \right] \\
&\quad - \frac{(d_{3,j}c_{3,j} - c_{3,j}\phi_{3,j} - \bar{c}_{2,j}\phi_{1,j} - \bar{d}_{2,j})^2}{c_{3,j}^2 + 1} + \phi_{3,j}^2 + \bar{c}_{2,j}^2\phi_{1,j}^2 + \bar{d}_{2,j}^2 + 2\bar{d}_{2,j}\bar{c}_{2,j}\phi_{1,j} - 2d_{3,j}\phi_{3,j} \\
&= (c_{3,j}^2 + 1) \cdot \left( \phi_{2,j} - \frac{c_{3,j}\phi_{3,j} + \bar{c}_{2,j}\phi_{1,j} + \bar{d}_{2,j} - d_{3,j}c_{3,j}}{c_{3,j}^2 + 1} \right)^2 \\
&\quad + \frac{c_{3,j}^2\phi_{3,j}^2 + \bar{c}_{2,j}^2c_{3,j}^2\phi_{1,j}^2 + \bar{d}_{2,j}^2c_{3,j}^2 + 2d_{3,j}\bar{c}_{2,j}c_{3,j}\phi_{1,j} - \bar{c}_{2,j}^2\phi_{1,j}^2 - \bar{d}_{2,j}^2 - 2\bar{d}_{2,j}\bar{c}_{2,j}\phi_{1,j}}{c_{3,j}^2 + 1} \\
&\quad + \frac{\phi_{3,j}^2 + \bar{c}_{2,j}^2\phi_{1,j}^2 + \bar{d}_{2,j}^2 + 2\bar{d}_{2,j}\bar{c}_{2,j}\phi_{1,j} - 2d_{3,j}\phi_{3,j} - c_{3,j}^2d_{3,j}^2 - c_{3,j}^2\phi_{3,j}^2 + 2c_{3,j}^2d_{3,j}\phi_{3,j}}{c_{3,j}^2 + 1} \\
&\quad + \frac{2c_{3,j}\bar{d}_{2,j}d_{3,j} - 2c_{3,j}\phi_{3,j}\bar{c}_{2,j}\phi_{1,j} - 2c_{3,j}\phi_{3,j}\bar{d}_{2,j} + 2\bar{d}_{2,j}\bar{c}_{2,j}c_{3,j}^2\phi_{1,j} - 2d_{3,j}c_{3,j}^2\phi_{3,j}}{c_{3,j}^2 + 1}.
\end{aligned}$$

68 Then we have

$$\begin{aligned}
S_1 &= (c_{3,j}^2 + 1) \cdot \left( \phi_{2,j} - \frac{c_{3,j}\phi_{3,j} + \bar{c}_{2,j}\phi_{1,j} + \bar{d}_{2,j} - d_{3,j}c_{3,j}}{c_{3,j}^2 + 1} \right)^2 \\
&\quad + \frac{\phi_{3,j}^2 - 2(\bar{c}_{2,j}c_{3,j}\phi_{1,j} + \bar{d}_{2,j}c_{3,j} + d_{3,j})\phi_{3,j}}{c_{3,j}^2 + 1} \\
&\quad + \frac{(\bar{c}_{2,j}c_{3,j}\phi_{1,j} + \bar{d}_{2,j}c_{3,j} + d_{3,j})^2 - d_{3,j}^2 \cdot (c_{3,j}^2 + 1)}{c_{3,j}^2 + 1} \\
&= (c_{3,j}^2 + 1) \cdot \left( \phi_{2,j} - \frac{c_{3,j}\phi_{3,j} + \bar{c}_{2,j}\phi_{1,j} + \bar{d}_{2,j} - d_{3,j}c_{3,j}}{c_{3,j}^2 + 1} \right)^2 \\
&\quad + \frac{1}{c_{3,j}^2 + 1} \cdot [\phi_{3,j} - (c_{3,j}\bar{c}_{2,j}\phi_{1,j} + \bar{d}_{2,j}c_{3,j} + d_{3,j})]^2 - d_{3,j}^2.
\end{aligned}$$

69 Therefore we have

$$\begin{aligned}
&q(\phi_{m,j}, \phi_{m-1,j}, \dots, \phi_{3,j} | \phi_{1,j}) \\
&= \frac{1}{\sqrt{2\pi}\sigma_{2,j}} \cdot \prod_{i=4}^m \frac{1}{\sqrt{2\pi}\sigma_{i,j}} \cdot e^{-\frac{1}{2\sigma_{i,j}^2}(\phi_{i,j} - \mu_{i,j})^2} \\
&\quad \left[ \underbrace{\int \frac{\sqrt{c_{3,j}^2 + 1}}{\sqrt{2\pi}\sigma_{3,j}} e^{-\frac{(c_{3,j}^2 + 1)}{2\sigma_{3,j}^2} \left( \phi_{2,j} - \frac{c_{3,j}\phi_{3,j} + \bar{c}_{2,j}\phi_{1,j} + \bar{d}_{2,j} - d_{3,j}c_{3,j}}{c_{3,j}^2 + 1} \right)^2} d\phi_{2,j}}_{S_2} \right] \\
&\quad e^{\frac{-1}{2\sigma_{3,j}^2 \cdot (c_{3,j}^2 + 1)} [\phi_{3,j} - (c_{3,j}\bar{c}_{2,j}\phi_{1,j} + \bar{d}_{2,j}c_{3,j} + d_{3,j})]^2} \cdot e^{\frac{d_{3,j}^2}{2\sigma_{3,j}^2}} \cdot \frac{1}{\sqrt{c_{3,j}^2 + 1}}
\end{aligned}$$

70 The  $S_2$  part is the integration form of  $\phi_{2,j} \sim \mathcal{N}\left(\frac{c_{3,j}\phi_{3,j} + \bar{c}_{2,j}\phi_{1,j} + \bar{d}_{2,j} - d_{3,j}c_{3,j}}{c_{3,j}^2 + 1}, \frac{\sigma_{3,j}^2}{c_{3,j}^2 + 1}\right)$ , which is  
71 equal to 1, so we have

$$\begin{aligned}
&q(\phi_{m,j}, \phi_{m-1,j}, \dots, \phi_{3,j} | \phi_{1,j}) \\
&= \frac{1}{\sqrt{2\pi}\sigma_{2,j}} \cdot \prod_{i=4}^m \frac{1}{\sqrt{2\pi}\sigma_{i,j}} \cdot e^{-\frac{1}{2\sigma_{i,j}^2}(\phi_{i,j} - \mu_{i,j})^2} \cdot e^{\frac{-1}{2\sigma_{3,j}^2 \cdot (c_{3,j}^2 + 1)} [\phi_{3,j} - (c_{3,j}\bar{c}_{2,j}\phi_{1,j} + \bar{d}_{2,j}c_{3,j} + d_{3,j})]^2} \\
&\quad e^{\frac{d_{3,j}^2}{2\sigma_{3,j}^2}} \cdot \frac{1}{\sqrt{c_{3,j}^2 + 1}}.
\end{aligned}$$

72 We denote  $\bar{c}_{3,j} = c_{3,j}\bar{c}_{2,j}$  and  $\bar{d}_{3,j} = \bar{d}_{2,j}c_{3,j} + d_{3,j}$ , and denote  $r_{2,j} = e^{\frac{d_{3,j}^2}{2\sigma_{3,j}^2}} \cdot \frac{1}{\sqrt{c_{3,j}^2 + 1}}$ , so we have

$$\begin{aligned}
&q(\phi_{m,j}, \phi_{m-1,j}, \dots, \phi_{3,j} | \phi_{1,j}) \\
&= \frac{1}{\sqrt{2\pi}\sigma_{2,j}} \cdot \prod_{i=4}^m \frac{1}{\sqrt{2\pi}\sigma_{i,j}} \cdot e^{-\frac{1}{2\sigma_{i,j}^2}(\phi_{i,j} - \mu_{i,j})^2} \cdot e^{\frac{-1}{2\sigma_{3,j}^2 \cdot (c_{3,j}^2 + 1)} [\phi_{3,j} - (\bar{c}_{3,j}\phi_{1,j} + \bar{d}_{3,j})]^2} \cdot r_{2,j} \\
&= \frac{1}{\sqrt{2\pi}\sigma_{2,j}} \cdot \prod_{i=5}^m \frac{1}{\sqrt{2\pi}\sigma_{i,j}} \cdot e^{-\frac{1}{2\sigma_{i,j}^2}(\phi_{i,j} - \mu_{i,j})^2} \cdot \frac{1}{\sqrt{2\pi}\sigma_4} \cdot e^{\frac{-1}{2\sigma_4^2}[\phi_4 - (c_4\phi_{3,j} + d_4)]^2} \\
&\quad e^{\frac{-1}{2\sigma_{3,j}^2 \cdot (c_{3,j}^2 + 1)} [\phi_{3,j} - (\bar{c}_{3,j}\phi_{1,j} + \bar{d}_{3,j})]^2} \cdot r_{2,j} \\
&= \frac{1}{\sqrt{2\pi}\sigma_{2,j}} \cdot \prod_{i=5}^m \frac{1}{\sqrt{2\pi}\sigma_{i,j}} \cdot e^{-\frac{1}{2\sigma_{i,j}^2}(\phi_{i,j} - \mu_{i,j})^2} \cdot \frac{1}{\sqrt{2\pi}\sigma_4} \cdot e^{\frac{-1}{2\sigma_4^2}[\phi_4 - (c_4\phi_{3,j} + d_4)]^2} \\
&\quad e^{\frac{-1}{2\sigma_4^2}[\phi_{3,j} - (\bar{c}_{3,j}\phi_{1,j} + \bar{d}_{3,j})]^2} \cdot r_{2,j}
\end{aligned}$$

73 Similarly, then we integrate  $\phi_{3,j}$

$$\begin{aligned}
& q(\phi_{m,j}, \phi_{m-1,j}, \dots, \phi_4 | \phi_{1,j}) \\
&= \int q(\phi_{m,j}, \phi_{m-1,j}, \dots, \phi_{3,j} | \phi_{1,j}) d\phi_{3,j} \\
&= \frac{r_{2,j}}{\sqrt{2\pi}\sigma_{2,j}} \cdot \prod_{i=5}^m \frac{1}{\sqrt{2\pi}\sigma_{i,j}} \cdot e^{\frac{-1}{2\sigma_{i,j}^2}(\phi_{i,j} - \mu_{i,j})^2} \cdot \int \frac{1}{\sqrt{2\pi}\sigma_4} \cdot e^{\frac{-1}{2\sigma_4^2}[\phi_4 - (c_4\phi_{3,j} + d_4)]^2} \\
& \quad e^{\frac{-1}{2\sigma_4^2}[\phi_{3,j} - (\bar{c}_{3,j}\phi_{1,j} + \bar{d}_{3,j})]^2} d\phi_{2,j}.
\end{aligned}$$

74 The formulation is similar to the previous one, so we can utilize the process above to integrate  
75 successively, and we finally obtain

$$q(\phi_{m,j} | \phi_{1,j}) = \frac{\prod_{i=2}^{m-1} r_{i,j}}{\sqrt{2\pi}\sigma_{2,j}} \cdot e^{\frac{-1}{2\sigma_{m,j}^2 \cdot (c_{m,j}^2 + 1)}[\phi_{m,j} - (\bar{c}_{m,j}\phi_{1,j} + \bar{d}_{m,j})]^2},$$

76 which indicates

$$p(\phi_{m,j} | \phi_{1,j}, D) \approx \frac{\prod_{i=2}^{m-1} r_{i,j}}{\sqrt{2\pi}\sigma_{2,j}} \cdot e^{\frac{-1}{2\sigma_{m,j}^2 \cdot (c_{m,j}^2 + 1)}[\phi_{m,j} - (\bar{c}_m\phi_{1,j} + \bar{d}_m)]^2},$$

77 where we have the iterative calculation by

$$\begin{aligned}
r_{i,j} &= e^{\frac{d_{i+1,j}^2}{2\sigma_{i+1,j}^2}} \cdot \frac{1}{\sqrt{c_{i+1,j}^2 + 1}}, \\
\bar{c}_{i,j} &= c_{i,j}\bar{c}_{i-1,j}, \quad (i \geq 3), \\
\bar{d}_{i,j} &= \bar{d}_{i-1,j}c_{i,j} + d_{i,j}, \quad (i \geq 3).
\end{aligned}$$

### 78 A.3 Sampling from $p(\phi_m | \phi_M, D)$

79 To sample from the distribution  $p(\phi_m | \phi_M, D)$ , we first obtain a sample  $\phi_1$  from the prior distribution  
80 (i.e.,  $\phi_1 \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ ), then get  $\phi_m$  from a consecutive sampling process:

$$\begin{aligned}
\phi_{M-1} &\sim p(\phi_{M-1} | \phi_M, \phi_1, D), \\
\phi_{M-2} &\sim p(\phi_{M-2} | \phi_{M-1}, \phi_1, D), \\
&\vdots \\
\phi_m &\sim p(\phi_m | \phi_{m+1}, \phi_1, D),
\end{aligned}$$

81 because of the Markov property in our cascaded model. So our target is obtaining the distributions  
82  $p(\phi_{i-1} | \phi_i, \phi_1, D)$ . For a certain  $p(\phi_{i-1} | \phi_i, \phi_1, D)$ , according to the Bayes rule, we have

$$p(\phi_{i-1} | \phi_i, \phi_1, D) = \frac{p(\phi_i | \phi_{i-1}, \phi_1, D) \cdot p(\phi_{i-1} | \phi_1, D)}{p(\phi_i | \phi_1, D)}.$$

83 Similarly with the last section, we use non-bold symbols to represent one dimension of the multi-  
84 dimension parameters, where they are able to transfer independently, and finally construct the eventual  
85 parameters by concatenating, that is,

$$p(\phi_{i-1} | \phi_i, \phi_1, D) = \prod_{j=1}^d p(\phi_{i,j} | \phi_{i-1,j}, \phi_{1,j}, D).$$

86 So our target can be converted to calculate the probability  $p(\phi_{i,j} | \phi_{i-1,j}, \phi_{1,j}, D)$  for all dimensions  
87  $\forall 1 \leq j \leq d$ . According to the Markov property and the transportation probability, we have

$$\begin{aligned}
p(\phi_{i,j} | \phi_{i-1,j}, \phi_{1,j}, D) &\approx q(\phi_{i,j} | \phi_{i-1,j}, \phi_{1,j}) \\
&= \frac{1}{\sqrt{2\pi}\sigma_{i,j}} \cdot e^{\frac{-1}{2\sigma_{i,j}^2}[\phi_{i,j} - (c_i\phi_{i-1,j} + d_{i,j})]^2}.
\end{aligned}$$

88 According to the previous section, we have

$$p(\phi_{i,j}|\phi_{1,j}, D) \approx q(\phi_{i,j}|\phi_{1,j}) = \frac{\prod_{i=2}^{i-1} r_{i,j}}{\sqrt{2\pi\sigma_{2,j}}} \cdot e^{\frac{-1}{2\sigma_{i,j}^2 \cdot (c_{i,j}^2+1)} [\phi_{i,j} - (\bar{c}_i \phi_{1,j} + \bar{d}_i)]^2},$$

$$p(\phi_{i-1,j}|\phi_{1,j}, D) \approx q(\phi_{i-1,j}|\phi_{1,j}) = \frac{\prod_{i=2}^{i-2} r_{i,j}}{\sqrt{2\pi\sigma_{2,j}}} \cdot e^{\frac{-1}{2\sigma_{i-1,j}^2 \cdot (c_{i-1,j}^2+1)} [\phi_{i-1,j} - (\bar{c}_{i-1,j} \phi_{1,j} + \bar{d}_{i-1,j})]^2}.$$

89 Then we have

$$p(\phi_{i-1,j}|\phi_{i,j}, \phi_{1,j}, D) \approx \frac{q(\phi_{i,j}|\phi_{i-1,j}, \phi_{1,j}) \cdot q(\phi_{i-1,j}|\phi_{1,j})}{q(\phi_{i,j}|\phi_{1,j})}$$

$$= \frac{1}{\sqrt{2\pi\sigma_{i,j}} \cdot r_{i-1}} \cdot e^{\frac{-1}{2\sigma_{i,j}^2} [\phi_{i,j} - (c_{i,j} \phi_{i-1,j} + d_{i,j})]^2} \cdot \frac{e^{\frac{[\phi_{i-1,j} - (\bar{c}_{i-1,j} \phi_{1,j} + \bar{d}_{i-1,j})]^2}{-2\sigma_{i-1,j}^2 \cdot (c_{i-1,j}^2+1)}}}{e^{\frac{[\phi_{i,j} - (c_{i,j} \phi_{i-1,j} + d_{i,j})]^2}{-2\sigma_{i,j}^2 \cdot (c_{i,j}^2+1)}}}$$

$$= \sqrt{\frac{c_{i,j} + 1}{2\pi\sigma_{i,j}^2}} \cdot e^{\frac{2\sigma_{i,j}^2}{d_{i,j}^2}} \cdot e^{\frac{[\phi_{i,j} - (c_{i,j} \phi_{i-1,j} + d_{i,j})]^2}{2\sigma_{i+1,j}^2}} \cdot e^{\frac{[\phi_{i,j} - (c_{i,j} \phi_{i-1,j} + d_{i,j})]^2}{-2\sigma_{i,j}^2}} \cdot e^{\frac{[\phi_{i-1,j} - (\bar{c}_{i-1,j} \phi_{1,j} + \bar{d}_{i-1,j})]^2}{-2\sigma_{i,j}^2}}$$

$$= \sqrt{\frac{c_{i,j} + 1}{2\pi\sigma_{i,j}^2}} \cdot e^{\frac{2\sigma_{i,j}^2}{d_{i,j}^2} + \frac{[\phi_{i,j} - (c_{i,j} \phi_{i-1,j} + d_{i,j})]^2}{2\sigma_{i+1,j}^2}} \cdot e^{\frac{-1}{-2\sigma_{i,j}^2} \cdot \underbrace{\{[\phi_{i,j} - (c_{i,j} \phi_{i-1,j} + d_{i,j})]^2 + [\phi_{i-1,j} - (\bar{c}_{i-1,j} \phi_{1,j} + \bar{d}_{i-1,j})]^2\}}_C}$$

90 Then we calculate the part  $C$  as

$$C = [\phi_{i,j} - (c_{i,j} \phi_{i-1,j} + d_{i,j})]^2 + [\phi_{i-1,j} - (\bar{c}_{i-1,j} \phi_{1,j} + \bar{d}_{i-1,j})]^2$$

$$= \phi_{i,j}^2 + (c_{i,j} \phi_{i-1,j} + d_{i,j})^2 - 2(c_{i,j} \phi_{i-1,j} + d_{i,j}) \phi_{i,j}$$

$$+ \phi_{i-1,j}^2 + (\bar{w}_{i-1} \phi_{1,j} + \bar{d}_{i-1,j})^2 - 2\phi_{i-1,j} (\bar{c}_{i-1,j} \phi_{1,j} + \bar{d}_{i-1,j})$$

$$= \phi_{i,j}^2 + c_{i,j}^2 \phi_{i-1,j}^2 + d_{i,j}^2 + 2d_{i,j} c_{i,j} \phi_{i-1,j} - 2c_{i,j} \phi_{i-1,j} \phi_{i,j} - 2d_{i,j} \phi_{i,j} + \phi_{i-1,j}^2$$

$$+ \bar{c}_{i-1,j}^2 \phi_{1,j}^2 + \bar{d}_{i-1,j}^2 + 2\bar{c}_{i-1,j} \bar{d}_{i-1,j} \phi_{1,j} - 2\bar{c}_{i-1,j} \phi_{1,j} \phi_{i-1,j} - 2\bar{d}_{i-1,j} \phi_{i-1,j}$$

$$= (\bar{w}_{i-1}^2 + 1) \cdot \left[ \phi_{i-1,j} - \frac{c_{i,j} \phi_{i,j} + \bar{c}_{i-1,j} \phi_{1,j} + \bar{d}_{i-1,j} - d_{i,j} c_{i,j}}{c_{i,j}^2 + 1} \right]^2 + B,$$

91 where  $B$  does not include  $\phi_i$ , which indicates

$$p(\phi_{i-1,j}|\phi_{i,j}, \phi_{1,j}, D) \sim \mathcal{N}\left(\frac{c_{i,j} \phi_{i,j} + \bar{c}_{i-1,j} \phi_{1,j} + \bar{d}_{i-1,j} - d_{i,j} w_{i-2}}{c_{i,j}^2 + 1}, \frac{\sigma_{i,j}^2}{c_{i,j}^2 + 1}\right).$$

#### 92 A.4 Calculation of $p(\mathbf{x}|\mathbf{e}, \phi_m)$

93 In this section, we will show how to calculate the graph probability  $p(\mathbf{x}|\mathbf{e}, \phi_m)$ . Remember the graph  
94 parameters  $\phi_m = [\boldsymbol{\theta}_m; \mathbf{S}_m; \mathbf{T}_m]$ , so we have

$$p(\mathbf{x}|\mathbf{e}, \phi_m) = \int_{\mathbf{E}} p(\mathbf{x}|\mathbf{e}, \boldsymbol{\theta}_m, \mathbf{E}) \cdot p(\mathbf{E}|\mathbf{e}, \mathbf{S}_m, \mathbf{T}_m) d\mathbf{E}$$

$$= \mathbb{E}_{\mathbf{E} \sim p(\mathbf{E}|\mathbf{S}_m, \mathbf{T}_m)} [p(\mathbf{x}|\mathbf{e}, \boldsymbol{\theta}_m, \mathbf{E})].$$

95 According to Monte Carlo sampling, we have

$$p(\mathbf{x}|\mathbf{e}, \phi_m) = \frac{1}{K} \cdot \sum_{l=1}^K p(\mathbf{x}|\mathbf{e}, \boldsymbol{\theta}_m, \mathbf{E}_l),$$

96 where  $\mathbf{E}_l[i, j] \sim \text{Bernoulli}(\sigma(\mathbf{S}_m^T[i] \cdot \mathbf{T}_m[j]))$ . In order to conduct intervention process, we change  
97 the  $j$ th column of  $\mathbf{E}_l$  to zeros, and represent it with  $\tilde{\mathbf{E}}_l$ . Moreover, we replace the  $j$ th element of  $\mathbf{x}$



98 with  $v$ , and get the result  $\tilde{x}$ . We change the  $j$ th element of  $\epsilon_m$  with zero, and get the result  $\tilde{\epsilon}_m$ . Then  
 99 according the definition of causal graphs, we have

$$p(\mathbf{x}|\mathbf{e}, \phi_m) = \frac{1}{K} \sum_{l=1}^K \mathcal{N}(\mathbf{x}; \mathbf{f}(\tilde{\mathbf{x}}; \tilde{\mathbf{E}}_l, \gamma_m), \tilde{\epsilon}_m),$$

100 where  $\mathbf{f}$  is the causal function that depends on the parameter  $\gamma_m$ .

## 101 **B Bayesian Optimization for Determining $(j^*, v^*, m^*)$**

102 We intend to find the best tuple for acquisition, that is,

$$(j^*, v^*, m^*) = \arg \max_{(j, v, m)} f(j, v, m).$$

103 We define the best interventional value  $v$  under interventional node  $j$  and fidelity  $m$  as

$$\begin{aligned} v^*(j, m) &= \arg \max_v f(j, v, m) \\ &= \arg \max_v f_{j, m}(v). \end{aligned}$$

104 where  $f_{j, m}(v)$  is rewritten from  $f(j, v, m)$  under given  $j, m$ . Therefore, our task is calculat-  
 105 ing  $v^*(j, m)$  for  $\forall j \in [d], m \in [M]$  with Bayesian optimization [1]. We utilize a Gaus-  
 106 sian Process (GP) [2] to model surrogate function distributions for each  $v^*(j, m)$ . We denote  
 107  $f \sim \mathcal{GP}(\mathbf{0}, \mathcal{K}(v_i, v_j))$ , and  $\mathcal{K}(v_i, v_j)$  is the kernel of GP. We sequentially find  $v_t$  and calculate  
 108  $f_{j, m}(v_t)$  to direct the process. According to GP, the previous  $t$  functions and the  $t + 1$  function are  
 109 multivariate Gaussian distribution,

$$\begin{bmatrix} \mathbf{F}_{1:t} \\ f_{t+1} \end{bmatrix} \sim \mathcal{N} \left( \mathbf{0}, \begin{bmatrix} \mathbf{K}_t & \mathbf{k}_{t+1} \\ \mathbf{k}_{t+1}^T & \mathcal{K}(v_{t+1}, v_{t+1}) \end{bmatrix} \right),$$

110 where we define

$$\begin{aligned} \mathbf{F}_{1:t} &= [f_1, f_2, \dots, f_t], \\ \mathbf{k}_{t+1} &= [\mathcal{K}(v_{t+1}, v_1), \mathcal{K}(v_{t+1}, v_2), \dots, \mathcal{K}(v_{t+1}, v_t)]^T, \end{aligned}$$

111

$$\mathbf{K}_t = \begin{bmatrix} \mathcal{K}(v_1, v_1) & \cdots & \mathcal{K}(v_t, v_1) \\ \vdots & \ddots & \vdots \\ \mathcal{K}(v_t, v_1) & \cdots & \mathcal{K}(v_t, v_t) \end{bmatrix}. \quad (1)$$

112 Given previous  $t$  steps, we have the posterior probability is

$$p(f_{t+1} | \{(v_i, f_{j, m}(v_i))\}_{i=1}^t, v_{t+1}) = \mathcal{N}(\mu_t(v_{t+1}), \sigma_t^2(v_{t+1})),$$

113 with the non-parametric means and variances

$$\mu_t(v_{t+1}) = \mathbf{k}_{t+1}^T (\mathbf{K} + \mathbf{I})^{-1} \mathbf{F}_{1:t}, \quad (2)$$

$$\sigma_t^2(v_{t+1}) = \mathcal{K}(v_{t+1}, v_{t+1}) - \mathbf{k}_{t+1}^T (\mathbf{K} + \mathbf{I})^{-1} \mathbf{k}_{t+1}. \quad (3)$$

114 We acquire the next  $v_{t+1}$  with GP-UCB [3] function

$$\begin{aligned} a_{t+1}(v) &= \mu_t(v) + \beta_{ac} \cdot \sqrt{\sigma_t^2(v)}, \\ v_{t+1} &= \arg \max_v a_{t+1}(v). \end{aligned}$$

115 where  $\beta_{ac}$  is a hyper-parameter. Suppose the maximum of steps is  $T$ , the final output of function  
 116  $v^*(j, m)$  is

$$v^*(j, m) = \arg \max_v \mu_T(v).$$

117 Then we choose the best interventional node  $j$  and fidelity  $m$  by their best values under  $\mathcal{O}(d \cdot M)$

$$\begin{aligned} j^*, m^* &= \arg \max_{j, m} v^*(j, m), \\ v^* &= v^*(j^*, m^*). \end{aligned}$$

## 118 C Detailed Training Process of ELBO

### 119 C.1 Derivation Process of ELBO

120 Because we use the distribution  $q(\phi_m)$  to approximate the distribution  $p(\phi_m)$ , then we intend to  
 121 minimize the distance between these two distributions optimize the parameters of  $q(\phi_m)$ , where we  
 122 utilize KL divergence to measure the distance, that is,

$$\Psi^* = \arg \min_{\Psi} \text{KL}[q(\Phi)||p(\Phi|D)].$$

123 According to the variational inference, we have

$$\begin{aligned} & \text{KL}[q(\Phi)||p(\Phi|D)] \\ &= \int q(\Phi) \log \frac{q(\Phi)}{p(\Phi|D)} d\Phi \\ &= \int q(\Phi) \log q(\Phi) d\Phi - \int q(\Phi) \log p(\Phi|D) d\Phi \\ &= \mathbb{E}_{\Phi \sim q(\Phi)} [\log q(\Phi)] - \int q(\Phi) \log \frac{p(\Phi, D)}{p(D)} d\Phi \\ &= \mathbb{E}_{\Phi \sim q(\Phi)} [\log q(\Phi)] - \int q(\Phi) \log p(\Phi, D) d\Phi + \int q(\Phi) \log p(D) d\Phi \\ &= \mathbb{E}_{\Phi \sim q(\Phi)} [\log q(\Phi)] - \mathbb{E}_{\Phi \sim q(\Phi)} [\log p(\Phi, D)] + \int q(\Phi) \log p(D) d\Phi \\ &= \underbrace{\mathbb{E}_{\Phi \sim q(\Phi)} [\log q(\Phi)] - \mathbb{E}_{\Phi \sim q(\Phi)} [\log p(\Phi, D)]}_{-\text{ELBO}} + \log p(D). \end{aligned}$$

124 Because  $\log p(D)$  is not related to  $\Psi$ , minimizing  $\text{KL}[q(\Phi)||p(\Phi|D)]$  is equivalent to maximizing  
 125 the ELBO part, and we have

$$\begin{aligned} \text{ELBO} &= \mathbb{E}_{\Phi \sim q(\Phi)} [\log p(\Phi, D)] - \mathbb{E}_{\Phi \sim q(\Phi)} [\log q(\Phi)] \\ &= \mathbb{E}_{\Phi \sim q(\Phi)} [\log p(D|\Phi)] + \mathbb{E}_{\Phi \sim q(\Phi)} [\log p(\Phi)] - \mathbb{E}_{\Phi \sim q(\Phi)} [\log q(\Phi)] \\ &= \mathbb{E}_{\Phi \sim q(\Phi)} [\log p(D|\Phi) - \log q(\Phi) + \log p(\Phi)] \end{aligned}$$

126 Above all, we can conclude that

$$\Psi^* = \arg \min_{\Psi} \text{KL}[q(\Phi)||p(\Phi|D)]$$

127 is equivalent to maximize evidence lower bound

$$\begin{aligned} \Psi^* &= \arg \max_{\Psi} \text{ELBO} \\ &= \arg \max_{\Psi} \mathbb{E}_{\Phi \sim q(\Phi)} [\log p(D|\Phi) - \log q(\Phi) + \log p(\Phi)]. \end{aligned}$$

### 128 C.2 Estimation of ELBO

129 We represent the equation of ELBO as

$$\begin{aligned} \text{ELBO} &= \mathbb{E}_{\Phi \sim q(\Phi)} [\log p(D|\Phi) - \log q(\Phi) + \log p(\Phi)] \\ &= \underbrace{\mathbb{E}_{\Phi \sim q(\Phi)} [\log p(D|\Phi)]}_A - \underbrace{\mathbb{E}_{\Phi \sim q(\Phi)} [\log q(\Phi) - \log p(\Phi)]}_B. \end{aligned}$$

130 For the part  $A$ , we have

$$A = \mathbb{E}_{\Phi \sim q(\Phi)} \left[ \log \prod_{i=1}^N p(\mathbf{x}^{(i)} | j^{(i)}, v^{(i)}, m^{(i)}, \Phi) \right],$$

131 where  $N$  is the current number of samples in buffer. Then we have

$$\begin{aligned}
A &= \mathbb{E}_{\Phi \sim q(\Phi)} \left[ \log \prod_{i=1}^N p(\mathbf{x}^{(i)} | j^{(i)}, v^{(i)}, m^{(i)}, \Phi) \right] \\
&= \mathbb{E}_{\Phi \sim q(\Phi)} \left[ \sum_{i=1}^N \log p(\mathbf{x}^{(i)} | j^{(i)}, v^{(i)}, m^{(i)}, \Phi) \right] \\
&= \sum_{i=1}^N \mathbb{E}_{\Phi \sim q(\Phi)} \left[ \log p(\mathbf{x}^{(i)} | j^{(i)}, v^{(i)}, m^{(i)}, \Phi) \right] \\
&= \sum_{i=1}^N \mathbb{E}_{\phi_{m^{(i)}} \sim q(\phi_{m^{(i)}})} \left[ \log p(\mathbf{x}^{(i)} | j^{(i)}, v^{(i)}, m^{(i)}, \phi_{m^{(i)}}) \right].
\end{aligned}$$

132 Using Monte Carlo sampling [4], we can calculate the expectation by  $N_S$  samples for each point.

$$A = \sum_{i=1}^N \sum_{j=1}^{N_S} \log p(\mathbf{x}^{(i)} | j^{(i)}, v^{(i)}, m^{(i)}, \phi_{m^{(i)}}^{(j)}),$$

133 where we sample  $\phi_{m^{(i)}}^{(j)} \sim q(\phi_{m^{(i)}})$  with size  $N_S$ .

134 Then we denote the distribution  $q(\Phi) = \mathcal{N}(\tilde{\boldsymbol{\mu}}_{all}, \tilde{\boldsymbol{\Sigma}}_{all})$ , and similarly, we have  $p(\Phi) =$   
135  $\prod_{m=1}^M e^{-\beta \cdot f(\mathbf{S}_m, \mathbf{T}_m)} \cdot \mathcal{N}(\boldsymbol{\mu}_{all}, \boldsymbol{\Sigma}_{all})$ . Both the parameter  $\tilde{\boldsymbol{\mu}}_{all}, \tilde{\boldsymbol{\Sigma}}_{all}$  can be represented by the  
136 parameters in  $\Psi$ , while  $\boldsymbol{\mu}_{all}$  and  $\boldsymbol{\Sigma}_{all}$  are constant. Then we calculate part  $B$

$$\begin{aligned}
B &= \mathbb{E}_{\Phi \sim q(\Phi)} [\log q(\Phi) - \log p(\Phi)] \\
&= \int_{\Phi} \mathcal{N}(\tilde{\boldsymbol{\mu}}_{all}, \tilde{\boldsymbol{\Sigma}}_{all}) \log \frac{\mathcal{N}(\tilde{\boldsymbol{\mu}}_{all}, \tilde{\boldsymbol{\Sigma}}_{all})}{\prod_{m=1}^M e^{-\beta \cdot f(\mathbf{S}_m, \mathbf{T}_m)} \cdot \mathcal{N}(\boldsymbol{\mu}_{all}, \boldsymbol{\Sigma}_{all})} d\Phi \\
&= \int_{\Phi} \mathcal{N}(\tilde{\boldsymbol{\mu}}_{all}, \tilde{\boldsymbol{\Sigma}}_{all}) \log \frac{\mathcal{N}(\tilde{\boldsymbol{\mu}}_{all}, \tilde{\boldsymbol{\Sigma}}_{all})}{\mathcal{N}(\boldsymbol{\mu}_{all}, \boldsymbol{\Sigma}_{all})} d\Phi + \int_{\Phi} \mathcal{N}(\tilde{\boldsymbol{\mu}}_{all}, \tilde{\boldsymbol{\Sigma}}_{all}) \log \frac{1}{\prod_{m=1}^M e^{-\beta \cdot f(\mathbf{S}_m, \mathbf{T}_m)}} d\Phi \\
&= \underbrace{\text{KL}[\mathcal{N}(\tilde{\boldsymbol{\mu}}_{all}, \tilde{\boldsymbol{\Sigma}}_{all}) || \mathcal{N}(\boldsymbol{\mu}_{all}, \boldsymbol{\Sigma}_{all})]}_C + \underbrace{\int_{\Phi} \mathcal{N}(\tilde{\boldsymbol{\mu}}_{all}, \tilde{\boldsymbol{\Sigma}}_{all}) \log \prod_{m=1}^M e^{\beta \cdot f(\mathbf{S}_m, \mathbf{T}_m)} d\Phi}_D.
\end{aligned}$$

137 According to KL divergence of Gaussian distribution, we can calculate  $C$  in a close-form.

$$\begin{aligned}
C &= \text{KL}[\mathcal{N}(\tilde{\boldsymbol{\mu}}_{all}, \tilde{\boldsymbol{\Sigma}}_{all}) || \mathcal{N}(\boldsymbol{\mu}_{all}, \boldsymbol{\Sigma}_{all})] \\
&= \frac{1}{2} \left[ \log \frac{|\boldsymbol{\Sigma}_{all}|}{|\tilde{\boldsymbol{\Sigma}}_{all}|} - d + \text{tr}(\boldsymbol{\Sigma}_{all}^{-1} \tilde{\boldsymbol{\Sigma}}_{all}) + (\tilde{\boldsymbol{\mu}}_{all} - \boldsymbol{\mu}_{all})^T \boldsymbol{\Sigma}_{all}^{-1} (\tilde{\boldsymbol{\mu}}_{all} - \boldsymbol{\mu}_{all}) \right].
\end{aligned}$$

138 Then we calculate  $D$  by the following steps:

$$\begin{aligned}
D &= \int_{\Phi} \mathcal{N}(\boldsymbol{\mu}_{all}, \boldsymbol{\Sigma}_{all}) \log \prod_{m=1}^M e^{\beta \cdot f(\mathbf{S}_m, \mathbf{T}_m)} d\Phi \\
&= \int_{\Phi} \mathcal{N}(\boldsymbol{\mu}_{all}, \boldsymbol{\Sigma}_{all}) \sum_{m=1}^M \log e^{\beta \cdot f(\mathbf{S}_m, \mathbf{T}_m)} d\Phi \\
&= \int_{\Phi} \mathcal{N}(\boldsymbol{\mu}_{all}, \boldsymbol{\Sigma}_{all}) \sum_{m=1}^M \beta \cdot f(\mathbf{S}_m, \mathbf{T}_m) d\Phi \\
&= \beta \cdot \mathbb{E}_{\Phi \sim \mathcal{N}(\boldsymbol{\mu}_{all}, \boldsymbol{\Sigma}_{all})} \left[ \sum_{m=1}^M f(\mathbf{S}_m, \mathbf{T}_m) \right].
\end{aligned}$$

139 Using Monte Carlo sampling, we can calculate the expectation by  $N_D$  samples for each point.

$$\begin{aligned}
D &= \beta \cdot \sum_{i=1}^{N_D} \sum_{m=1}^M f(\mathbf{S}_m^{(i)}, \mathbf{T}_m^{(i)}). \\
&= \beta \cdot \sum_{i=1}^{N_D} \sum_{m=1}^M \mathbb{E}_{p(\mathbf{E}|\mathbf{S}_m^{(i)}, \mathbf{T}_m^{(i)})} [\lambda_1 \cdot [\text{tr}(e^{\mathbf{E}}) - d] + \lambda_2 \cdot \|\mathbf{E}\|],
\end{aligned}$$

140 where we samples  $\Phi^{(i)} \sim \mathcal{N}(\boldsymbol{\mu}_{all}, \boldsymbol{\Sigma}_{all})$  with size  $N_D$ . Using Monte Carlo sampling again, we can  
141 calculate the expectation by  $N_E$  samples.

$$D = \beta \cdot \sum_{i=1}^{N_D} \sum_{m=1}^M \sum_{j=1}^{N_E} [\lambda_1 \cdot [\text{tr}(e^{\mathbf{E}}) - d] + \lambda_2 \cdot \|\mathbf{E}\|],$$

142 where we samples  $\mathbf{E}^{(j)} \sim p(\mathbf{E}|\mathbf{S}_m^{(i)}, \mathbf{T}_m^{(i)})$  with size  $N_E$ .

143 Finally, we obtain the estimation

$$\begin{aligned}
\text{ELBO} &= \sum_{i=1}^N \sum_{j=1}^{N_S} \log p(x^{(i)} | j^{(i)}, v^{(i)}, m^{(i)}, \phi_{m^{(i)}}^{(j)}) \\
&\quad - \frac{1}{2} \left[ \log \frac{\|\boldsymbol{\Sigma}_{all}\|}{\|\tilde{\boldsymbol{\Sigma}}_{all}\|} - d + \text{tr}(\boldsymbol{\Sigma}_{all}^{-1} \tilde{\boldsymbol{\Sigma}}_{all}) + (\tilde{\boldsymbol{\mu}}_{all} - \boldsymbol{\mu}_{all})^T \boldsymbol{\Sigma}_{all}^{-1} (\tilde{\boldsymbol{\mu}}_{all} - \boldsymbol{\mu}_{all}) \right] \\
&\quad - \beta \cdot \sum_{i=1}^{N_D} \sum_{m=1}^M \sum_{j=1}^{N_E} [\lambda_1 \cdot [\text{tr}(e^{\mathbf{E}}) - d] + \lambda_2 \cdot \|\mathbf{E}\|].
\end{aligned}$$

### 144 C.3 Gaussian Reparameterization Trick

145 In the last section, we derive the objection function for optimizing the model parameters, where we  
146 can use methods of the gradient decent to solve it. However, a significant problem rises due to the  
147 sampling process, because the gradient of model parameters can not pass backward from the naive  
148 sampling process(*i.e.*, untraceable). Therefore, we use Gaussian reparameterization trick to make the  
149 Gaussian sampling process traceable.

150 In specific, we will demonstrate the traceable calculation of  $\phi$  by Gaussian reparameterization  
151 trick. In order to sample  $\phi \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ , we first sample  $\delta \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  instead, and then obtain  
152  $\phi = \boldsymbol{\mu} + \delta \odot \boldsymbol{\Sigma}$ . Therefore, the gradient can be traced from  $\phi$  to  $\boldsymbol{\mu}$  and  $\boldsymbol{\Sigma}$ . In specific, both  $\boldsymbol{\mu}$  and  $\boldsymbol{\Sigma}$   
153 can be represented with the function of learnable parameter  $\Psi$ .

### 154 C.4 Gumbel-softmax Reparameterization Trick

155 Besides of the Gaussian sampling process, the Bernoulli sampling in our equation is not traceable  
156 either, so we utilize Gumbel-softmax reparameterization trick to make it traceable.

157 We demonstrate the traceable calculation of  $\mathbf{E} \sim p(\mathbf{E}|\mathbf{S}, \mathbf{T})$  by Gumbel-max reparameterization  
158 trick. According to Gumbel-max [5], we have

$$\text{Bernoulli}(p) \iff \mathbf{1}[G_1 + \log p > G_0 + \log(1 - p)], \quad G_0, G_1 \sim \text{Gumbel}(0, 1).$$

159 Instead of using unit step function, we utilize sigmoid function

$$\sigma(G_1 + \log p > G_0 + \log(1 - p)).$$

160 Therefore, we have

$$\mathbf{E}_{i,j} = \sigma(\mathbf{L}_{i,j} + \mathbf{S}_i^T \cdot \mathbf{T}_j),$$

161 where  $\mathbf{L}_{i,j} \sim L(0, 1)$ . Therefore, we sample  $\mathbf{L}_{i,j} \sim L(0, 1)$  instead, where  $L(0, 1)$  is logistic  
162 distribution, and calculate  $\mathbf{E}_{i,j} = \sigma(\mathbf{L}_{i,j} + \mathbf{S}_i^T \cdot \mathbf{T}_j)$  to trace gradients. Specifically, both  $\mathbf{S}_i$  and  $\mathbf{T}_i$   
163 can be represented with the function of learnable parameter  $\Psi$ .

164 **C.5 Optimization of ELBO**

165 With the estimation and reparameterization trick, we are able to conduct gradient descent methods to  
 166 optimize our parameters with the objection function

$$\Psi^* = \arg \max_{\Psi} \text{ELBO}.$$

167 The format of stochastic gradient descent (SGD) is

$$\Psi \leftarrow \Psi + \gamma \cdot \frac{\partial \text{ELBO}}{\partial \Psi},$$

168 where  $\gamma$  is the learning rate.

169 **D Training Process of Constraint based ELBO**

170 We intend to optimize our parameter with

$$\begin{aligned} \Psi^* &= \arg \max_{\Psi} \mathbb{E}_{\Phi \sim q(\Phi)} [\log p(D|\Phi) - \log q(\Phi) + \log p(\Phi)], \\ \text{s.t. } &\sum_{\{e_s, e_t\}} I(\mathbf{x}_s; \mathbf{x}_t | \phi_M, \{e_s, e_t\}, D) \leq \epsilon. \end{aligned}$$

171 However, the objection has a constraint, which is hard to optimize with gradient descent methods. So  
 172 we utilize Lagrange multiplier [6] to convert it to a constraint-free method:

$$\Psi^* = \arg \max_{\Psi} \mathbb{E}_{\Phi \sim q(\Phi)} [\log p(D|\Phi) - \log q(\Phi) + \log p(\Phi)] + \lambda \cdot \sum_{\{e_s, e_t\}} I(\mathbf{x}_s; \mathbf{x}_t | \phi_M, \{e_s, e_t\}, D),$$

173 where  $\lambda$  is the Lagrange multiplier. Then, we intend to calculate the constraint part.

174 First of all, we have

$$\begin{aligned} &I(\mathbf{x}_s; \mathbf{x}_t | \phi_M, \{e_s, e_t\}, D) \\ &= H(\mathbf{x}_s | \phi_M, \{e_s, e_t\}, D) + H(\mathbf{x}_t | \phi_M, \{e_s, e_t\}, D) - H(\mathbf{x}_s, \mathbf{x}_t | \phi_M, \{e_s, e_t\}, D) \\ &= H(\mathbf{x}_s | \phi_M, e_s, D) + H(\mathbf{x}_t | \phi_M, e_t, D) - H(\mathbf{x}_s, \mathbf{x}_t | \phi_M, \{e_s, e_t\}, D), \end{aligned}$$

175 For the term  $H(\mathbf{x}_s | \phi_M, e_s, D)$ , we have

$$\begin{aligned} H(\mathbf{x}_s | \phi_M, e_s, D) &= - \int p(\mathbf{x}_s | \phi_M, e_s, D) \log p(\mathbf{x}_s | \phi_M, e_s, D) d\mathbf{x}_s \\ &= - \mathbb{E}_{p(\mathbf{x}_s | \phi_M, e_s, D)} [\log p(\mathbf{x}_s | \phi_M, e_s, D)]. \end{aligned}$$

176 We use Monte Carlo sampling to estimate  $H(\mathbf{x}_s | \phi_M, e_s, D)$ , and we have

$$H(\mathbf{x}_s | \phi_M, e_s, D) \approx \frac{1}{K_1 \cdot K_2} \sum_{k_1=1}^{K_1} \sum_{k_2=1}^{K_2} \log p(\mathbf{x}^{(k_1, k_2)} | e^s, \phi_M),$$

177 where we sample graphs  $\phi_m^{k_1} \sim q(\phi_m | e^s, \phi_M)$ , and obtain samples  $\mathbf{x}^{(k_1, k_2)} \sim p(\mathbf{x} | e^s, \phi_m^{k_1})$ .  
 178 Similarly, we can calculate

$$H(\mathbf{x}_t | \phi_M, e_t, D) \approx \frac{1}{K_1 \cdot K_2} \sum_{k_1=1}^{K_1} \sum_{k_2=1}^{K_2} \log p(\mathbf{x}^{(k_1, k_2)} | e^t, \phi_M),$$

179 where we sample graphs  $\phi_m^{k_1} \sim q(\phi_m | e^t, \phi_M)$ , and obtain samples  $\mathbf{x}^{(k_1, k_2)} \sim p(\mathbf{x} | e^t, \phi_m^{k_1})$ .

180 And we have

$$H(\mathbf{x}_s, \mathbf{x}_t | \phi_M, \{e_s, e_t\}, D) \approx \frac{1}{K_1 \cdot K_2 \cdot K_3} \sum_{k_1=1}^{K_1} \sum_{k_2^1=1}^{K_2} \sum_{k_2^2=1}^{K_2} \log p(\mathbf{x}^{(k_1, k_2^1)}, \mathbf{x}^{(k_1, k_2^2)} | \{e_s, e_t\}, \phi_M),$$

181 where we sample graphs  $\phi_m^{k_1} \sim q(\phi_m | \{e_s, e_t\}, \phi_M)$ , obtain samples  $\mathbf{x}^{(k_1, k_2^1)} \sim p(\mathbf{x} | e^s, \phi_m^{k_1})$ , and  
 182 obtain samples  $\mathbf{x}^{(k_1, k_2^2)} \sim p(\mathbf{x} | e^t, \phi_m^{k_1})$ .

183 Therefore, we add constraint on the original loss function to obtained the estimation of constraint  
 184 based ELBO, that is,

$$\begin{aligned} \text{ELBO} = & \sum_{i=1}^N \sum_{j=1}^{N_S} \log p(x^{(i)} | j^{(i)}, v^{(i)}, m^{(i)}, \phi_{m^{(i)}}^{(j)}) \\ & - \frac{1}{2} \left[ \log \frac{\|\Sigma_{all}\|}{\|\tilde{\Sigma}_{all}\|} - d + \text{tr}(\Sigma_{all}^{-1} \tilde{\Sigma}_{all}) + (\tilde{\mu}_{all} - \mu_{all})^T \Sigma_{all}^{-1} (\tilde{\mu}_{all} - \mu_{all}) \right] \\ & - \beta \cdot \sum_{i=1}^{N_D} \sum_{m=1}^M \sum_{j=1}^{N_E} [\lambda_1 \cdot [\text{tr}(e^{\mathbf{E}}) - d] + \lambda_2 \cdot \|\mathbf{E}\|] + \lambda \cdot [I(\mathbf{x}_s; \mathbf{x}_t | \phi_M, \{e_s, e_t\}, D)]. \end{aligned}$$

### 185 E Proof of Theory 3

186 *Proof.* To begin with, we introduce two anchor variables  $\mathbf{x}, e$ , indicating existing samples and  
 187 experiments in the system, which are independent with the following experiments. Since  $\mathbf{x}_s, \mathbf{x}_t$  are  
 188  $\epsilon$ -independent given  $\phi_M, \{e_s, e_t\}$  and  $D$ , we have:

$$\begin{aligned} I(\mathbf{x}_s; \mathbf{x}_t | \phi_M, \{e_s, e_t\}, D) &= I(\mathbf{x}_s; \mathbf{x}_t | \phi_M, \mathbf{x}, e \cup \{e_s, e_t\}, D) \leq \epsilon \\ \Leftrightarrow H(\mathbf{x}_s | \phi_M, \mathbf{x}, e \cup \{e_s, e_t\}, D) &+ H(\mathbf{x}_t | \phi_M, \mathbf{x}, e \cup \{e_s, e_t\}, D) \\ &- H(\mathbf{x}_s, \mathbf{x}_t | \phi_M, \mathbf{x}, e \cup \{e_s, e_t\}, D) \leq \epsilon, \end{aligned}$$

189 Since

$$\begin{aligned} I(\mathbf{x}_s; \phi_M | \mathbf{x}, e \cup \{e_s, e_t\}, D) &= H(\mathbf{x}_s | \mathbf{x}, e \cup \{e_s, e_t\}, D) - H(\mathbf{x}_s | \phi_M, \mathbf{x}, e \cup \{e_s, e_t\}, D) \\ I(\mathbf{x}_t; \phi_M | \mathbf{x}, e \cup \{e_s, e_t\}, D) &= H(\mathbf{x}_t | \mathbf{x}, e \cup \{e_t, e_t\}, D) - H(\mathbf{x}_t | \phi_M, \mathbf{x}, e \cup \{e_s, e_t\}, D) \end{aligned}$$

190 We have:

$$\begin{aligned} & I(\mathbf{x}_s; \phi_M | \mathbf{x}, e \cup \{e_s, e_t\}, D) + I(\mathbf{x}_t; \phi_M | \mathbf{x}, e \cup \{e_s, e_t\}, D) \\ &= H(\mathbf{x}_s | \mathbf{x}, e \cup \{e_s, e_t\}, D) + H(\mathbf{x}_t | \mathbf{x}, e \cup \{e_t, e_t\}, D) \\ &\quad - H(\mathbf{x}_s | \phi_M, \mathbf{x}, e \cup \{e_s, e_t\}, D) - H(\mathbf{x}_t | \phi_M, \mathbf{x}, e \cup \{e_s, e_t\}, D) \\ &\geq H(\mathbf{x}_s, \mathbf{x}_t | \mathbf{x}, e \cup \{e_s, e_t\}, D) - H(\mathbf{x}_s, \mathbf{x}_t | \phi_M, \mathbf{x}, e \cup \{e_s, e_t\}, D) - \epsilon \\ &= I(\mathbf{x}_s, \mathbf{x}_t; \phi_M | \mathbf{x}, e \cup \{e_s, e_t\}, D) - \epsilon. \end{aligned}$$

191 According to the basic mutual information property  $I(A, B; C) - I(B; C) = I(A; C | B)$ , we have:

$$\begin{aligned} & I(\mathbf{x} \cup \mathbf{x}_s; \phi_M | e \cup \{e_s, e_t\}, D) - I(\mathbf{x}; \phi_M | e \cup \{e_s, e_t\}, D) \\ &+ I(\mathbf{x} \cup \mathbf{x}_t; \phi_M | e \cup \{e_s, e_t\}, D) - I(\mathbf{x}; \phi_M | e \cup \{e_s, e_t\}, D) \\ &\geq I(\mathbf{x} \cup \{\mathbf{x}_t, \mathbf{x}_s\}; \phi_M | e \cup \{e_s, e_t\}, D) - I(\mathbf{x}; \phi_M | e \cup \{e_s, e_t\}, D) - \epsilon. \end{aligned}$$

192 Thus, we have:

$$\begin{aligned} & I(\mathbf{x} \cup \mathbf{x}_s; \phi_M | e \cup \{e_s, e_t\}, D) + I(\mathbf{x} \cup \mathbf{x}_t; \phi_M | e \cup \{e_s, e_t\}, D) \\ &\geq I(\mathbf{x} \cup \{\mathbf{x}_t, \mathbf{x}_s\}; \phi_M | e \cup \{e_s, e_t\}, D) + I(\mathbf{x}; \phi_M | e \cup \{e_s, e_t\}, D) - \epsilon. \end{aligned}$$

193 Since different experiments are independent, we have:

$$\begin{aligned} & I(\mathbf{x} \cup \mathbf{x}_s; \phi_M | e \cup \{e_s\}, D) + I(\mathbf{x} \cup \mathbf{x}_t; \phi_M | e \cup \{e_t\}, D) \\ &\geq I(\mathbf{x} \cup \{\mathbf{x}_t, \mathbf{x}_s\}; \phi_M | e \cup \{e_s, e_t\}, D) + I(\mathbf{x}; \phi_M | e, D) - \epsilon. \end{aligned}$$

194 Thus,  $I(\cdot; \phi_M | \cdot, D)$  is  $\epsilon$ -submodular.  $\square$

### 195 F Proof of Theory 4

196 For clear presentation, we denote  $g(\{e_i\}_{i=1}^n) = I(\{\mathbf{x}_i\}_{i=1}^n; \phi_M | \{e_i\}_{i=1}^n, D)$ , then we need to solve  
 197 the following problem:

$$\arg \max_{\{e_i\}_{i=1}^n} g(\{e_i\}_{i=1}^n), \quad (4)$$

198 Suppose  $S^* = \{e_i^*\}_{i=1}^n$  is the optimal solution for objective (4), and the results of the greedy method  
 199 is  $S = \{e_i\}_{i=1}^n$ , where the experiments are sequentially determined from  $e_1$  to  $e_n$ . We denote  
 200  $S_{1:j} = \{e_i\}_{i=1}^j$ , and  $\Delta(e|S_{1:j}) = g(S_{1:j} \cup e) - g(S_{1:j})$ , according to the greedy method, we have:

$$e_{j+1} = \arg \max_e \frac{\Delta(e|S_{1:j})}{\lambda_e},$$

201 where  $\lambda_e$  is the cost of experiment  $e$ .

202 Based on all the above notations, we have:

$$\begin{aligned} g(S^*) &\leq g(S^* \cup S_{1:j}) \\ &= g(S_{1:j}) + g(S_{1:j} \cup e_1^*) - g(S_{1:j}) \\ &\quad + g(S_{1:j} \cup e_1^* \cup e_2^*) - g(S_{1:j} \cup e_1^*) \\ &\quad + \dots \\ &\quad + g(S_{1:j} \cup \{e_1^*, \dots, e_n^*\}) - g(S_{1:j} \cup \{e_1^*, \dots, e_{n-1}^*\}) \\ &= g(S_{1:j}) + \sum_{k=1}^n [g(S_{1:j} \cup \{e_1^*, \dots, e_k^*\}) - g(S_{1:j} \cup \{e_1^*, \dots, e_{k-1}^*\})] \\ &\leq g(S_{1:j}) + \sum_{k=1}^n [g(S_{1:j} \cup \{e_k^*\}) - g(S_{1:j}) + \epsilon] \\ &= g(S_{1:j}) + \sum_{k=1}^n [\Delta(\{e_k^*\}|S_{1:j}) + \epsilon], \end{aligned}$$

203 where the first inequality holds because of the non-decreasing property, and the second inequality  
 204 holds because of the  $\epsilon$ -submodular property.

205 Since  $e_{j+1} = \arg \max_e \frac{\Delta(e|S_{1:j})}{\lambda_e}$ , we have  $\frac{\Delta(e|S_{1:j})}{\lambda_e} \leq \frac{\Delta(e_{j+1}|S_{1:j})}{\lambda_{e_{j+1}}}$  for any  $e$ , thus  $\Delta(e|S_{1:j}) \leq$   
 206  $\frac{\lambda_e}{\lambda_{e_{j+1}}} \Delta(e_{j+1}|S_{1:j}) \leq B_\lambda \Delta(e_{j+1}|S_{1:j})$ . By bringing this result into the above equation, we have:

$$\begin{aligned} g(S^*) &\leq g(S_{1:j}) + \sum_{k=1}^n [\Delta(\{e_k^*\}|S_{1:j}) + \epsilon] \\ &\leq g(S_{1:j}) + \sum_{k=1}^n [B_\lambda \Delta(e_{j+1}|S_{1:j}) + \epsilon] \\ &= g(S_{1:j}) + nB_\lambda \Delta(e_{j+1}|S_{1:j}) + n\epsilon \end{aligned}$$

207 Let  $T_j = g(S^*) - g(S_{1:j})$ , we have:

$$T_j - T_{j+1} = g(S_{1:j+1}) - g(S_{1:j}) = \Delta(e_{j+1}|S_{1:j}) \geq \frac{T_j - n\epsilon}{nB_\lambda}$$

208 Then

$$\begin{aligned} T_n &\leq (1 - \frac{1}{nB_\lambda})T_{n-1} + \frac{\epsilon}{B_\lambda} \leq [(1 - \frac{1}{nB_\lambda})]^2 T_{n-2} + (1 - \frac{1}{nB_\lambda}) \frac{\epsilon}{B_\lambda} + \frac{\epsilon}{B_\lambda} \\ &\leq \dots \leq [(1 - \frac{1}{nB_\lambda})]^n T_0 + [(1 - \frac{1}{nB_\lambda})]^{n-1} \frac{\epsilon}{B_\lambda} + \dots + \frac{\epsilon}{B_\lambda} \end{aligned}$$

209 Let  $B = [(1 - \frac{1}{nB_\lambda})]^{n-1} \frac{\epsilon}{B_\lambda} + \dots + \frac{\epsilon}{B_\lambda} = \frac{\epsilon}{B_\lambda} \sum_{i=1}^n [(1 - \frac{1}{nB_\lambda})]^{i-1}$ , and considering that  $[(1 -$   
 210  $\frac{1}{nB_\lambda})]^n = e^{-\frac{1}{B_\lambda}}$ , we have:

$$g(S^*) - g(S_{1:n}) \leq e^{-\frac{1}{B_\lambda}} g(S^*) + B$$

211 Thus, we have  $g(S_{1:n}) \geq (1 - e^{-\frac{1}{B_\lambda}})g(S^*) - B$ .

## 212 G Algorithm

213 The algorithm for Licence method for single-target intervention scenario is shown in Algorithm 1.  
 214 Moreover, the algorithm for Licence method for multi-target intervention scenario is shown in  
 215 Algorithm 2.

---

**Algorithm 1:** Algorithm of Licence for Single-target Intervention Scenario

---

**Input:** Variable set  $X_V$ , number of oracles  $M$ , cost of oracles  $\Lambda$ , observational data  $D^O$ , total budget  $C$ , and learning rate  $\eta$ .  
**Output:** Causal graph  $\phi_M$ .

- 1 Initialize the model parameter  $\Psi$ .
- 2 Optimize  $\Psi$  with the training process of ELBO under  $D^O$ .
- 3 Initialize  $D^I = \emptyset$ .
- 4 **while** Budget  $C$  does not run out **do**
- 5     Initialize  $j^*, m^*, v^*$  and let  $\zeta^* = -\infty$ .
- 6     **for**  $(j, m)$  in  $\{1, 2, \dots, d\} \times \{1, 2, \dots, M\}$  **do**
- 7         Calculate  $v^*(j, m)$  with BO.
- 8         **if**  $f(j, v^*(j, m), m) > \zeta^*$  **then**
- 9             Update  $j^* \leftarrow j, m^* \leftarrow m$  and  $v^* \leftarrow v^*(j, m)$ .
- 10            Update  $\zeta^* \leftarrow f(j, v^*(j, m), m)$ .
- 11         **end**
- 12     **end**
- 13     Subtract the budget with  $C \leftarrow C - \lambda_{m^*}$ .
- 14     Acquire  $(j^*, v^*, m^*)$  towards the true causal graph to obtain  $\mathbf{x}^* \sim p_m(X_V | do(X_j = v))$ .
- 15     Update  $D^I \leftarrow D^I \cup \{\mathbf{x}^*\}$ .
- 16     Optimize  $\Psi$  with training process of ELBO under  $D^O \cup D^I$ .
- 17 **end**
- 18 Sample  $\phi_M$  from  $p(\phi_M | D)$
- 19 **return** Causal graph  $\phi_M$ .

---

## 216 H More Experiments

### 217 H.1 Experimental Settings

#### 218 H.1.1 Datasets

219 The details of our experimental datasets are presented as follows:

220 • **Erdős-Rényi (ER)** [7] graph is a random graph introduced by Paul Erdős and Alfréd Rényi. For  
221 ER graph, a graph with  $n$  vertices is generated by connecting each pair of vertices with a probability  
222  $p$ .

223 • **Scale-Free (SF)** [8] graph is a type of random graph that has a degree distribution following power  
224 law. A small number of vertices in SF graph own a large number of edges, while the vast majority of  
225 vertices have relatively few edges.

226 • **DREAM** [9] is the abbreviation for Dialogue for Reverse Engineering Assessments and Methods,  
227 which can estimate the reverse quality that causal discovery methods perform. Specifically, we use a  
228 biological graph generator GeneNetWeaver for our experiments, which is a real-word public dataset.

#### 229 H.1.2 Baselines

230 The details of experimental baselines are demonstrated as follows. We utilize DiBS [10] as our basic  
231 graph representation component. For acquisition methods, we use AIT and CBED and obtain the  
232 query tuples of node and value.

233 • **AIT** [11] is an active learning method that utilize f-score to select intervention queries.

234 • **CBED** [12] is based on the calculation of mutual information (MI), which intend to select interven-  
235 tion queries with maximal MI scores after obtaining new samples under current queries.

236 For the multi-target intervention scenario, we extend above methods with greedy strategy, which can  
237 promise an lower bound for approximation with submodular property. For choosing the fidelities to  
238 query, we use two circumstances, *i.e.*, REAL and RANDOM.



---

**Algorithm 2:** Algorithm of Licence for Multi-target Intervention Scenario

---

**Input:** Variable set  $X_V$ , number of oracles  $M$ , cost of oracles  $\Lambda$ , observational data  $D^O$ , total multi-target experiment step  $T$ , total budget  $C$ , and learning rate  $\eta$ .

**Output:** Causal graph  $\phi_M$ .

```
1 Initialize the model parameter  $\Psi$ .
2 Optimize  $\Psi$  with training process of constraint based ELBO under  $D^O$ .
3 Initialize  $B^I = \emptyset$ 
4 for  $t$  in  $1, 2, \dots, T$  do
5   while Budget  $C$  does not run out do
6     Initialize  $j^*, m^*, v^*$  and let  $\zeta^* = -\infty$ .
7     for  $(j, m)$  in  $\{1, 2, \dots, d\} \times \{1, 2, \dots, M\}$  do
8       Calculate  $v^*(j, m)$  with BO.
9       if  $f(j, v^*(j, m), m) > \zeta^*$  then
10        Update  $j^* \leftarrow j, m^* \leftarrow m$  and  $v^* \leftarrow v^*(j, m)$ .
11        Update  $\zeta^* \leftarrow f(j, v^*(j, m), m)$ .
12      end
13    end
14    Subtract the budget with  $C \leftarrow C - \lambda_{m^*}$ .
15    Update  $B^I \leftarrow B^I \cup \{(j^*, v^*, m^*)\}$ .
16  end
17  Acquire  $B^I$  towards the true causal graph to obtain
     $\{\mathbf{x}^* \sim p_m(X_V | do(X_j = v))\}_{(j,v,m) \in B^I}$ .
18  Update  $D^I \leftarrow D^I \cup \{\mathbf{x}^*\}_{(j,v,m) \in B^I}$ .
19  Optimize  $\Psi$  with training process of constraint based ELBO under  $D^O \cup D^I$ .
20 end
21 Sample  $\phi_M$  from  $p(\phi_M | D)$ 
22 return Causal graph  $\phi_M$ .
```

---

239 • **REAL** fidelity means the model always choose the highest fidelity to conduct experiments. This  
240 strategy is aligned with classic causal discovery under active learning paradigm without multi-fidelity  
241 settings, which can just choose the most accurate samples to conduct discovery process.

242 • **RANDOM** fidelity means the model choose different fidelities randomly with uniform probability.

### 243 H.1.3 Metrics

244 The details of experimental metrics are demonstrated as follows. We utilize SHD and AUPRC to  
245 reflect the topological structure discovering performance, and design MSE to reflex the predicting  
246 performance of functional relations.

247 • **SHD** [13] is the abbreviation for Structural Hamming Distance, and it estimate the topological  
248 structure by counting the number of different edges on adjacency matrix. We calculate the expectation  
249 of SHD under multiple graph samplings.

250 • **AUPRC** [14] is the area under precision-recall curve, where we consider entities on the adjacency  
251 matrix as binary classification problem. The AUPRC is also under the expectation for multiple graph  
252 sampling.

253 • **MSE** is designed for estimating the performance of grasping functional relations. We obtain several  
254 samples from the true causal graph, and let our model and the true causal function to conduct forward  
255 process respectively, then calculate the MSE between the two results. We calculate MSE by sampling  
256 graphs for multiple times.

## 257 H.2 Details of Configurations and Computation

258 The details of the configurations of device and platform are demonstrate in Table 1(left). We will  
259 show the details of the time cost on computation. We measure the time cost on the generation of each  
260 intervention per fidelity for all models, and the results are shown in Figure 1(right). We find that our

Table 1: The left table demonstrate the details of the configuration of device and platform. The right table shows the details of time cost on computation.

Name		Details	Model	Time (secs)
CPU	Intel Xeon Platinum 8350C	2.60GHz	AIT-REAL	7.686
GPU	RTX A5000	(24GB)	AIT-RANDOM	7.451
Memory	42GB	RAM	CBED-REAL	7.998
Python	Version 3.8		CBED-RANDOM	7.989
Java	Version 1.8.0	(Necessary for DREAM)	Licence	8.320

Table 2: The details of experimental settings.

Name	Explanation	Value
budget	The total budget for interventional experiments, ( <i>i.e.</i> , $C$ ).	10/20/30/40/50
oracle number	The number of oracles, ( <i>i.e.</i> , $M$ )	3
oracle cost	The cost for each oracle, ( <i>i.e.</i> , $\Lambda$ )	2, 8, 32
oracle noise	The extra additive noise for each oracle.	0.04, 0.02, 0.00
observation number	The number of observational samples.	1000
expect edge number	The number of expect edges.	2
additive noise	The value of additive noise during data generations.	0.01

261 method cost a little more than the baselines, which is probably due to the more complex sampling  
 262 process in our model.

263 We also show the details of experimental settings for our overall experiments in Table 2. We carefully  
 264 tune the hyper-parameters for baselines and our model, and the final values can be obtained in the  
 265 configuration file in our codes.

### 266 H.3 Experiments on DREAM Dataset

267 We conduct experiments on a real-world biological dataset, called DREAM. Note that, DREAM does  
 268 not support the calculation of MSE, because of the lack of interface in this real-world dataset. We use  
 269 two sub-datasets *Ecoli* and *Yeast* as our true causal graphs. The results are shown in Figure 1. We  
 270 find that our model outperforms that other baselines on both *Ecoli* and *Yeast*, and both single-target  
 271 and multi-target intervention scenario.

### 272 H.4 Experiments on More Nodes

273 In this section, we conduct further experiments on datasets with more nodes. We extend the number  
 274 of nodes from 10 to 20, and experiment on the ER graph. The results are shown in Figure 3. We find  
 275 that our model is still effective on the scenario of more nodes, and is better than baselines.

## 276 I Potentially Negative Social Impact

277 Causal discovery focuses on understanding causal relationships between variables. While causal  
 278 discovery has the potential to bring about positive social impacts, it is important to consider both the  
 279 positive and negative implications of its applications. In this response, I will focus on the negative  
 280 impact of causal discovery.

281 • **Reductionism and Oversimplification.** Causal discovery techniques often aim to identify simple  
 282 cause-and-effect relationships. However, complex social phenomena often involve a multitude of  
 283 interconnected factors, making it difficult to capture the full complexity of the system. Relying  
 284 solely on causal discovery may lead to oversimplification and reductionism, neglecting the nuanced  
 285 interactions between variables.

286 • **Ethical Concerns.** Causal discovery can involve analyzing sensitive data, such as personal  
 287 information or medical records. If not handled carefully, the use of this data can raise significant  
 288 ethical concerns related to privacy, consent, and potential discrimination. Improper handling of data  
 289 could lead to violations of privacy and unfair treatment of individuals or groups.

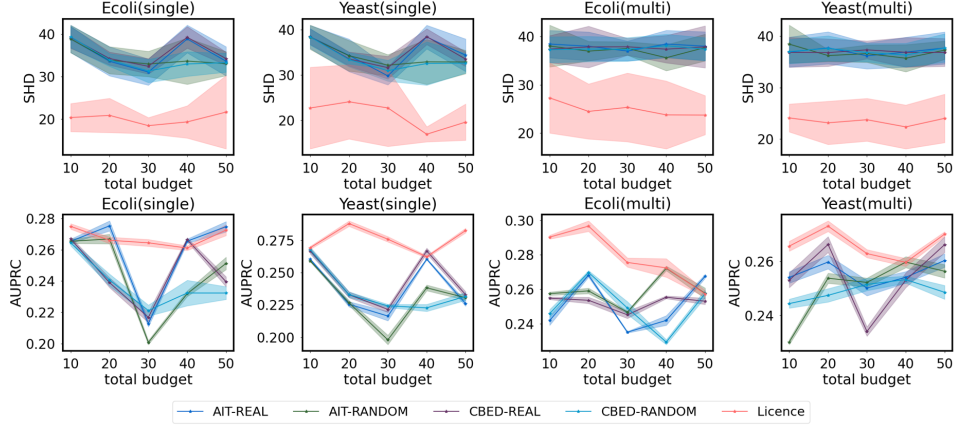


Figure 1: The performance among models on DREAM datasets with different datasets and budgets. Lower SHD, MSE indicate better performances. We conduct each experiment for ten times, and report the average performances and error bars.

Table 3: SHD results of 20 nodes graphs on different budgets. Lower SHD indicates better performances. We conduct each experiment for ten times, and report average performances and error bars.

Model	Budget(10)	Budget(20)	Budget(30)	Budget(40)	Budget(50)
AIT-REAL	63.36±4.89	64.36±5.18	64.53±6.83	63.28±4.86	64.35±5.19
AIT-RANDOM	63.62±4.61	62.16±5.75	64.60±5.23	66.87±6.47	63.53±5.27
DiBS-REAL	63.58±6.35	61.50±7.69	63.50±6.86	63.56±6.34	61.45±7.69
DiBS-RANDOM	63.68±6.77	65.07±6.41	63.91±7.14	63.99±4.46	63.86±3.00
Licence	49.67±11.64	49.61±8.08	55.68±8.63	51.34±11.24	51.36±9.11

290 • **Overreliance on Correlation.** Causal discovery often relies on identifying statistical correlations  
 291 between variables. However, correlation does not imply causation, and there is a risk of mistakenly  
 292 inferring causal relationships based solely on correlation. Overreliance on such methods can lead to  
 293 erroneous conclusions, leading to misguided decision-making and ineffective interventions.

294 • **Social Bias and Inequality.** Causal discovery relies on the data used for analysis, which can reflect  
 295 existing biases and inequalities present in society. If the data used is biased, the causal relationships  
 296 discovered may perpetuate or exacerbate existing social inequalities. Causal discovery methods need  
 297 to be sensitive to potential biases and strive for fairness and inclusivity in both data collection and  
 298 analysis.

299 In conclusion, while causal discovery holds promise in understanding complex systems, it is crucial  
 300 to consider its potential negative impacts. Oversimplification, ethical concerns, overreliance on  
 301 correlation, and social bias are all factors that need to be addressed to ensure responsible and  
 302 beneficial applications of causal discovery. It is essential to approach this field with caution and  
 303 incorporate broader societal considerations to mitigate the negative impacts and harness its potential  
 304 for positive social change.

305 **References**

- 306 [1] Harold J Kushner. A new method of locating the maximum point of an arbitrary multipeak  
307 curve in the presence of noise. 1964.
- 308 [2] Carl Edward Rasmussen. Gaussian processes in machine learning. In *Advanced Lectures on*  
309 *Machine Learning: ML Summer Schools 2003, Canberra, Australia, February 2-14, 2003,*  
310 *Tübingen, Germany, August 4-16, 2003, Revised Lectures*, pages 63–71. Springer, 2004.
- 311 [3] Niranjan Srinivas, Andreas Krause, Sham M Kakade, and Matthias Seeger. Gaussian process  
312 optimization in the bandit setting: No regret and experimental design. *arXiv preprint*  
313 *arXiv:0912.3995*, 2009.
- 314 [4] Reuven Y Rubinstein and Dirk P Kroese. *Simulation and the Monte Carlo method*. John Wiley  
315 & Sons, 2016.
- 316 [5] Chris J Maddison, Andriy Mnih, and Yee Whye Teh. The concrete distribution: A continuous  
317 relaxation of discrete random variables. *arXiv preprint arXiv:1611.00712*, 2016.
- 318 [6] Dimitri P Bertsekas. *Constrained optimization and Lagrange multiplier methods*. Academic  
319 press, 2014.
- 320 [7] Paul Erdős, Alfréd Rényi, et al. On the evolution of random graphs. *Publ. Math. Inst. Hung.*  
321 *Acad. Sci*, 5(1):17–60, 1960.
- 322 [8] Lun Li, David Alderson, John C Doyle, and Walter Willinger. Towards a theory of scale-free  
323 graphs: Definition, properties, and implications. *Internet Mathematics*, 2(4):431–523, 2005.
- 324 [9] Thomas Schaffter, Daniel Marbach, and Dario Floreano. Genenetweaver: in silico benchmark  
325 generation and performance profiling of network inference methods. *Bioinformatics*, 27(16):  
326 2263–2270, 2011.
- 327 [10] Lars Lorch, Jonas Rothfuss, Bernhard Schölkopf, and Andreas Krause. Dibs: Differentiable  
328 bayesian structure learning. *Advances in Neural Information Processing Systems*, 34:24111–  
329 24123, 2021.
- 330 [11] Yashas Annadani, Jonas Rothfuss, Alexandre Lacoste, Nino Scherrer, Anirudh Goyal, Yoshua  
331 Bengio, and Stefan Bauer. Variational causal networks: Approximate bayesian inference over  
332 causal structures. *arXiv preprint arXiv:2106.07635*, 2021.
- 333 [12] Panagiotis Tigas, Yashas Annadani, Andrew Jesson, Bernhard Schölkopf, Yarin Gal, and Stefan  
334 Bauer. Interventions, where and how? experimental design for causal models at scale. In  
335 *Advances in Neural Information Processing Systems*.
- 336 [13] Ioannis Tsamardinos, Laura E Brown, and Constantin F Aliferis. The max-min hill-climbing  
337 bayesian network structure learning algorithm. *Machine learning*, 65:31–78, 2006.
- 338 [14] Jesse Davis and Mark Goadrich. The relationship between precision-recall and roc curves. In  
339 *Proceedings of the 23rd international conference on Machine learning*, pages 233–240, 2006.