

Multi-task Learning of Order-Consistent Causal Graphs

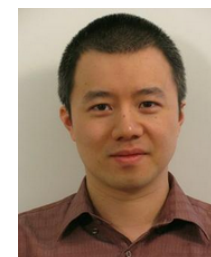
Xinshi Chen¹

Haoran Sun¹

Caleb Ellington²

Eric Xing^{2,3}

Le Song^{3,4}



Georgia
Tech



Carnegie
Mellon
University



MOHAMED BIN ZAYED
UNIVERSITY OF
ARTIFICIAL INTELLIGENCE



BioMap
百图生科

1 *Georgia Institute of Technology*

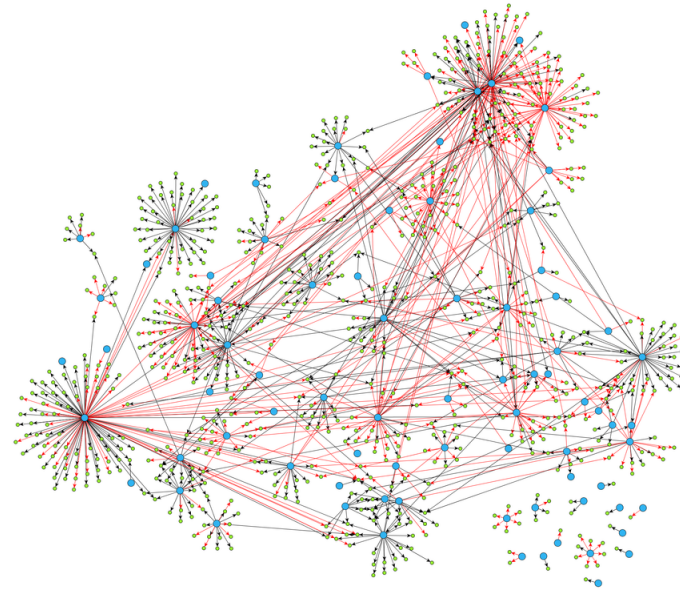
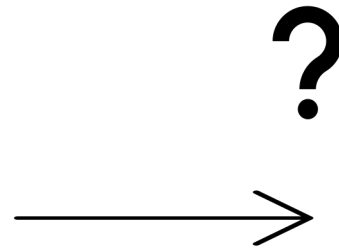
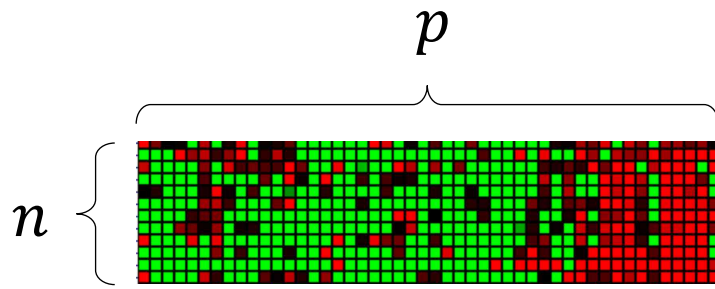
2 *Carnegie Mellon University*

3 *Mohamed bin Zayed University of Artificial Intelligence*

4 *BioMap*

Challenge 1 – Small Data

- Very few samples are observed for reconstructing the graph
- $n \ll p$



Challenge 2 – Non-identifiability

- Even with **infinite samples**, a DAG can be **non-identifiable**.

- Example:

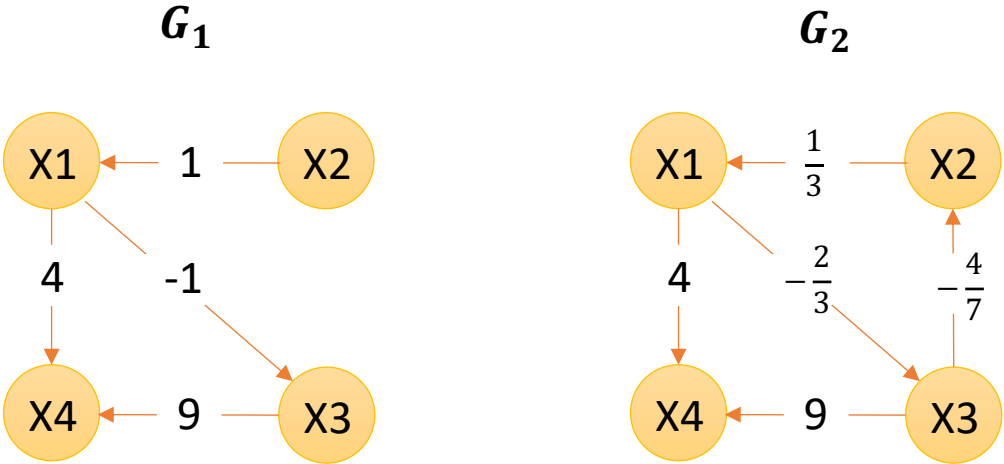
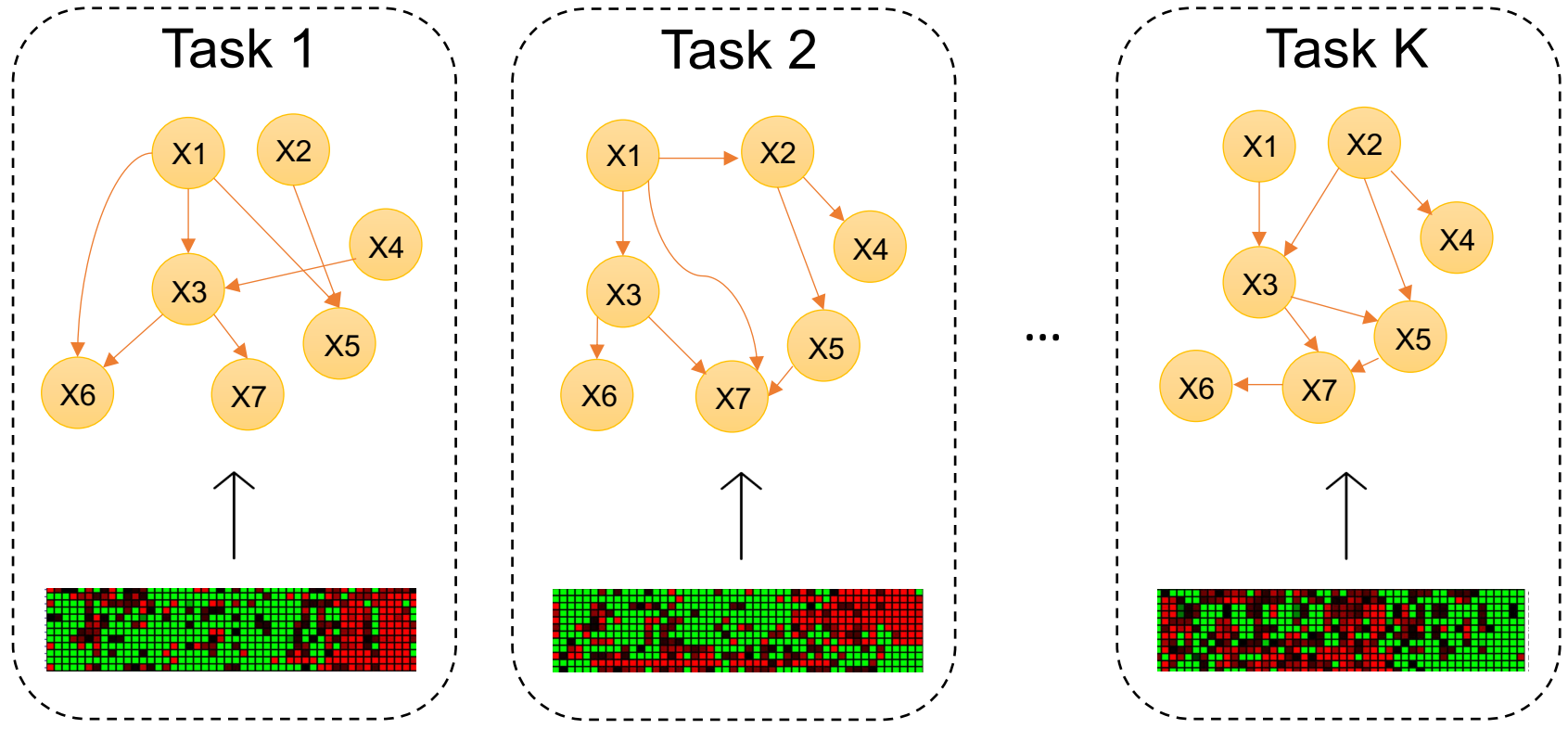


Figure 1 . (From [Aragam, 2015]).

Linear SEM $\begin{cases} X = G_1^T X + W_1 \\ X = G_2^T X + W_2 \end{cases} \implies \begin{cases} X \sim \mathcal{N}(0, \Sigma_1) \\ X \sim \mathcal{N}(0, \Sigma_2) \end{cases} \text{ with } \Sigma_1 = \Sigma_2 \quad \text{!!!}$

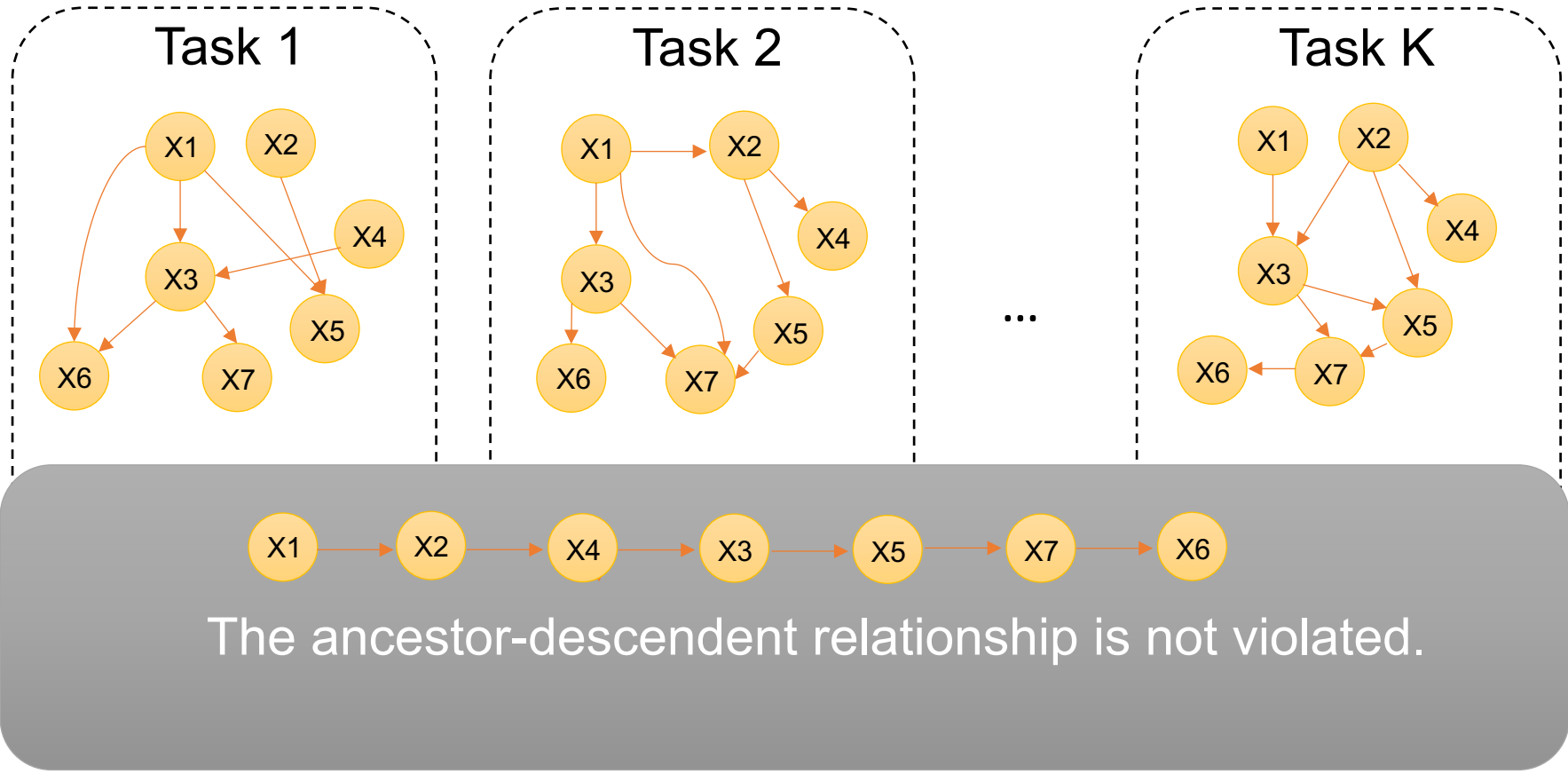
What can we utilize?

Similarity among multiple tasks!



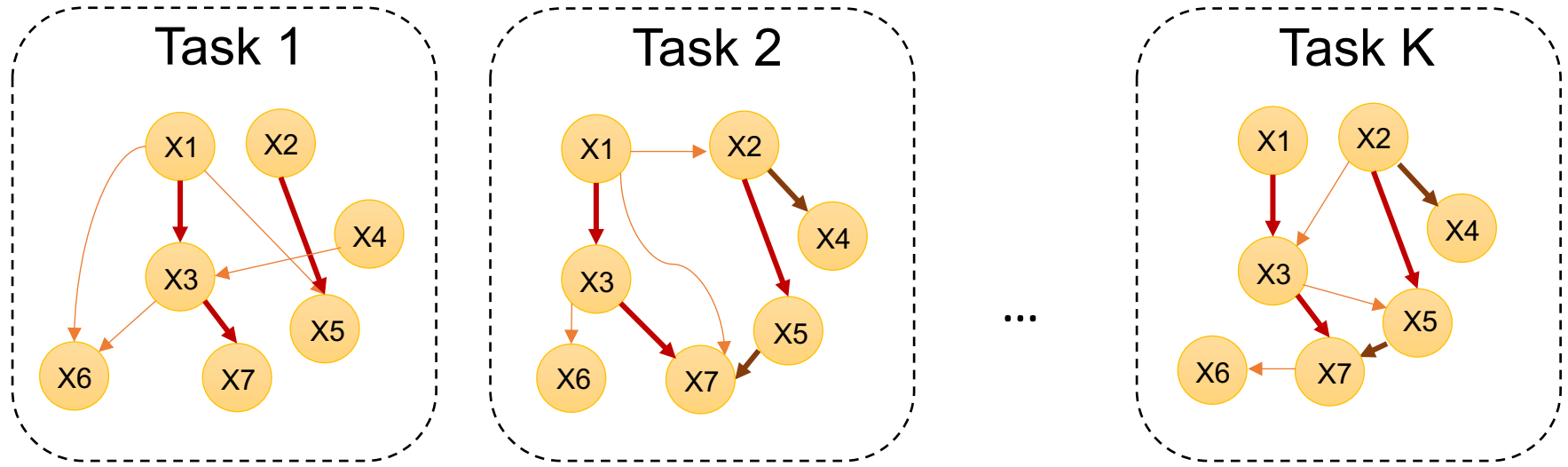
What can we utilize?

- Assumption 1 – Consistent Causal Ordering (Topological Ordering)



What can we utilize?

- Assumption 2 – Sparsity Pattern



Size of the support union of edges $|S| = s$

Multi-task Learning Setting

- K linear SEM (structural equation model)

$$\text{for } k = 1, \dots, K, \quad \mathbf{X}^{(k)} = \underbrace{G_0^{(k)\top}}_{\text{DAG}} \mathbf{X}^{(k)} + \underbrace{W^{(k)}}_{\sim \mathcal{N}(0, \Omega^{(k)})}$$

- Assume each task has n samples.

$$\underbrace{\mathbf{X}^{(1)}, \mathbf{X}^{(2)}, \dots, \mathbf{X}^{(K)}} \quad \text{where} \quad \mathbf{X}^{(k)} \in \mathbb{R}^{n \times p}$$

Jointly recover?

$$G_0^{(1)}, G_0^{(2)}, \dots, G_0^{(K)}$$

Joint Estimator

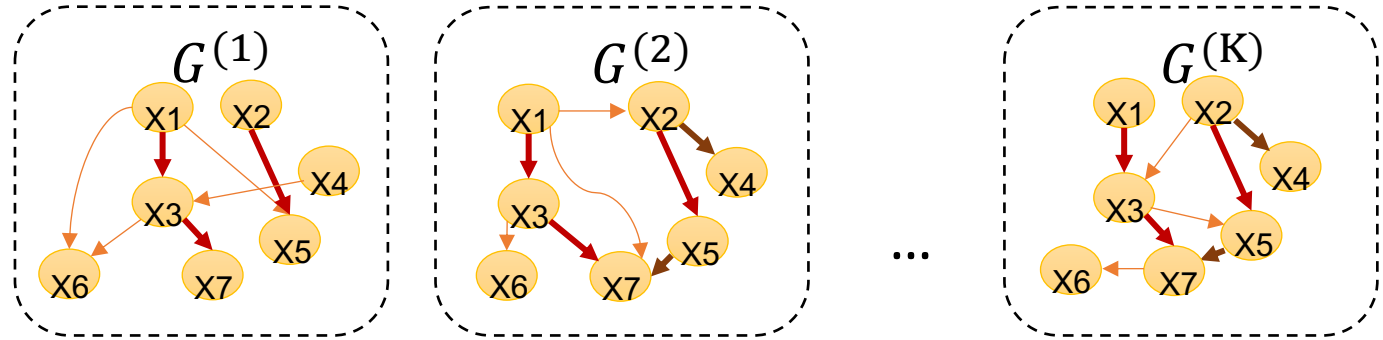
$$\min_{\pi, \{G^{(k)}\}_k^K} \sum_{k=1}^K \frac{1}{2n} \|\mathbf{X}^{(k)} - \mathbf{X}^{(k)} G^{(k)}\|_F^2 + \lambda \|G^{(1:K)}\|_{l_1/l_2}$$

s.t. $\begin{cases} 1. \pi \in \mathcal{S}_p \\ 2. G^{(k)} \in \text{DAG}(\pi) \end{cases}$

1. $\pi \in \mathcal{S}_p$ represents the causal order (i.e., topological order)



2. $G^{(k)} \in \text{DAG}(\pi)$ is a DAG whose topological order is **consistent with π** .



Joint Estimator

$$\min_{\pi, \{G^{(k)}\}_k^K} \sum_{k=1}^K \frac{1}{2n} \|\mathbf{X}^{(k)} - \mathbf{X}^{(k)} G^{(k)}\|_F^2 + \lambda \|G^{(1:K)}\|_{l_1/l_2}$$

s.t. $\begin{cases} 1. \pi \in \mathcal{S}_p \\ 2. G^{(k)} \in \text{DAG}(\pi) \end{cases}$

Different from separate estimation:

- It optimizes a single π *shared* across DAGs.
- The *group norm* $\|G^{(1:K)}\|_{l_1/l_2}$ penalizes the size of union support softly.

Questions

Theoretical questions:

- Is this joint estimator leading to an *improved sample complexity*?
- Can this joint estimator help to recover *non-identifiable DAGs*?

Practical question:

- How to compute the minimizer $\pi, \{G^{(k)}\}_k^K$ *efficiently*?

Main Result – Identifiable Case

- Assume for each k , $G_0^{(k)}$ is a unique minimum-trace DAG
- *(Theorem 3.1) Recovering the true causal order π_0*

A sample complexity measure: $\theta(n, K, p, s) = \frac{p}{s} \sqrt{\frac{nK}{p \log p}}$

the rate at which the sample size must grow

- *(Theorem 3.1) Recovering the DAGs*

Averaged error: $\frac{1}{K} \sum_{k=1}^K \|G^{(k)} - G_0^{(k)}\|_F^2 = \mathcal{O}\left(s \sqrt{\frac{p \log p}{nK}}\right)$

Main Result – Non-identifiable Case

- Assume K' DAGs $\{G_0^{(1)}, \dots, G_0^{(K')}\}$ are *identifiable*.
- The other $K - K'$ DAGs are *non-identifiable*.
- *Recovering the true causal order π_0*

A sample complexity measure: ~~$\theta(n, K, p, s) = \frac{p}{s} \sqrt{\frac{nK}{p \log p}}$~~

$$\theta(n, K, p, s) = \frac{p}{s} \sqrt{\frac{1}{p \log p} \frac{nK'^2}{K}} \quad \checkmark$$

Practical Algorithm

$$\min_{\pi, \{G^{(k)}\}_k^K} \sum_{k=1}^K \frac{1}{2n} \|\mathbf{X}^{(k)} - \mathbf{X}^{(k)} G^{(k)}\|_F^2 + \lambda \|G^{(1:K)}\|_{l_1/l_2}$$

s.t. $\begin{cases} 1. \pi \in \mathcal{S}_p \\ 2. G^{(k)} \in \text{DAG}(\pi) \end{cases}$

- How to compute the optimal solution $\pi, \{G^{(k)}\}_k^K$ *efficiently*?

Practical Algorithm

Key: an equivalent continuous formulation.

$$\min_{\substack{T \in \mathbb{R}^{p \times p} \\ G^{(1)}, \dots, G^{(K)} \in \mathbb{R}^{p \times p}}} \sum_{k=1}^K \frac{1}{2n} \left\| \mathbf{X}^{(k)} - \mathbf{X}^{(k)} \overline{G}^{(k)} \right\|_F^2 + \lambda \|\overline{G}^{(1:K)}\|_{l_1/l_2} + \rho \|\mathbf{1}_{p \times p} - T\|_F^2$$

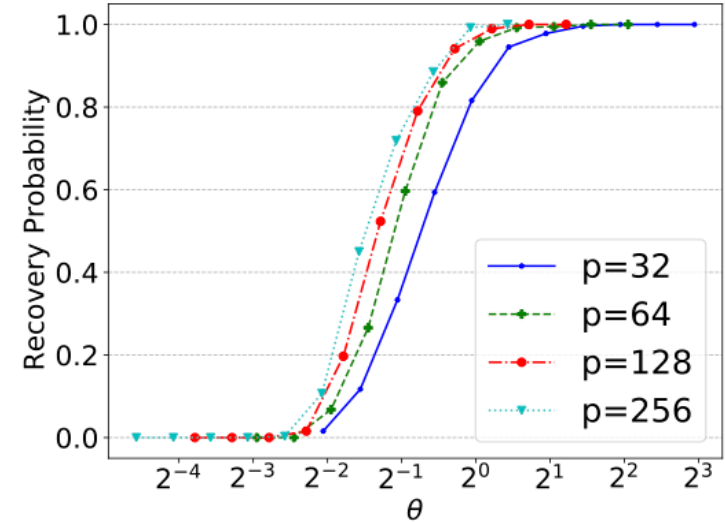
subject to $h(T) := \text{trace}(e^{T \circ T}) - p = 0,$

where $\overline{G}^{(k)} := G^{(k)} \circ T$ is element-wise multiplication between $G^{(k)}$ and T

Synthetic Experiment: Linear SEM



- Order recovery probability versus theoretical sample complexity: $\theta = p/s\sqrt{nK/(p\log p)}$
- Order recovery probability under different problem sizes, number of tasks, and number of samples per task.



$p = 32$

Number of Samples per Task	1	2	4	8	16	32
10	0	0	0	0.14	0.27	0.67
20	0.078	0.17	0.41	0.77	0.92	0.98
40	0.39	0.56	0.81	0.94	0.91	1
80	0.7	0.91	0.95	0.97	1	1
160	0.59	0.8	0.94	1	1	1
320	0.69	0.84	0.95	1	1	1

$p = 64$

Number of Samples per Task	1	2	4	8	16	32
10	0	0	0	0	0.047	0.48
20	0	0	0.062	0.47	0.84	0.95
40	0.062	0.23	0.45	0.91	0.98	1
80	0.31	0.52	0.89	0.97	1	1
160	0.5	0.58	0.88	0.97	0.98	1
320	0.62	0.83	0.95	1	1	1

$p = 128$

Number of Samples per Task	1	2	4	8	16	32
10	0	0	0	0	0	0
20	0	0	0	0.031	0.41	0.61
40	0	0	0.14	0.61	0.84	0.92
80	0	0.078	0.42	0.67	0.94	1
160	0.17	0.3	0.62	0.94	1	1
320	0.25	0.36	0.7	0.91	1	1

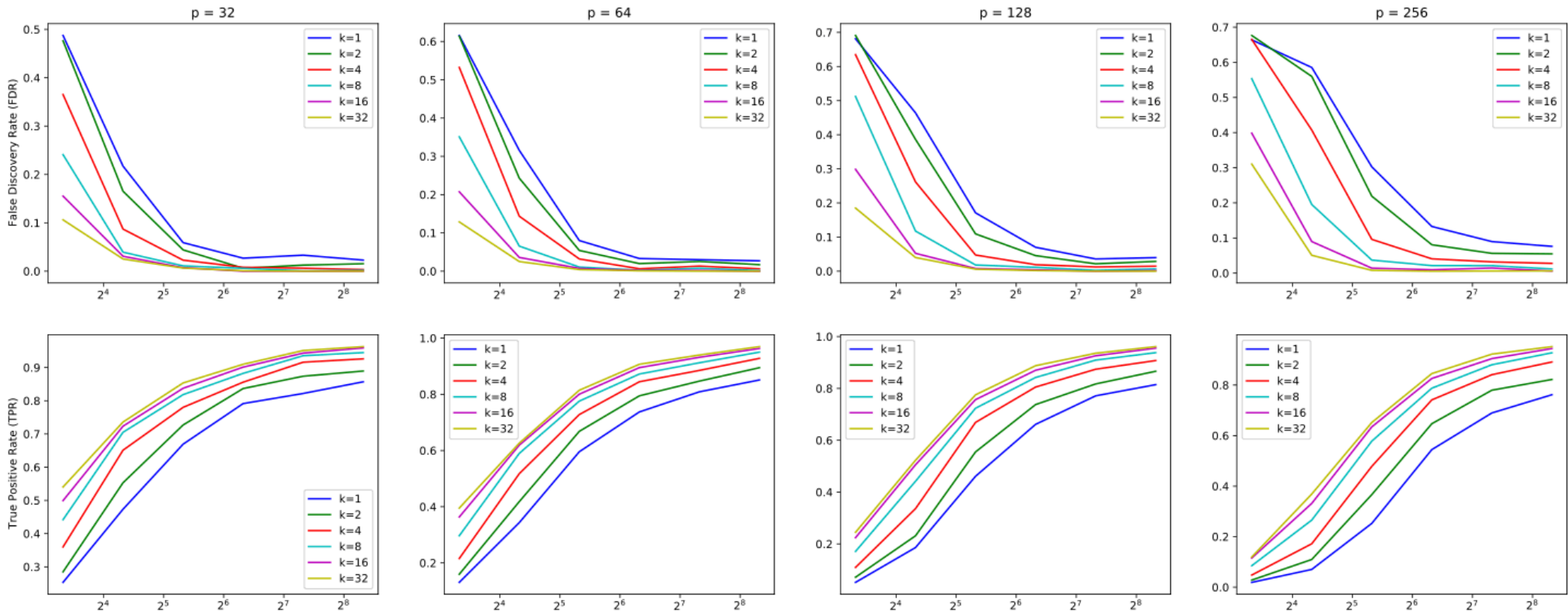
$p = 256$

Number of Samples per Task	1	2	4	8	16	32
10	0	0	0	0	0	0
20	0	0	0	0	0.031	0.27
40	0	0	0	0.094	0.53	0.86
80	0	0	0.016	0.34	0.84	0.92
160	0	0	0.078	0.39	0.7	0.95
320	0.016	0.047	0.19	0.72	0.94	1

Synthetic Experiment: Linear SEM



- Structure recovery quality with different numbers of tasks in False Discovery Rate (FDR), and True Positive Rate (TPR)



Gene Expression Experiment using SERGIO



Rate of gene i expression

$$\frac{dx_i}{dt} = P_i x_i - \beta(x_i)$$

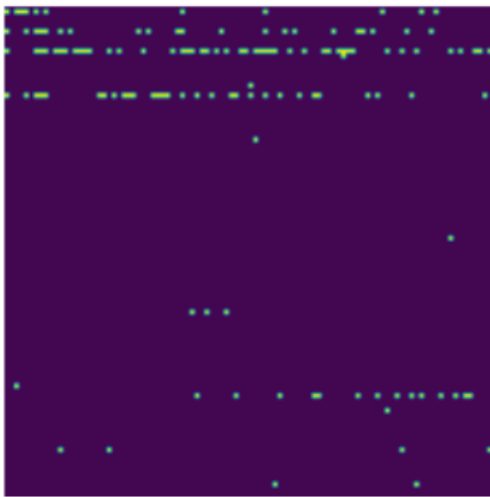
Total production rate of gene i

$$P_i = \sum_{j \in R_i} p_{ij} + b_i$$

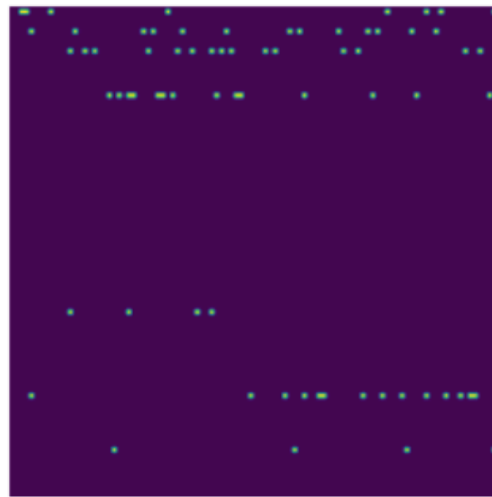
Effect of regulator j on gene i

$$p_{ij} = K_{ij} \frac{x_j^{n_{ij}}}{h_{ij}^{n_{ij}} + x_j^{n_{ij}}}$$

(a) True: Ecoli 100

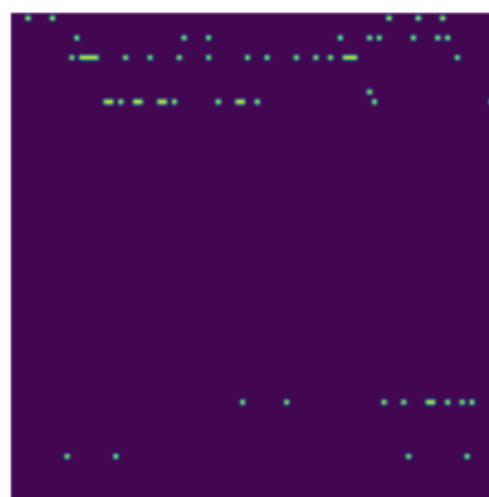


(b) $K = 1, n = 1000$



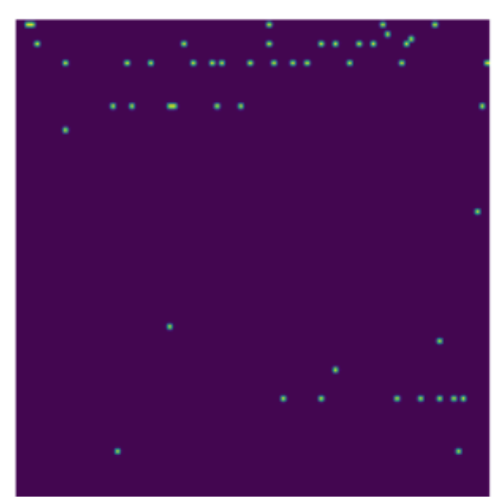
FDR: 0.11, TPR: 0.47

(c) $K = 10, n = 100$



FDR: 0.06, TPR: 0.41

(d) $K = 1, n = 100$



FDR: 0.18, TPR: 0.29