

Identifying Causal Direction via Variational Bayesian Compression

Quang-Duy Tran Bao Duong Phuoc Nguyen Thin Nguyen

Deakin Applied Artificial Intelligence Initiative, Deakin University, Geelong, Australia

42nd International Conference on Machine Learning
July 2025



DEAKIN
APPLIED ARTIFICIAL
INTELLIGENCE INITIATIVE

Problem Setup: Cause-Effect Identification from Observational Data

Given X and Y , there are 4 possible cases for their relationship



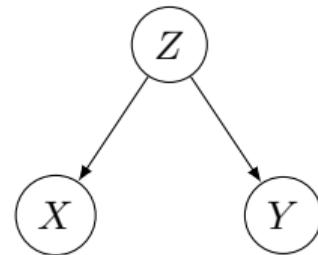
X causes Y
 $(X \rightarrow Y)$



Y causes X
 $(Y \rightarrow X)$



X and Y are independent
 $(X \perp\!\!\!\perp Y)$



X and Y have a hidden confounder Z ($X \leftarrow Z \rightarrow Y$)



DEAKIN
APPLIED ARTIFICIAL
INTELLIGENCE INITIATIVE

Problem Setup: Cause-Effect Identification from Observational Data

Given X and Y , there are 4 possible cases for their relationship



DEAKIN
APPLIED ARTIFICIAL
INTELLIGENCE INITIATIVE



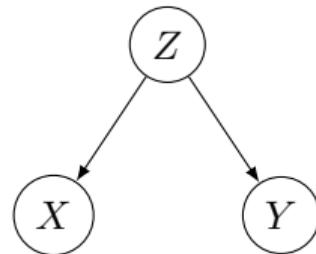
X causes Y
 $(X \rightarrow Y)$



Y causes X
 $(Y \rightarrow X)$



X and Y are independent
 $(X \perp\!\!\!\perp Y)$



X and Y have a hidden confounder Z ($X \leftarrow Z \rightarrow Y$)

Causal Asymmetry Interpretation via Kolmogorov Complexity



DEAKIN
APPLIED ARTIFICIAL
INTELLIGENCE INITIATIVE

① Algorithmic independence of conditionals¹

$$K(p(X, Y)) \stackrel{+}{=} K(p(X)) + K(p(Y | X))$$

where $K(\cdot)$ is the Kolmogorov/algorithmic complexity²

② Causal asymmetry in terms of Kolmogorov complexity³

$$\underbrace{K(p(X)) + K(p(Y | X))}_{\Delta_{X \rightarrow Y}} \stackrel{+}{\leq} \underbrace{K(p(Y)) + K(p(X | Y))}_{\Delta_{Y \rightarrow X}}$$

where $\Delta_{X \rightarrow Y}$ and $\Delta_{Y \rightarrow X}$ are causal indicator scores of $X \rightarrow Y$ and $Y \rightarrow X$ respectively

¹ D. Janzing and B. Schölkopf, "Causal inference using the algorithmic Markov condition," *IEEE Trans. Inf. Theory*, 2010.

² M. Li and P. Vitányi, *An Introduction to Kolmogorov Complexity and Its Applications*. Springer, 2019.

³ J. M. Mooij et al., "Probabilistic latent variable models for distinguishing between cause and effect," in *NIPS*, 2010.

Incomputability of Causal Indicator Scores

Issues of the Kolmogorov complexity-based scores⁴:



DEAKIN
APPLIED ARTIFICIAL
INTELLIGENCE INITIATIVE

- ① The ground-truth distributions are not available
 - ▶ Estimate the joint complexity $K(x, p_X)$ and $K(y, p_{Y|X} \mid x)$ ⁵
- ② The Kolmogorov complexity is not computable in practice
 - ▶ Approximate with the MDL principle⁶ by choosing a model from **a predefined class \mathcal{M}** that minimizes

$$L_M^{2\text{-p}}(D) := \underbrace{L_1(D \mid M)}_{\text{fitness}} + \underbrace{L_2(M)}_{\text{complexity}}, M \in \mathcal{M}$$

$$\hat{\Delta}_{X \rightarrow Y}^{2\text{-p}} := L_{M_X^*}^{2\text{-p}}(X) + L_{M_{Y|X}^*}^{2\text{-p}}(Y \mid X)$$

where M_X^* and $M_{Y|X}^*$ are models that minimize the corresponding codelengths.

⁴D. Kaltenpoth and J. Vreeken, "Causal discovery with hidden confounders using the algorithmic Markov condition," in *UAI*, 2023.

⁵A. Marx and J. Vreeken, "Formally justifying MDL-based inference of cause and effect," in *AAAI ITCI'22 Workshop*, 2022.

⁶P. D. Grünwald, *The minimum description length principle*. The MIT Press, 2007.

Classes of Models for Modeling the Conditionals



DEAKIN
APPLIED ARTIFICIAL
INTELLIGENCE INITIATIVE

Classes of models in related cause-effect identification methods:

- A set of basis functions⁷
- Splines⁸
- GP⁹ or GPLVM¹⁰
- Neural networks?
 - ▶ Bayesian Compression of Neural Networks for Identifying Causal Direction (COMIC)

⁷ A. Marx and J. Vreeken, "Telling cause from effect by local and global regression," *Knowl. Inf. Syst.*, 2019.

⁸ A. Marx and J. Vreeken, "Identifiability of cause and effect using regularized regression," in *KDD*, 2019.

⁹ J. M. Mooij et al., "Probabilistic latent variable models for distinguishing between cause and effect," in *NIPS*, 2010.

¹⁰ A. Dhir et al., "Bivariate causal discovery using Bayesian model selection," in *ICML*, 2024.

COMIC: Variational Bayesian Codelengths of Neural Networks

- Variational Bayesian (VB) codelength of a neural network^{11,12,13}:



DEAKIN
APPLIED ARTIFICIAL
INTELLIGENCE INITIATIVE

$$L_{p(\boldsymbol{\theta})}^{2\text{-p}} \left(y^{(1:N)} \mid x^{(1:N)} \right) := \underbrace{-\log p \left(y^{(1:N)} \mid x^{(1:N)}, \boldsymbol{\theta} \right)}_{\text{fitness}} - \underbrace{\log p(\boldsymbol{\theta})}_{\text{complexity}}$$

$$\begin{aligned} L_{q_\phi(\boldsymbol{\theta})}^{\text{var}} \left(y^{(1:N)} \mid x^{(1:N)} \right) &:= \mathbb{E}_{q_\phi(w)} \left[L_{p(\boldsymbol{\theta})}^{2\text{-p}} \left(y^{(1:N)} \mid x^{(1:N)}, \boldsymbol{\theta} \right) \right] - H(q_\phi(\boldsymbol{\theta})) \\ &:= \underbrace{-\mathbb{E}_{q_\phi(w)} \left[\log p \left(y^{(1:N)} \mid x^{(1:N)}, \boldsymbol{\theta} \right) \right]}_{\text{fitness}} + \underbrace{\text{KL}(q_\phi(\boldsymbol{\theta}) \parallel p(\boldsymbol{\theta}))}_{\text{complexity}} \end{aligned}$$

- The VB codelength is the upper bound of the marginal Bayesian codelength:

$$L_{p(\boldsymbol{\theta})}^{\text{Bayes}} \left(y^{(1:N)} \mid x^{(1:N)} \right) := -\log \int p \left(y^{(1:N)} \mid x^{(1:N)}, \boldsymbol{\theta} \right) p(\boldsymbol{\theta}) d\boldsymbol{\theta}$$

¹¹A. Honkela and H. Valpola, "Variational learning and bits-back coding: An information-theoretic view to Bayesian learning," *IEEE Trans. Neural. Netw.*, 2004.

¹²C. Louizos et al., "Bayesian compression for deep learning," in *NeurIPS*, 2017.

¹³L. Blier and Y. Ollivier, "The description length of deep learning models," in *NeurIPS*, 2018.

COMIC: Identifying Causal Direction via the VB Codelengths



DEAKIN
APPLIED ARTIFICIAL
INTELLIGENCE INITIATIVE

- Likelihood distribution for the conditionals $p(y | x, \theta)$:

$$p\left(y^{(1:N)} | x^{(1:N)}, \theta\right) := \prod_{i=1}^N \mathcal{N}\left(y^{(i)} | \mu\left(x^{(i)}; \theta\right), \sigma^2\left(x^{(i)}; \theta\right)\right)$$

where $\mu(\cdot; \theta)$ and $\sigma(\cdot; \theta)$ are modeled via a neural network $\mathbf{f}_Y : \mathbb{R} \times \Theta \rightarrow \mathbb{R}^2$ as

$$\mu(\cdot; \theta) = f_{Y,1}(\cdot; \theta) \text{ and } \sigma(\cdot; \theta) = \zeta(f_{Y,2}(\cdot; \theta)), \zeta : \mathbb{R} \rightarrow (0, +\infty)$$

- Distribution for encoding the marginals $p(x)$:

$$p(x) := \mathcal{N}(x | 0, 1), L_{\mathcal{N}}(x^{(1:N)}) := -\sum_{i=1}^N \log p\left(x^{(i)}\right)$$

COMIC: Causal Identifiability



DEAKIN
APPLIED ARTIFICIAL
INTELLIGENCE INITIATIVE

- Approximated causal indicator score:

$$\hat{\Delta}_{X \rightarrow Y}^{\text{var}} (\mathcal{D}^N) := L_{\mathcal{N}} \left(x^{(1:N)} \right) + L_{q_{\phi^*}(\boldsymbol{\theta})}^{\text{var}} \left(y^{(1:N)} \mid x^{(1:N)} \right)$$

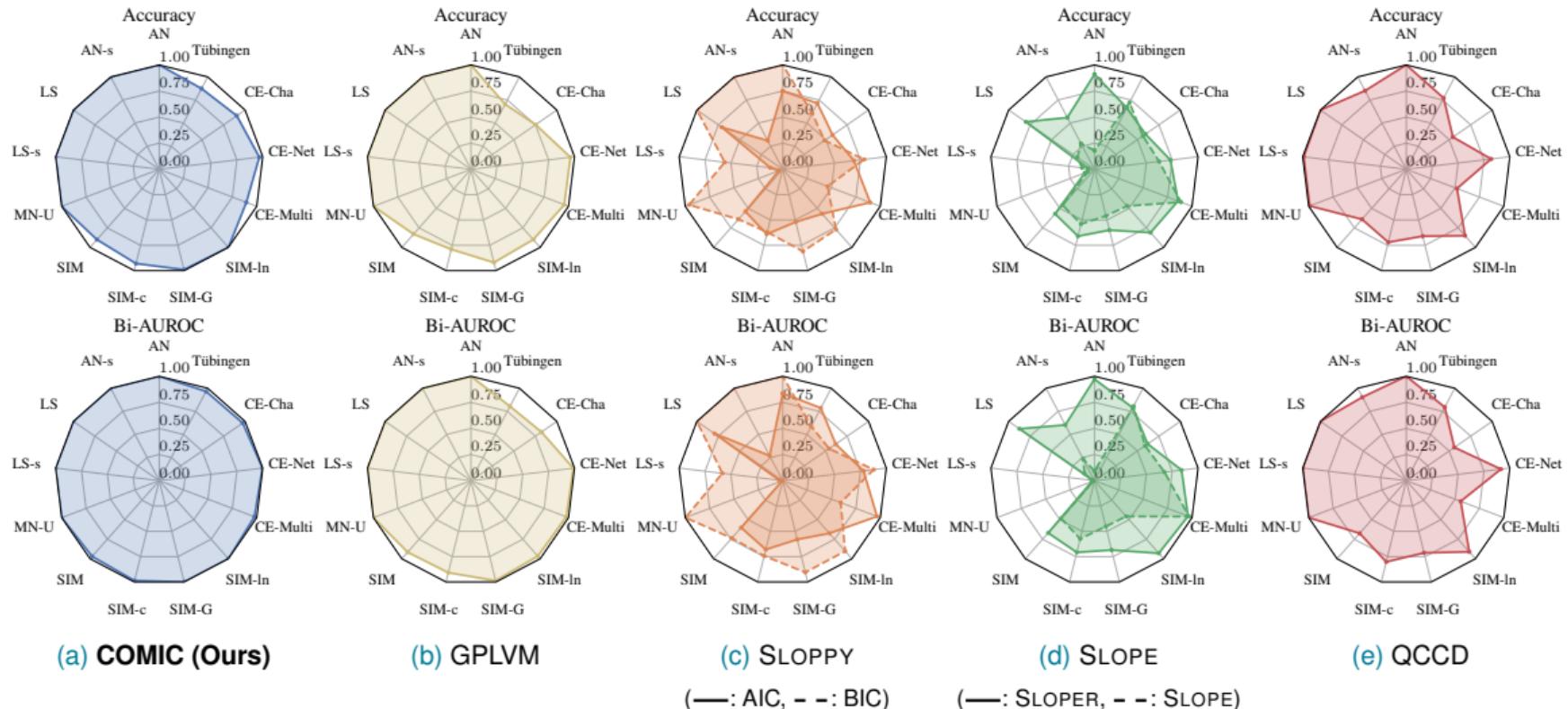
- If $q_{\phi^*}(\boldsymbol{\theta})$ converges to $p(\boldsymbol{\theta} \mid x^{(1:N)}, y^{(1:N)})$ and $\mathcal{N}(0, 1)$ is the ground-truth distribution of $p(X)$, $\hat{\Delta}_{X \rightarrow Y}^{\text{var}} (\mathcal{D}^N)$ will become

$$\begin{aligned}\Delta_{X \rightarrow Y}^{\text{Bayes}} (\mathcal{D}^N) &:= L_{\mathcal{N}} \left(x^{(1:N)} \right) + L_{p(\boldsymbol{\theta})}^{\text{Bayes}} \left(y^{(1:N)} \mid x^{(1:N)} \right) \\ &:= - \underbrace{\log p(\mathcal{D}^N \mid M_{X \rightarrow Y})}_{\text{marginal likelihood}}\end{aligned}$$

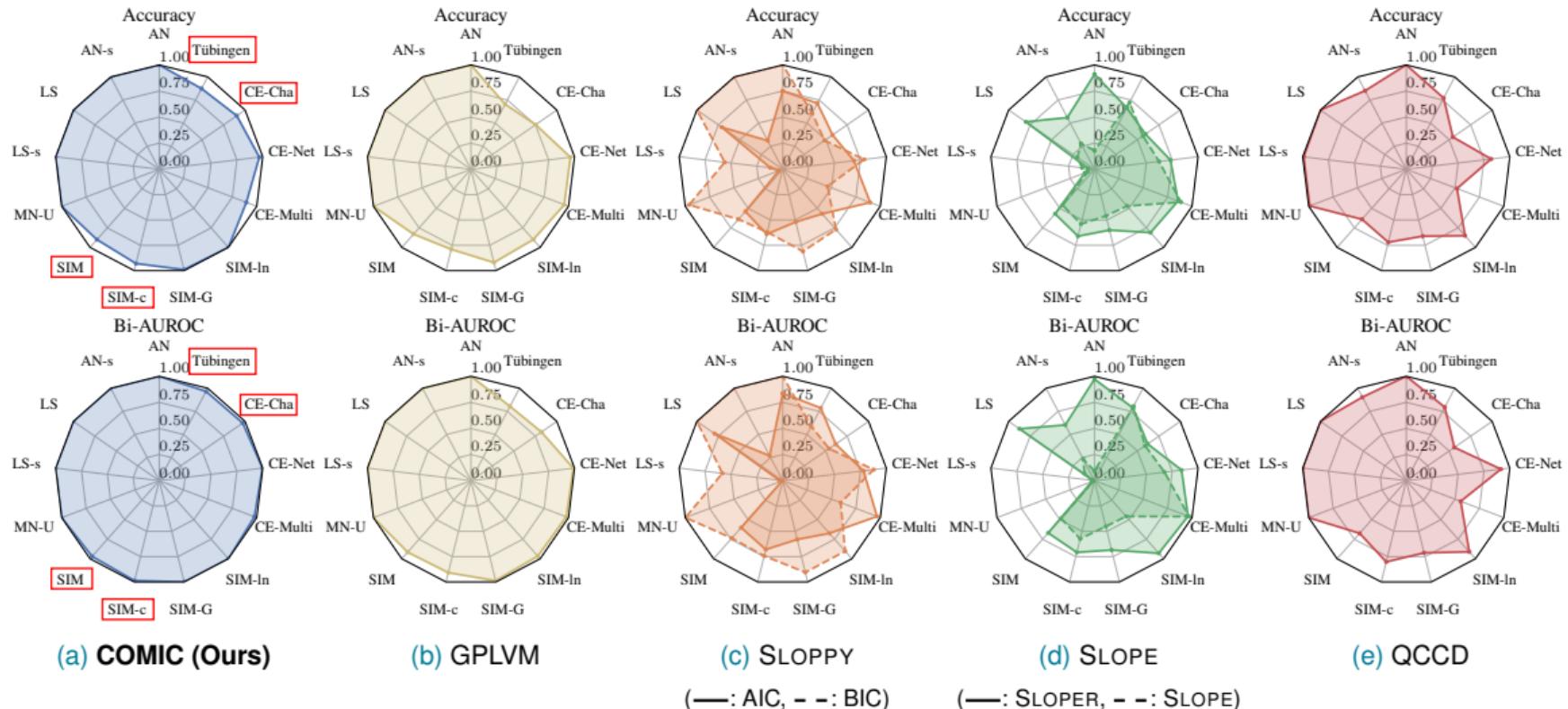
- The identifiability of COMIC is verifiable via the identifiability results of marginal likelihoods¹⁴

¹⁴ A. Dhir et al., "Bivariate causal discovery using Bayesian model selection," in ICML, 2024.

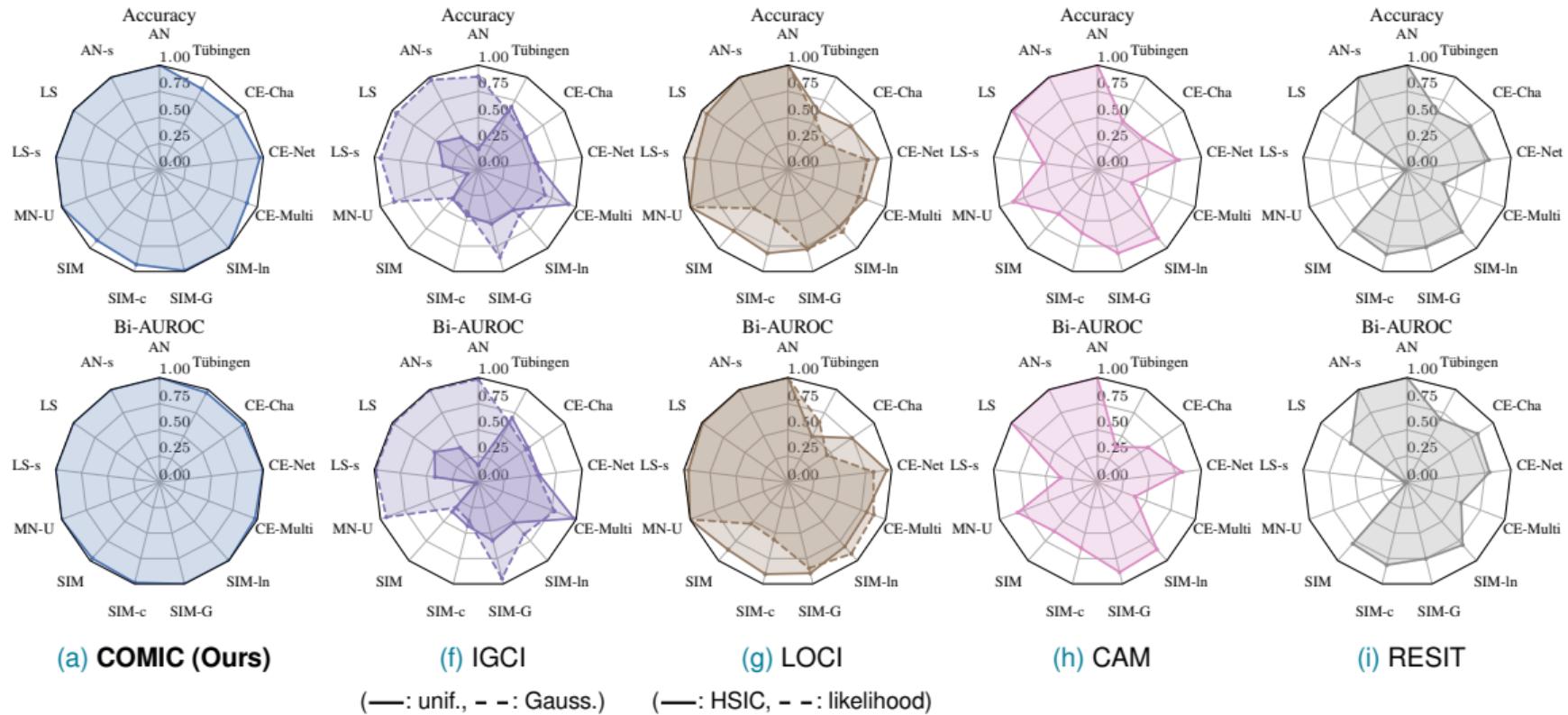
Empirical Performance on Synthetic & Real-World Benchmarks I



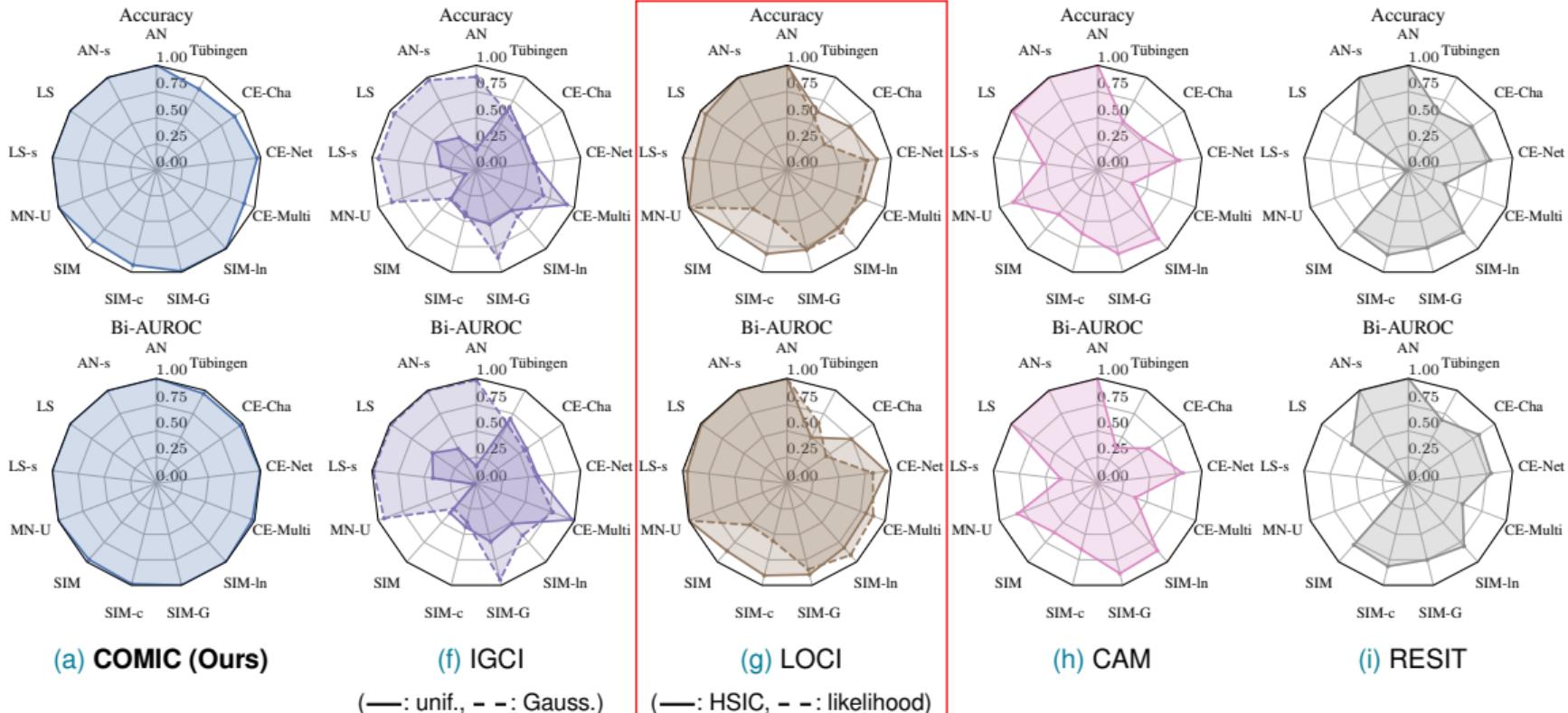
Empirical Performance on Synthetic & Real-World Benchmarks I



Empirical Performance on Synthetic & Real-World Benchmarks II



Empirical Performance on Synthetic & Real-World Benchmarks II



Summary

- ✓ We introduce COMIC—an approach based on variational Bayesian compression of neural networks for cause-effect identification
- ✓ The variational Bayesian codelengths in COMIC allow the identifiability to be justified from the marginal likelihood perspective
- ✓ The effectiveness of COMIC is demonstrated through comprehensive experiments on 13 benchmarks, which deliver promising results compared to related methods

Thank you!



- Paper: [Identifying Causal Direction via Variational Bayesian Compression](#)
- Contact: [Quang-Duy Tran <q.tran@deakin.edu.au>](mailto:Quang-Duy.Tran@deakin.edu.au)
- Code: <https://github.com/quangdzuytran/COMIC>



Preprint



Code

References I

- [1] L. Blier and Y. Ollivier, "The description length of deep learning models," in *NeurIPS*, 2018.
- [2] A. Dhir, S. Power, and M. van der Wilk, "Bivariate causal discovery using Bayesian model selection," in *ICML*, 2024.
- [3] P. D. Grünwald, *The minimum description length principle*. The MIT Press, 2007.
- [4] A. Honkela and H. Valpola, "Variational learning and bits-back coding: An information-theoretic view to Bayesian learning," *IEEE Trans. Neural. Netw.*, 2004.
- [5] D. Janzing and B. Schölkopf, "Causal inference using the algorithmic Markov condition," *IEEE Trans. Inf. Theory*, 2010.
- [6] D. Kaltenpoth and J. Vreeken, "Causal discovery with hidden confounders using the algorithmic Markov condition," in *UAI*, 2023.
- [7] M. Li and P. Vitányi, *An Introduction to Kolmogorov Complexity and Its Applications*. Springer, 2019.
- [8] C. Louizos, K. Ullrich, and M. Welling, "Bayesian compression for deep learning," in *NeurIPS*, 2017.

References II

- [9] A. Marx and J. Vreeken, "Identifiability of cause and effect using regularized regression," in *KDD*, 2019.
- [10] A. Marx and J. Vreeken, "Telling cause from effect by local and global regression," *Knowl. Inf. Syst.*, 2019.
- [11] A. Marx and J. Vreeken, "Formally justifying MDL-based inference of cause and effect," in *AAAI ITCI'22 Workshop*, 2022.
- [12] J. M. Mooij, O. Stegle, D. Janzing, K. Zhang, and Schölkopf, "Probabilistic latent variable models for distinguishing between cause and effect," in *NIPS*, 2010.