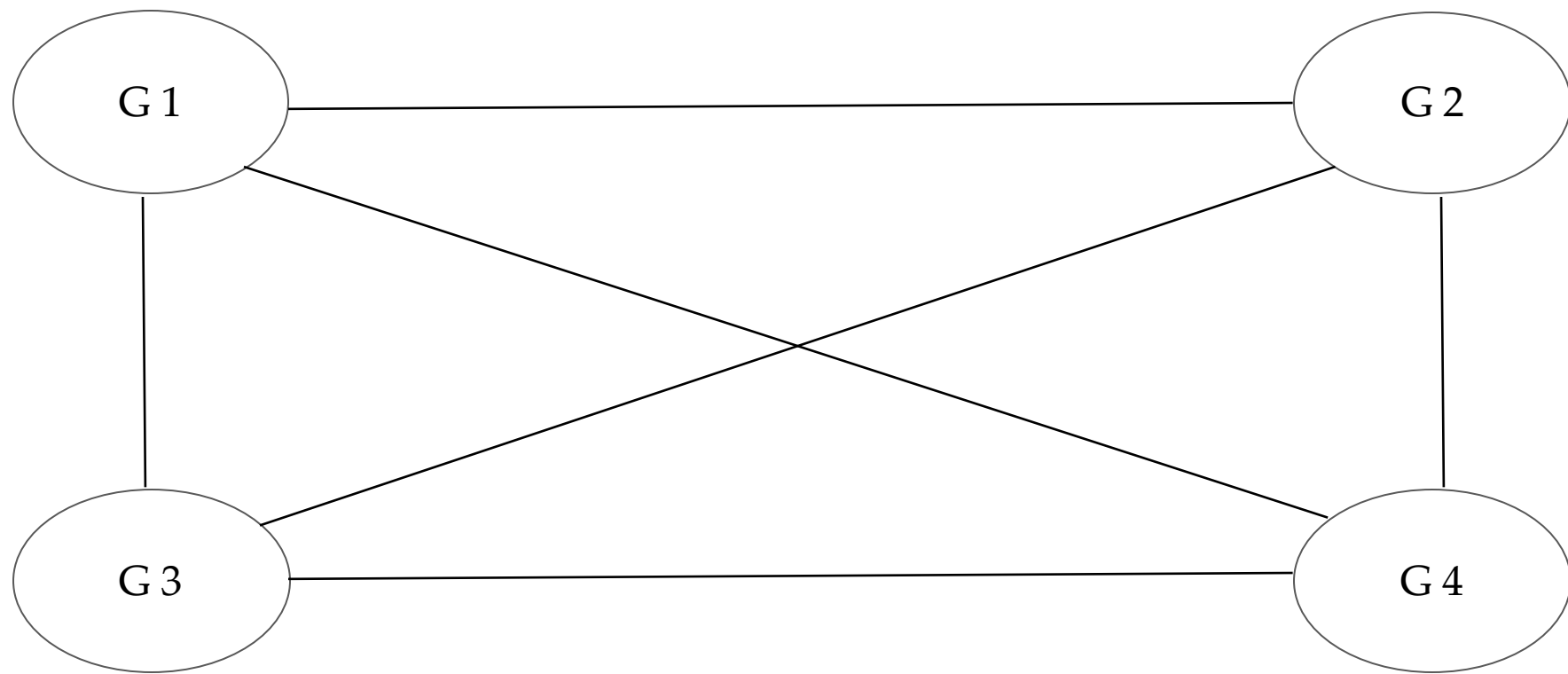# INFERRING GENE REGULATORY NETWORKS USING THE IMPROVED MARKOV BLANKET DISCOVERY ALGORITHM
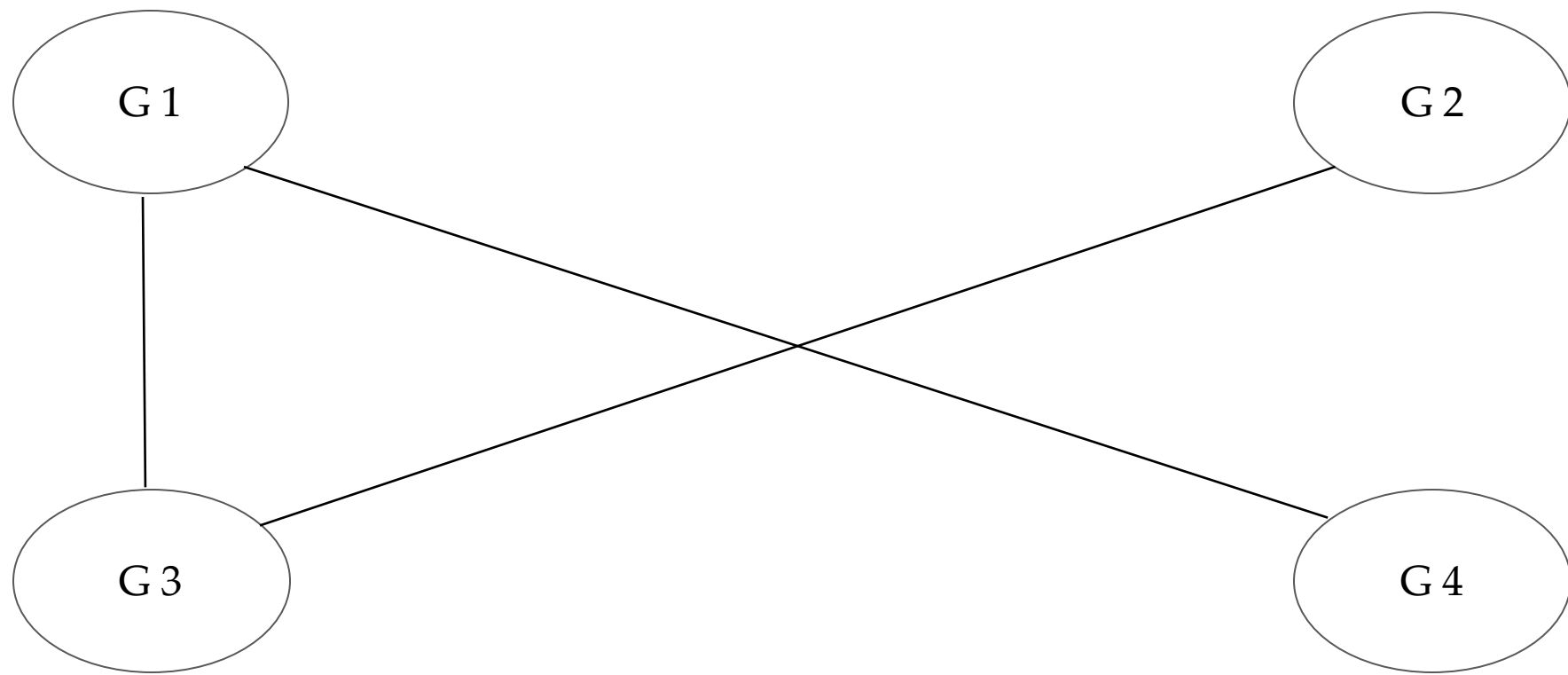
W. Liu, Y. Jiang, L. Peng, X. Sun, W. Gan, Q. Zhao · Interdisciplinary Sciences: Computational Life Sciences· 2022
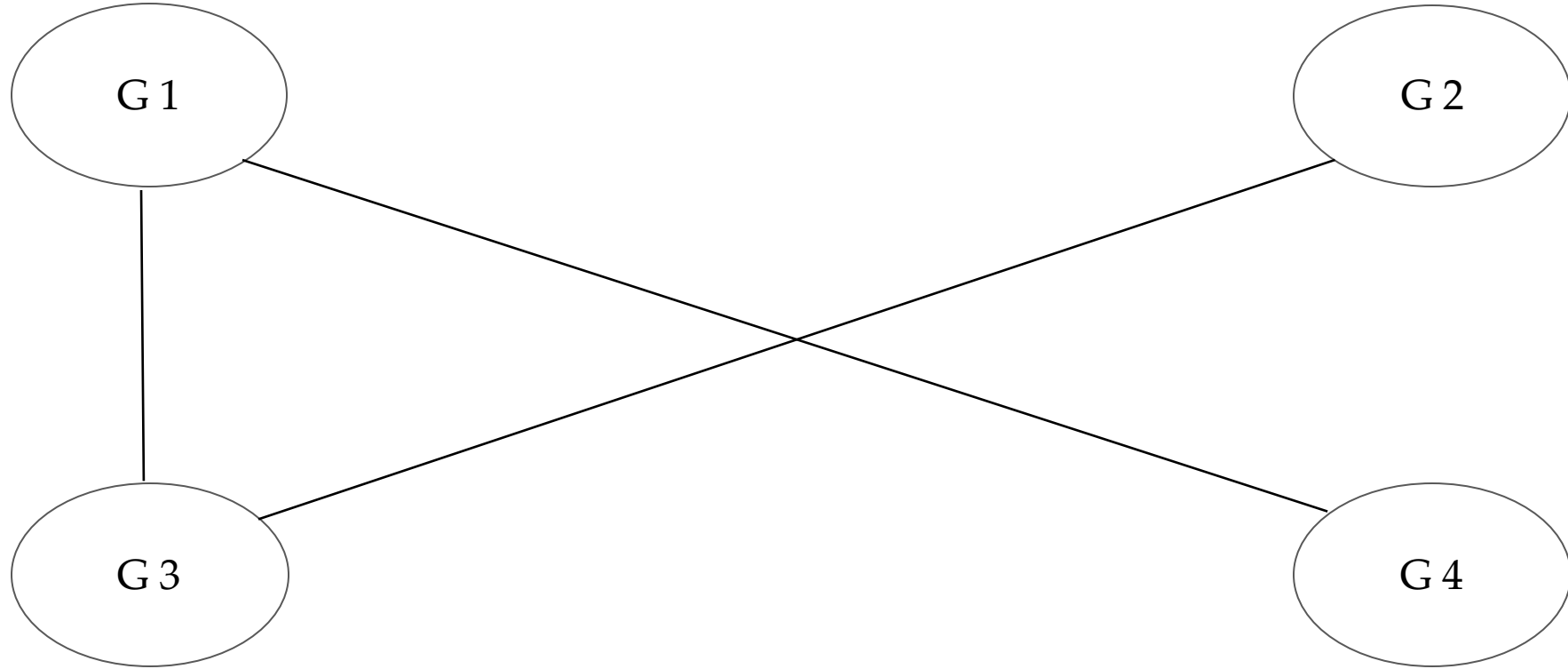
By Soel Micheletti · March 2023
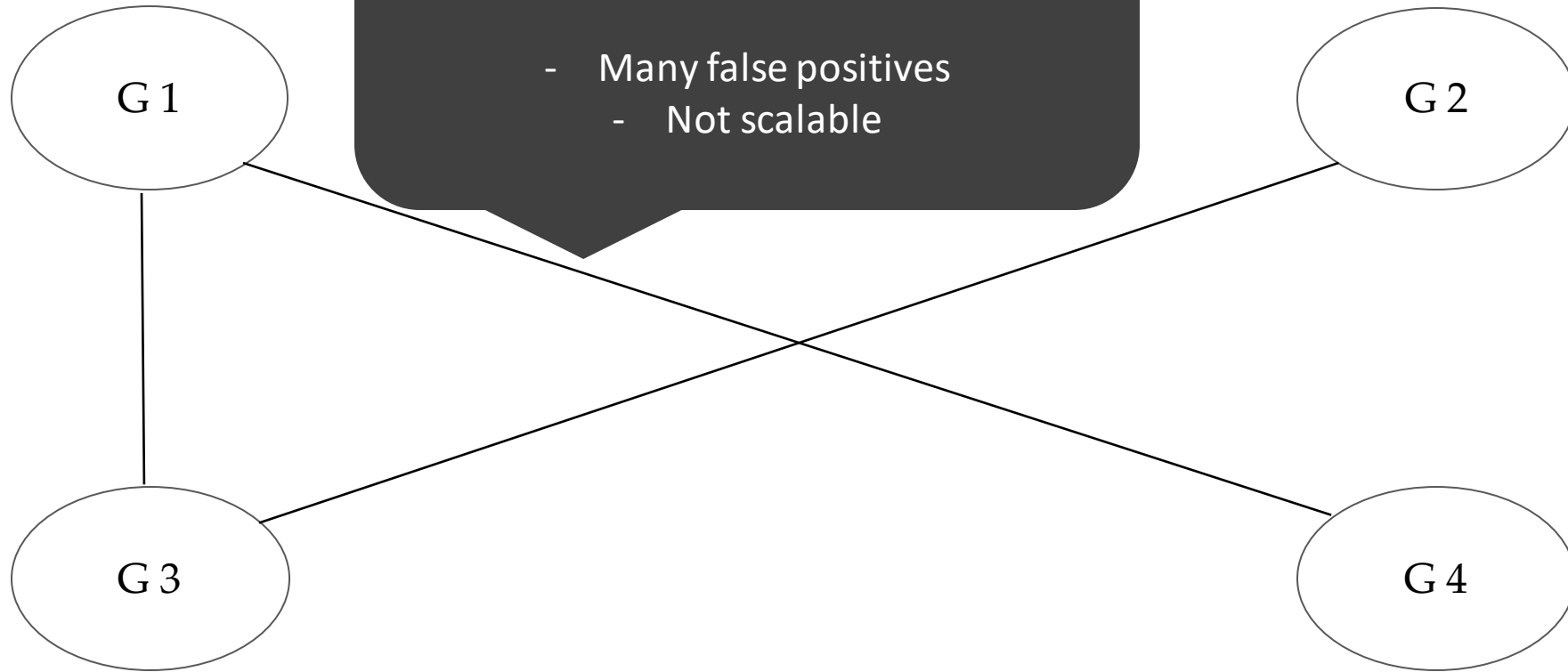
HARVARD **T.H. CHAN**
SCHOOL OF PUBLIC HEALTH

# Undirected, unweighted gene co-expression

# Undirected, unweighted gene co-expression

# IMBDANET

**Step 1: IMBDA** · Step 2: pruning · Step 3: fine tuning

# IMBDANET

**Step 1: IMBDA** · Step 2: pruning · Step 3: fine tuning

For each gene, it uses an information theoretic score to compute its Markov Blanket.

# IMBDANET

For each gene, it uses an information theoretic score to compute its Markov Blanket.

**Definition 5.** *We define a minimal subset of active predictors as a minimal subset* $S^* \subset \{1,\ldots,p\} =: F$, *such that*
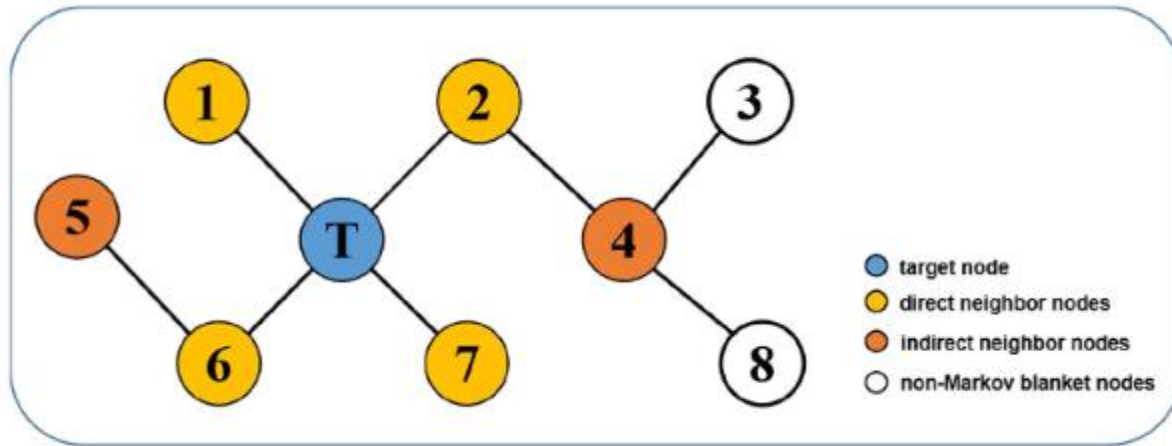
$$I(Y; X_{S^*}) = \max_{T \subseteq F} I(Y; X_T) = I(Y; X_F), \tag{11}$$

# IMBDANET

For each gene, it uses an information theoretic score to compute its Markov Blanket.

# IMBDANET

**Step 1: IMBDA** · Step 2: pruning · Step 3: fine tuning



target node
direct neighbor nodes
indirect neighbor nodes
non-Markov blanket nodes

For each gene, it uses an information theoretic score to compute its Markov Blanket.

Iterative approach (PC style)

Algorithm IMBDANET

Input: Gene expression data $G = \{g_1, g_2, \cdots, g_n\}$
Output: GRNs
1 for each target gene $T$ in $G$ do
2 ADJ_$T$ = RecogADJ ($T$),
3 for each gene $X$ in ADJ_$T$ do
4 ADJ_$X$ = RecogADJ ($X$)
5 compare ADJ_$X$ with ADJ_$T$ to determine the MB of $T$,
6 end
7 remove the IDRS from MB according to definition 3 and DPI, and get DRS of $T$
8 infer GRNs using DRS;
9 end
10 use IDS to process isolated nodes and get the final GRNs

Algorithm: recogADJ (recognize the ADJ)

Input target gene $T$, gene set $G$, threshold $\lambda$
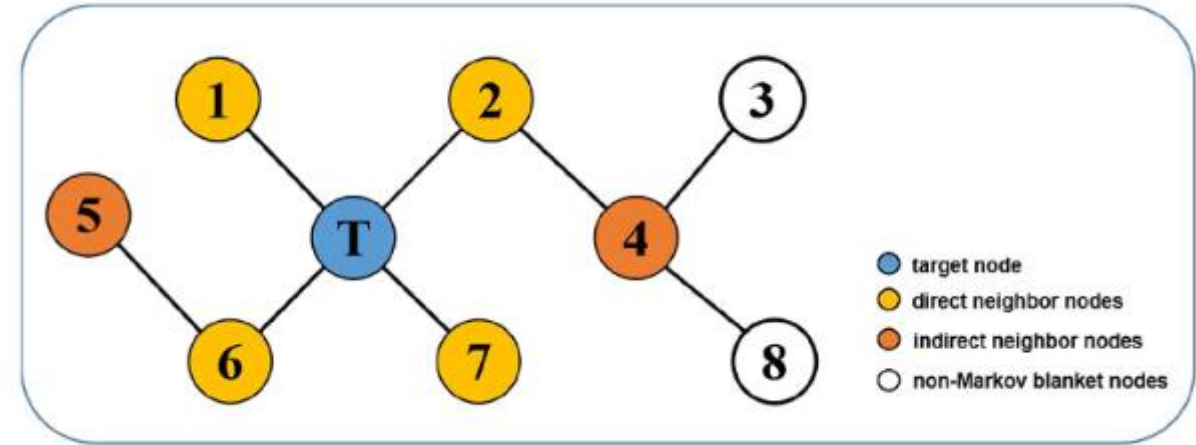Output ADJ
1 ADJ=G−T
2 for each gene $X$ in ADJ do
3 conditional independent test between $X$ and $T$
4 genes are deleted according to definition 2 and 3
5 end
6 return ADJ

$$\lambda = \min(\text{MIM}) + \text{var(MIM)} \times (\text{mean(MIM)} - \min(\text{MIM}))$$

# IMBDANET

The MB contains both directly and indirectly related nodes. To remove indirectly related nodes, they use
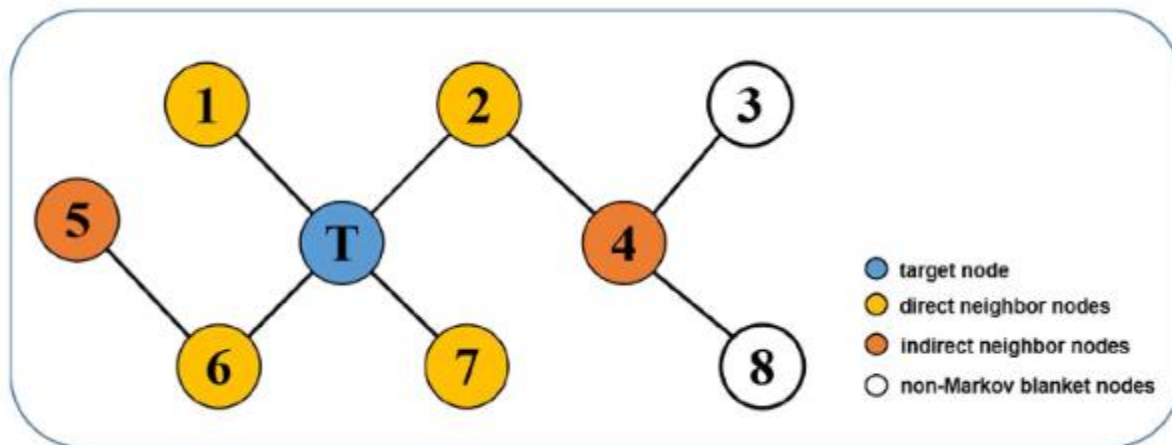
1. Data Processing Inequality
2. Definition of indirectly related nodes (dependency only when conditioning on collider)

# IMBDANET

The MB contains both directly and indirectly related nodes. To remove indirectly related nodes, they use

1. Data Processing Inequality
2. Definition of indirectly related nodes (dependency only when conditioning on collider)
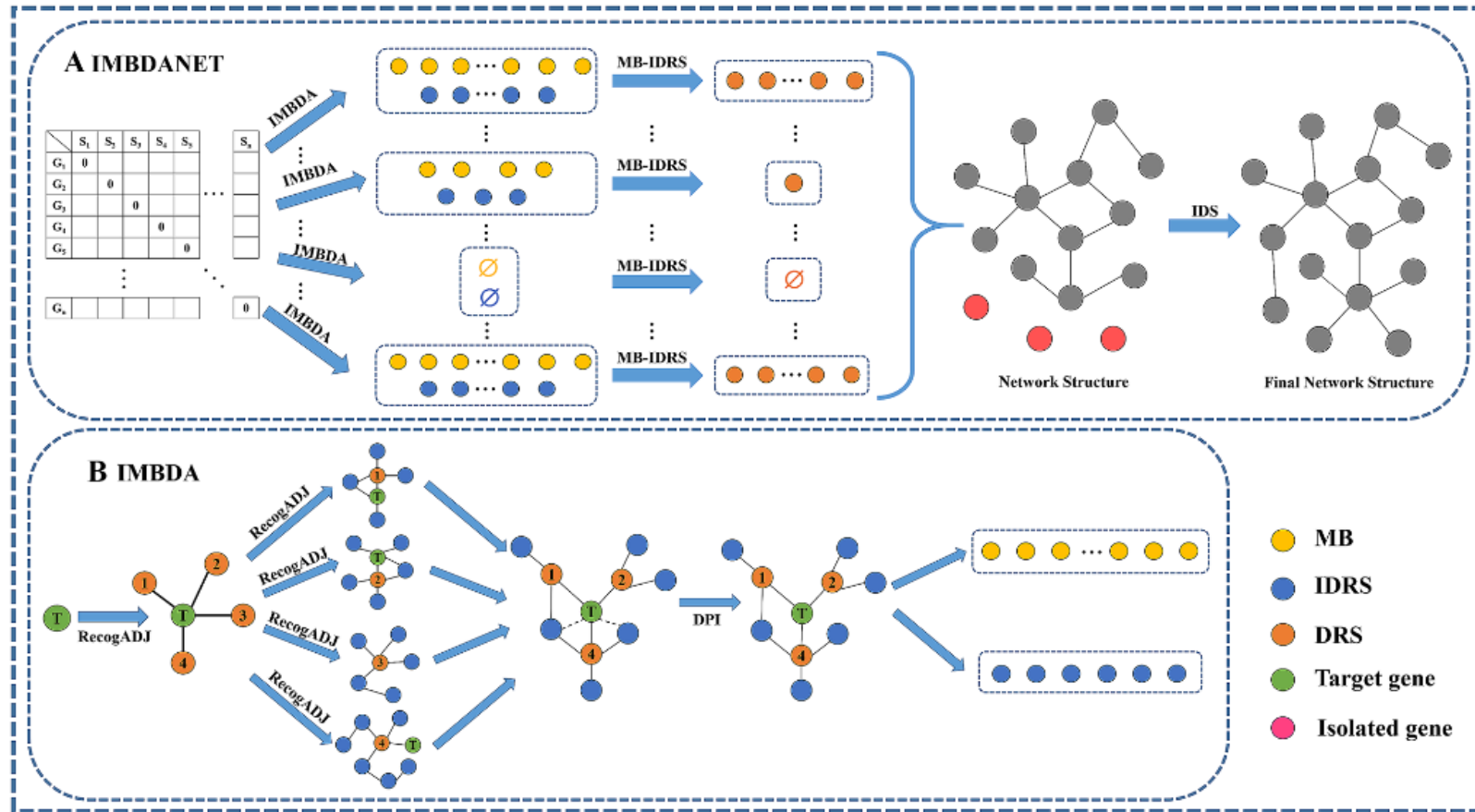
# IMBDANET

To avoid isolated genes, genes with no neighbors after pruning are connected to the "most relevant genes" according to the IDS metric.

$$\text{IDS}\left(g_i, g_c\right) = \sum_{g_c \in ISO; g_i, g_j \in G}^{i \neq j} \left[H\left(g_c | g_i\right) - H\left(g_c | g_i, g_j\right)\right]^2 \text{MI}\left(g_i, g_j\right)$$

# IMBDANET

# Undirected, unweighted gene co-expression

True Network    IMBDANET    CLR    ARACNE    MRNET    RRMRNET    MRMSn    PCA-PMI    CMI2NI

True Network    IMBDANET    RRMRNET

MRMSn    PCA-PMI    CMI2NI
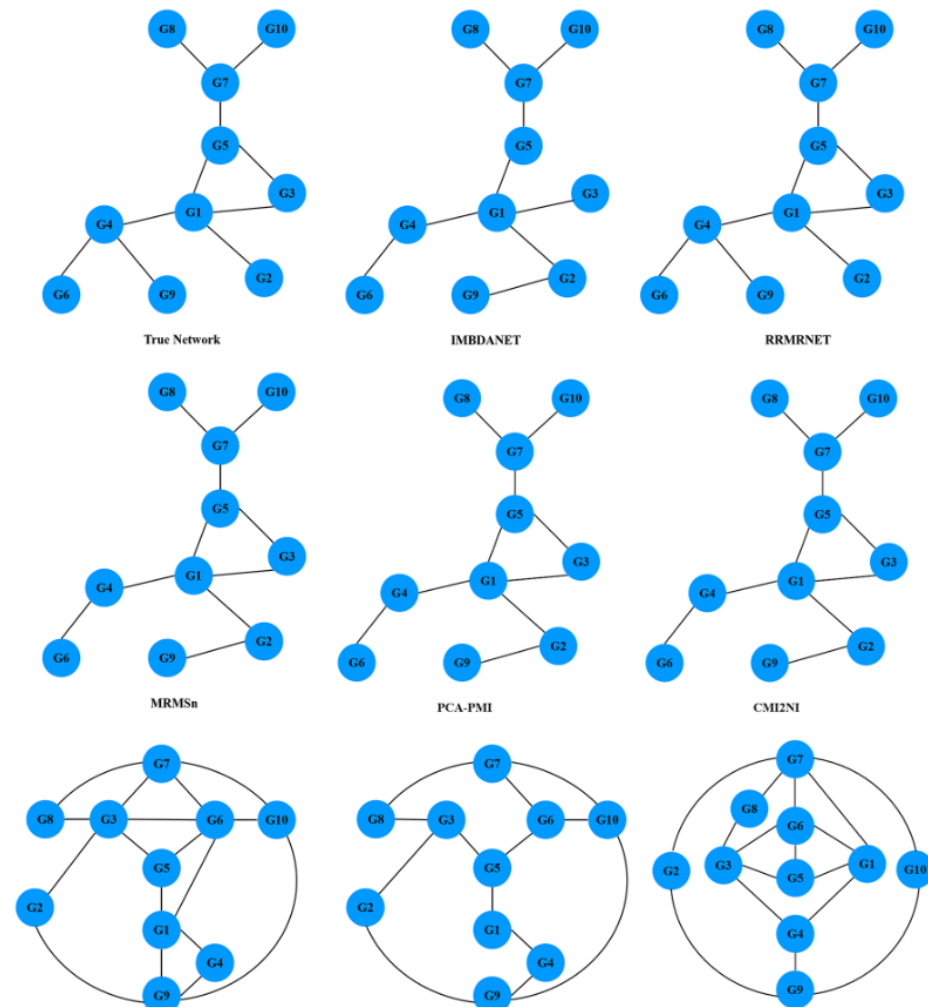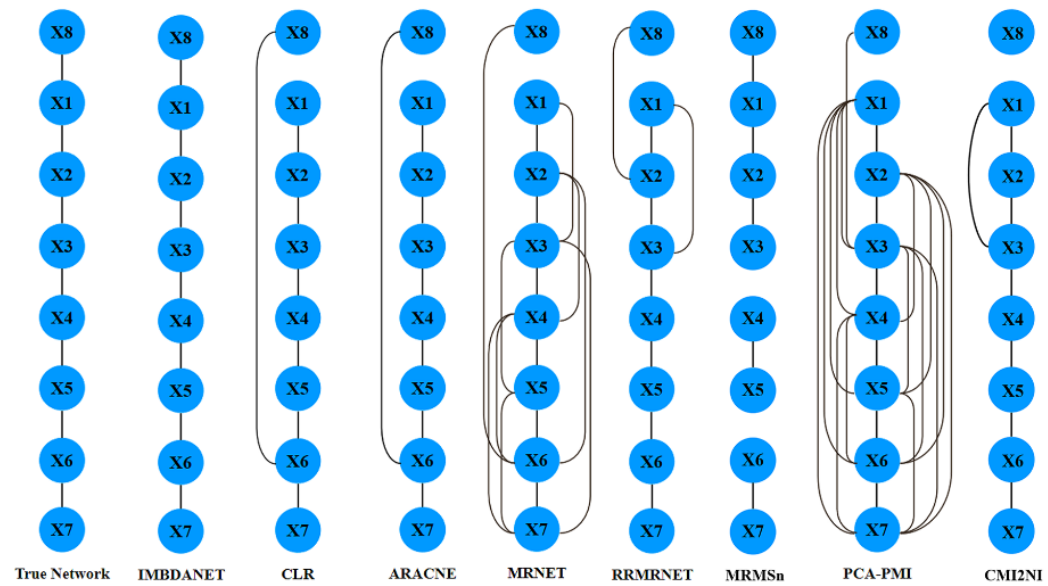
# Undirected, unweighted gene co-expression

Information theoretic scores suffer mainly from two drawbacks:
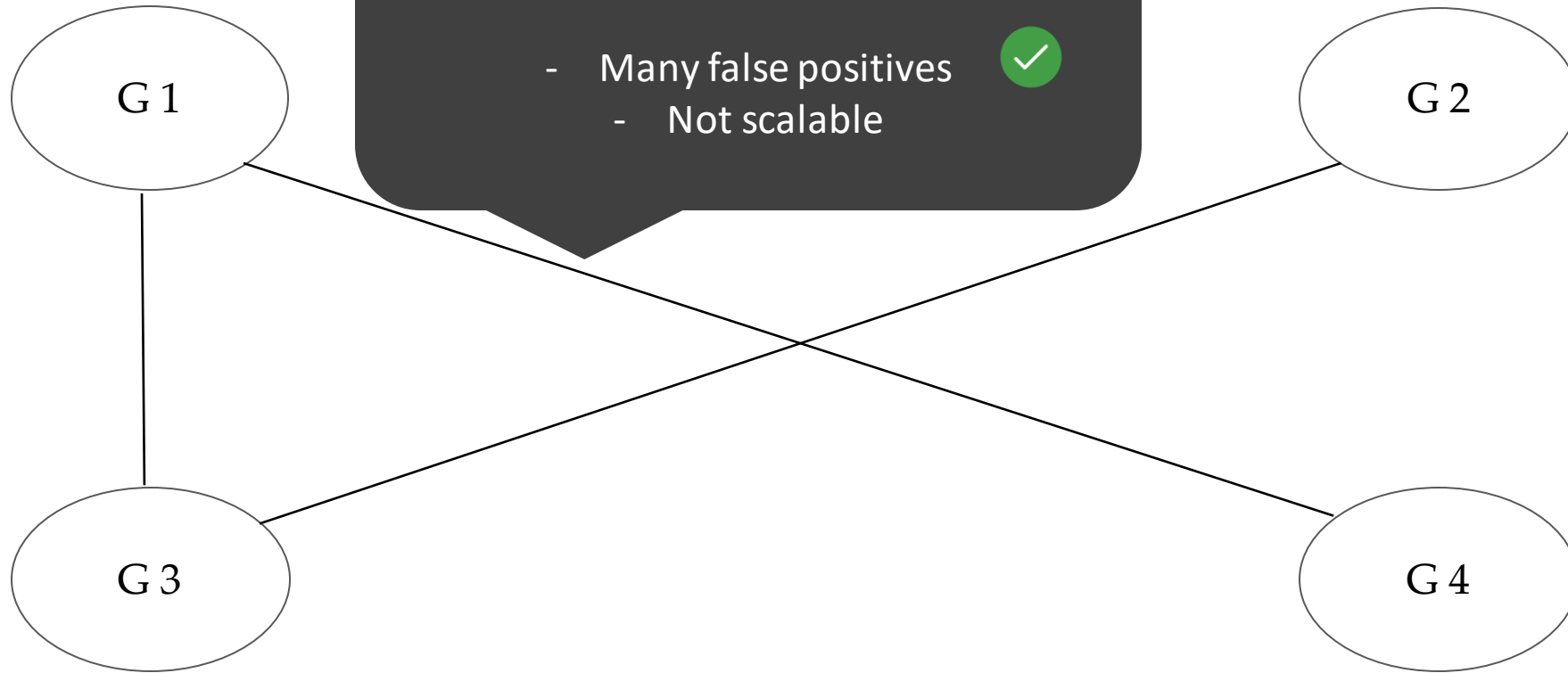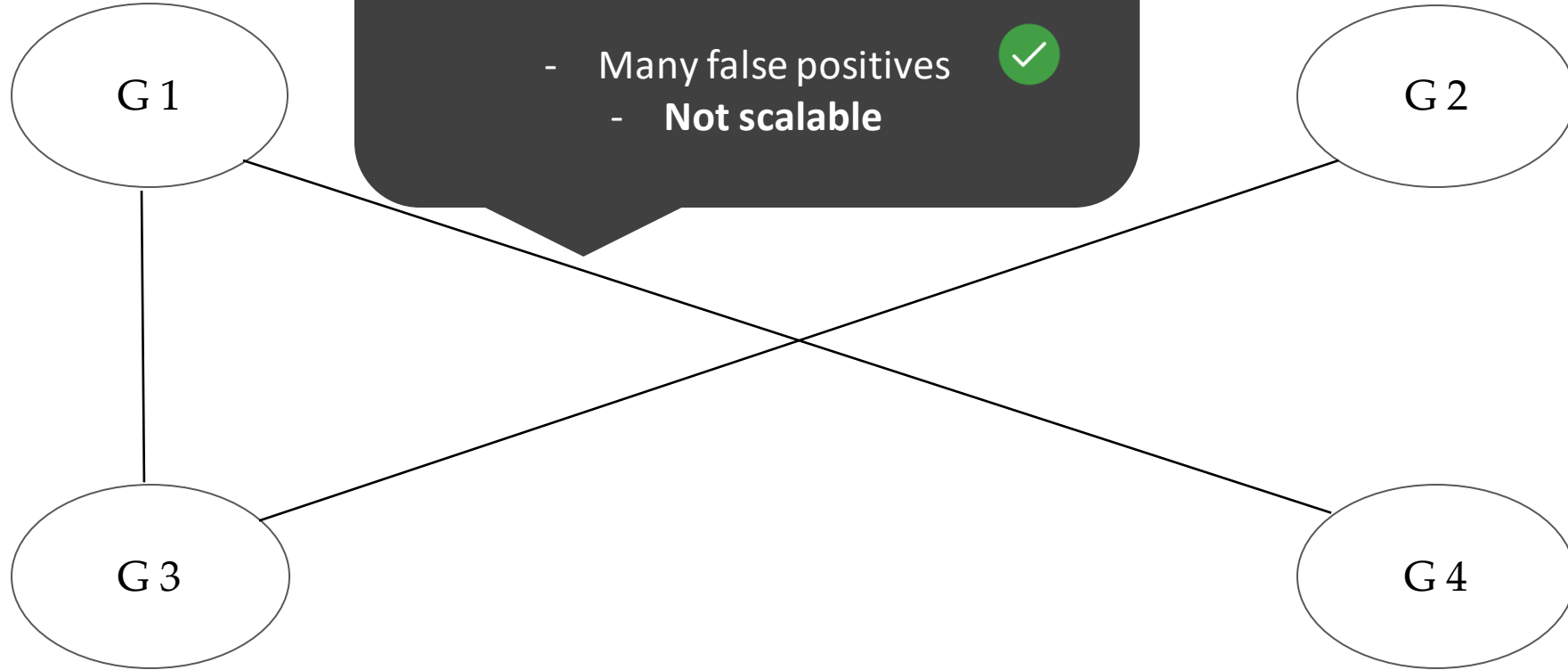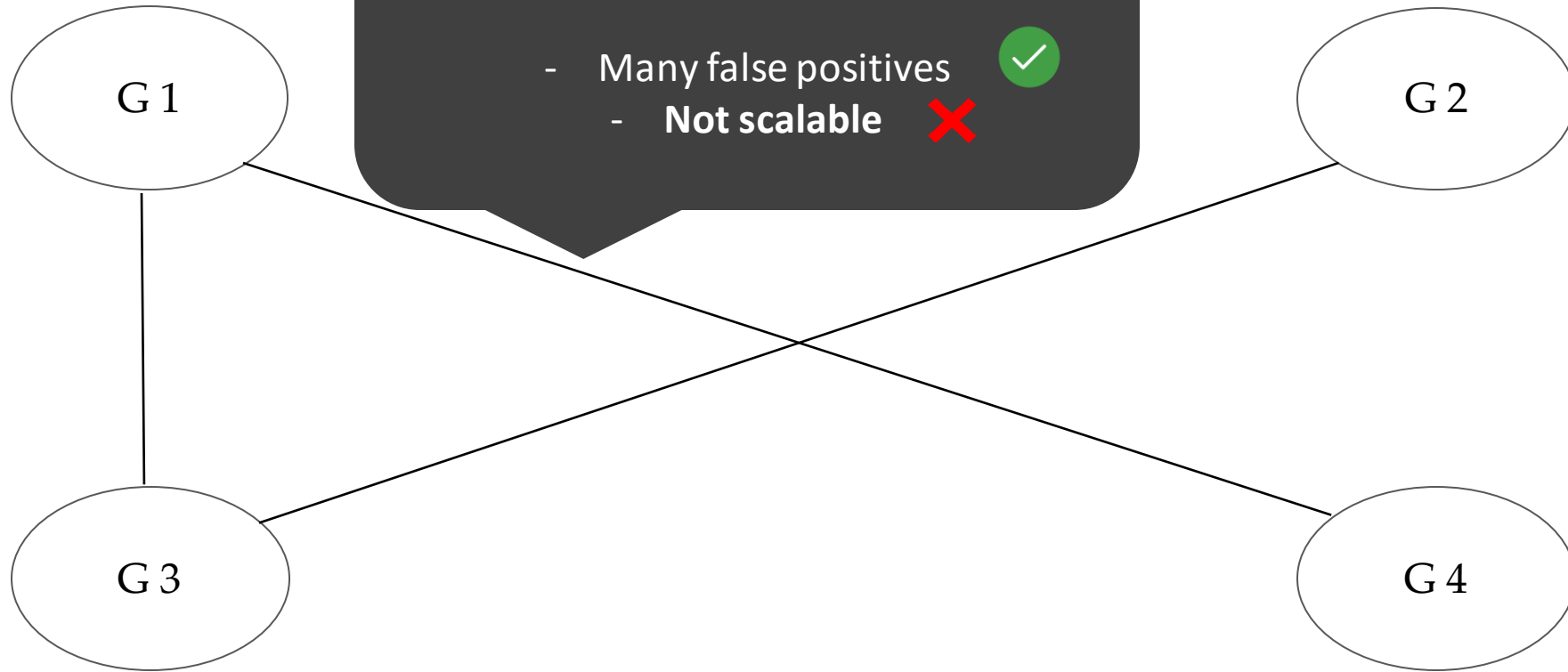
- Many false positives ✅
  - **Not scalable**

G 1

G 2

G 3

G 4

| Datasets | Variables | Samples | Type | Nodes | Edges |
|---|---|---|---|---|---|
| Reaction chain with four species | 4 | 100 | Simulated | 4 | 3 |
| Reaction chain with eight species | 8 | 250 | Simulated | 8 | 7 |
| InSilicoSize10-Yeast1-null-mutants | 10 | 10 | Simulated | 10 | 10 |
| InSilicoSize50-Yeast1-null-mutants | 50 | 50 | Simulated | 50 | 77 |
| InSilicoSize100-Yeast1-null-mutants | 100 | 100 | Simulated | 100 | 166 |
| SOS | 9 | 9 | Real | 9 | 24 |

# Undirected, unweighted gene co-expression

# My takes from the paper

- Interesting application of MBs to GRN

- Some details about the algorithms are not explicitly stated in the paper

- Lack of theoretical insights