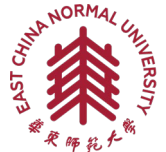


Exogenous Isomorphism for Counterfactual Identifiability

Yikang Chen^{1,†}, Dehui Du^{1,*}



ECNU

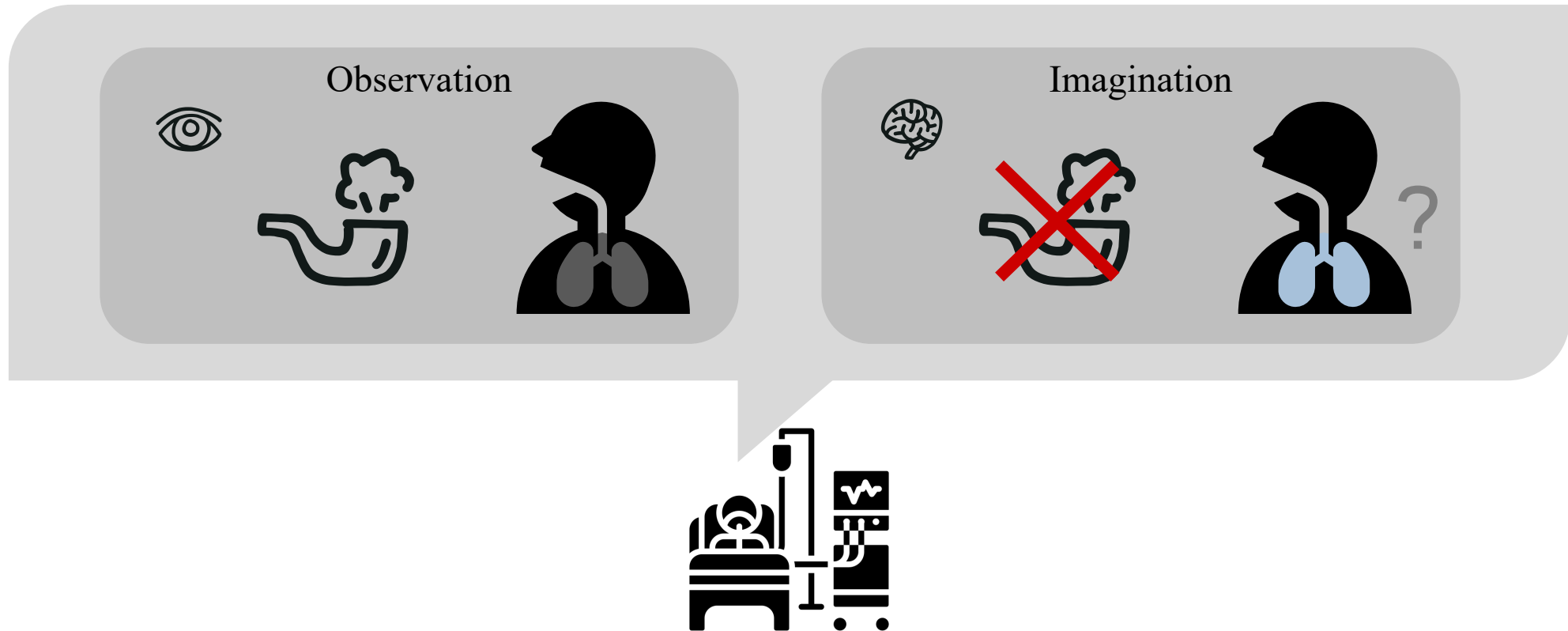
¹ Shanghai Key Laboratory of Trustworthy Computing, East China Normal University

* Corresponding author (dhdu@sei.ecnu.edu.cn)

† Presenter (yikangc620@gmail.com)

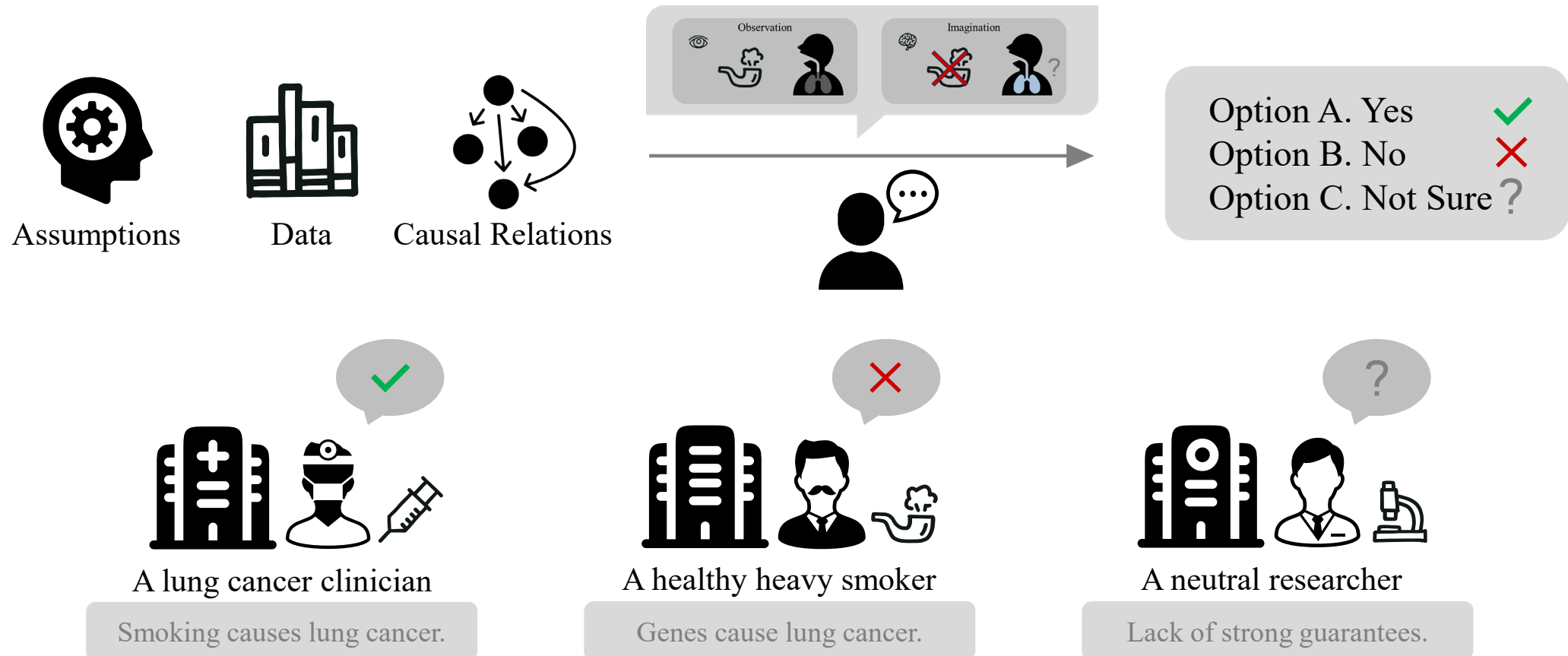
Counterfactual Identifiability

- Counterfactual is about answering “what-if” questions



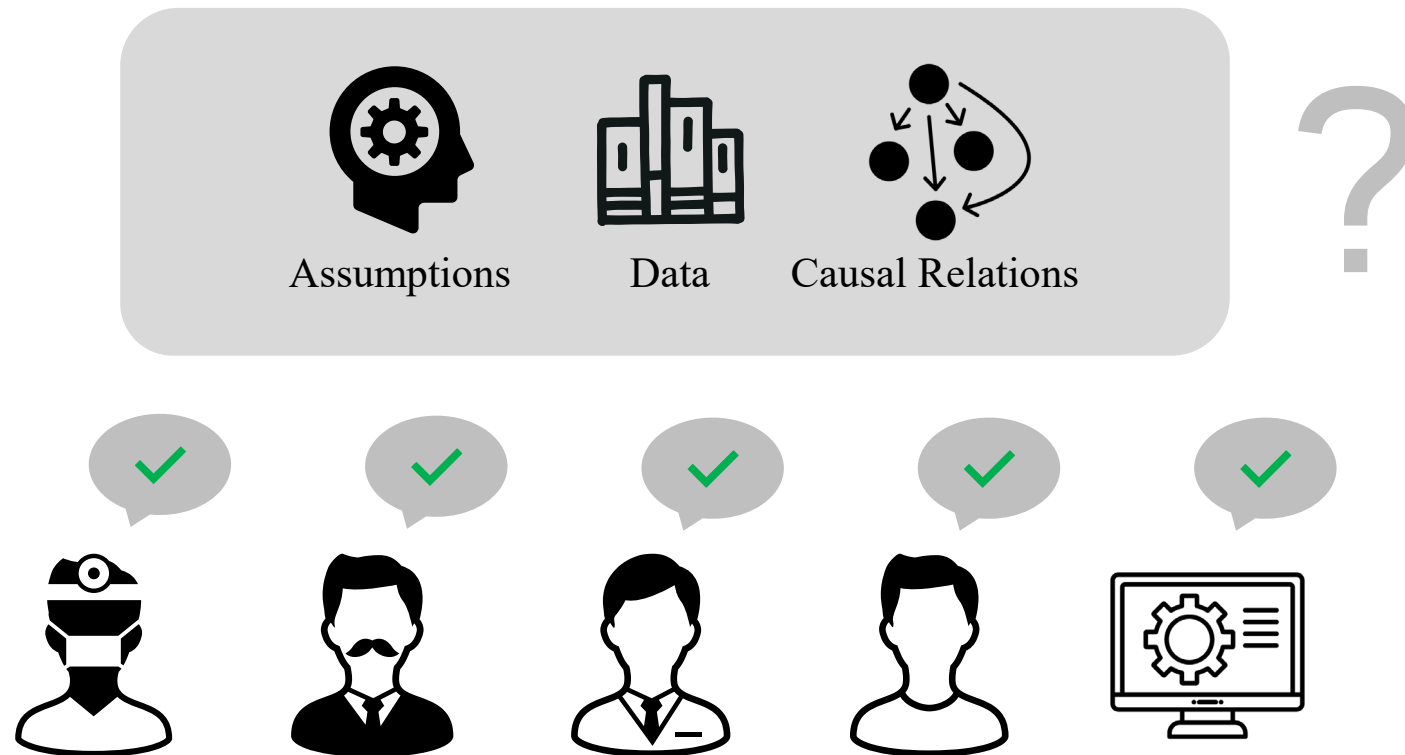
Counterfactual Identifiability

- Inconsistent answers may be produced.



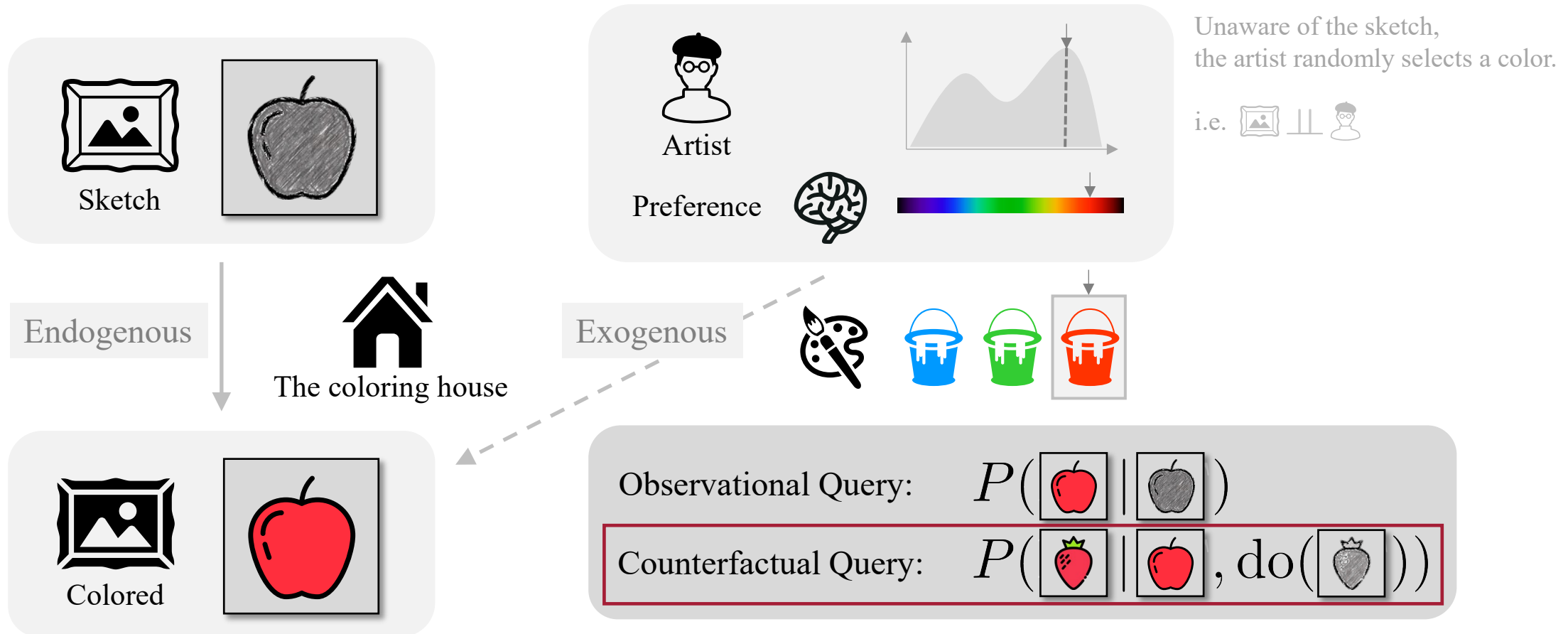
Counterfactual Identifiability

- Counterfactual Identifiability
 - How strong must the guarantees be for consistent answers?



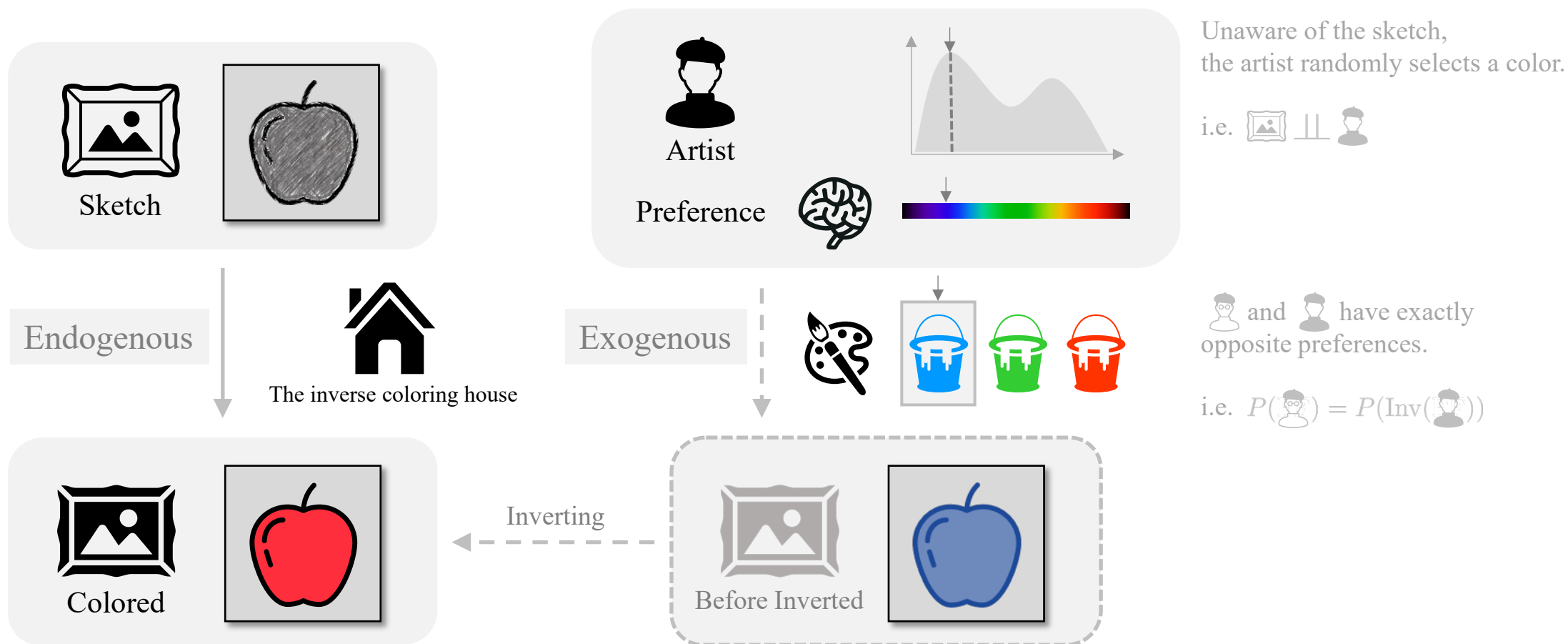
A Thought Experiment

- The coloring house  \leftarrow  ( , )



A Thought Experiment

- The inverse coloring house  \leftarrow  (, )

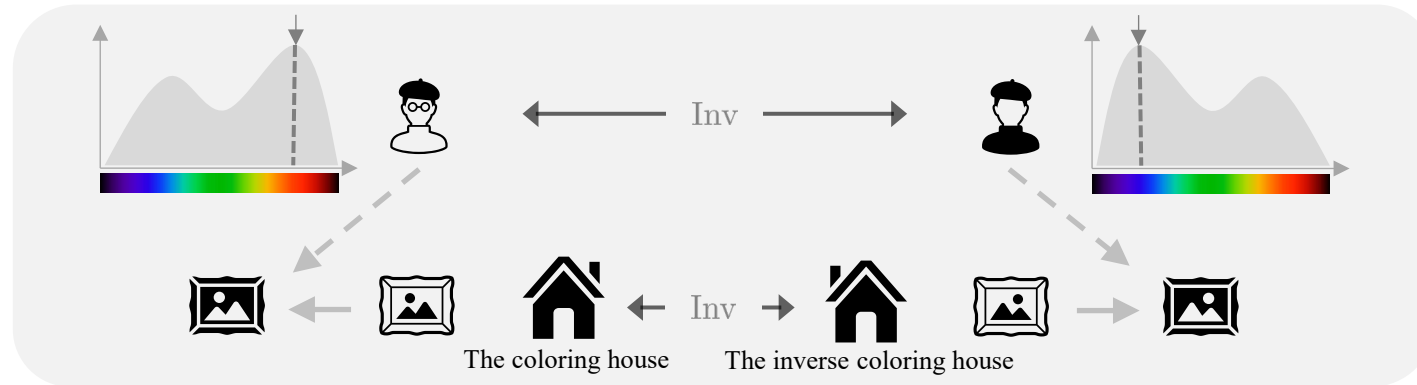


A Thought Experiment

- The consistency between two coloring houses

$$\begin{array}{l} \text{Observational Query:} \quad \begin{array}{c} \text{Image} \leftarrow \text{House}(\text{Image}, \text{Person}) \\ P^{\hat{H}}(\text{Image} | \text{Image}) \end{array} = P^{\hat{H}}(\text{Image} | \text{Image}) \end{array}$$

$$\text{Counterfactual Query:} \quad P^{\hat{H}}(\text{Image}' | \text{Image}, \text{do}(\text{Image}')) = P^{\hat{H}}(\text{Image}' | \text{Image}, \text{do}(\text{Image}'))$$

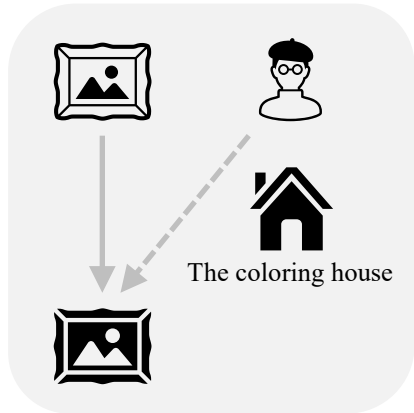


$$\text{House}(\text{Image}, \text{Person}) = \text{House}(\text{Image}, \text{Inv}(\text{Person}))$$

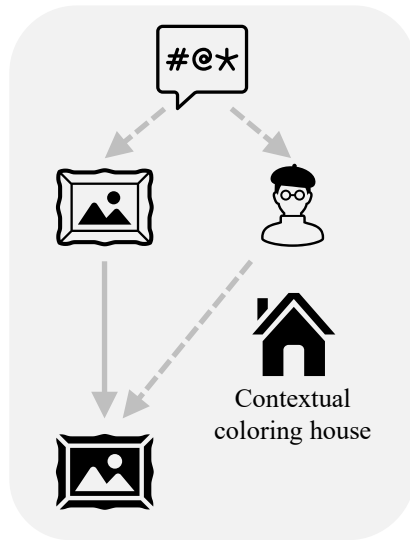
$$P_{\text{Person}} = \text{Inv}_{\#} P_{\text{Person}}$$

A Thought Experiment

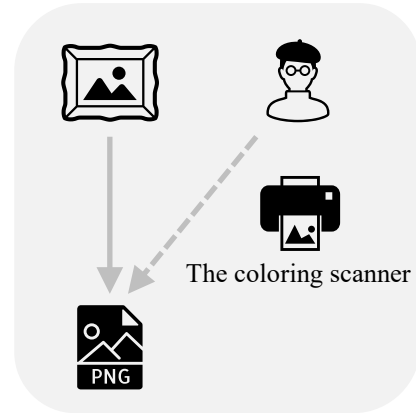
- Other cases



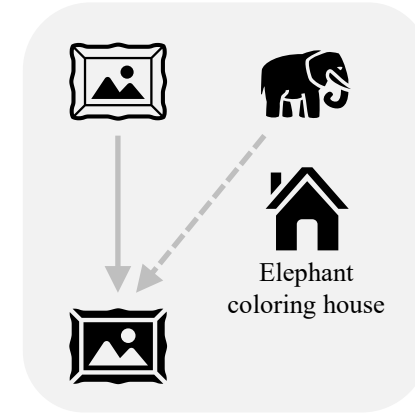
Markovian



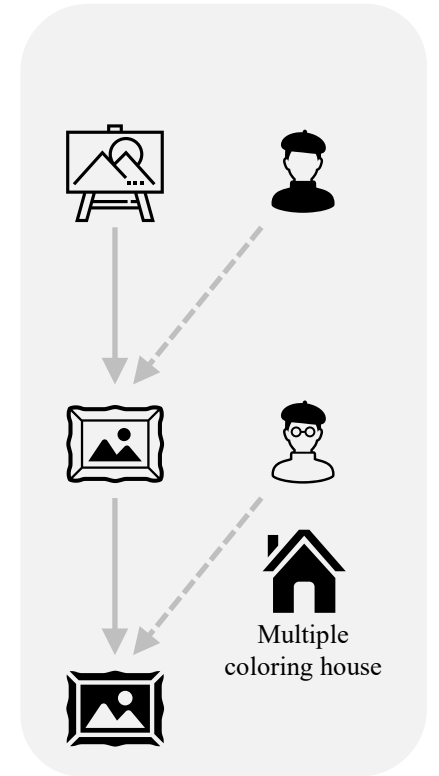
Not Markovian



Different Endogenous
Domain



Different Exogenous
Domain



Multiple Variables

Identifiability from Isomorphisms

$$\text{house}(\text{picture}, \text{person}) = \text{house}(\text{picture}, \text{Iso}(\text{person}))$$

$$P_{\text{person}} = \text{Iso}_{\#} P_{\text{person}}$$

Isomorphism	Mechanism Isomorphism	Distributional Isomorphism	Mechanistic Requirement	Structural Requirement
Counterfactual Equivalence (Peters et al., 2017)	$f^{(1)}(\mathbf{v}_{\text{pa}^{(1)}}, u^{(1)}) = f^{(2)}(\mathbf{v}_{\text{pa}^{(2)}}, u^{(1)})$	$P_{\mathbf{U}}^{(2)} = P_{\mathbf{U}}^{(1)}$	-	order
BGM Equivalence (Nasr-Esfahany et al. 2023)	$f^{(1)}(\mathbf{v}_{\text{pa}^{(1)}}, u^{(1)}) = f^{(2)}(\mathbf{v}_{\text{pa}^{(2)}}, g(u^{(1)}))$	-	bijection	graph
LCM Isomorphism (Brehmer et al. 2022)	$f^{(1)}(\mathbf{v}_{\text{pa}^{(1)}}, u^{(1)}) = f^{(2)}(\mathbf{v}_{\text{pa}^{(2)}}, \varphi(u^{(1)}))$	$P_{\mathbf{U}}^{(2)} = \varphi_{\#} P_{\mathbf{U}}^{(1)}$	diffeomorphism	graph isomorphism
Domain Counterfactual Equivalence (Zhou et al. 2024)	$f^{(1)}(\mathbf{v}_{\text{pa}^{(1)}}, h_1(u^{(1)})) = f^{(2)}(\mathbf{v}_{\text{pa}^{(2)}}, h_2(u^{(2)}))$	$P_{\mathbf{U}}^{(k)} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$	bijection	domain label
Exogenous Isomorphism (Ours)	$f^{(1)}(\mathbf{v}_{\text{pa}^{(1)}}, u^{(1)}) = f^{(2)}(\mathbf{v}_{\text{pa}^{(2)}}, h(u^{(1)}))$	$P_{\mathbf{U}}^{(2)} = h_{\#} P_{\mathbf{U}}^{(1)}$	-	order

[1] Peters, J., Janzing, D., & Scholkopf, B. (2017). *Elements of causal inference: Foundations and Learning Algorithms*. MIT Press.

[2] Nasr-Esfahany, A., Alizadeh, M., & Shah, D. (2023). Counterfactual Identifiability of Bijective Causal Models. In *Proceedings of the 40th International Conference on Machine Learning* (pp. 25733–25754). PMLR.

[3] Brehmer, J., Haan, P., Lippe, P., & Cohen, T. (2022). Weakly supervised causal representation learning. In *Advances in Neural Information Processing Systems* (pp. 38319–38331). Curran Associates, Inc..

[4] Zhou, Z., Bai, R., Kulinski, S., Kocaoglu, M., & Inouye, D. (2024). Towards Characterizing Domain Counterfactuals for Invertible Latent Causal Models. In *The Twelfth International Conference on Learning Representations*.

Identifiability from Isomorphisms

- Exogenous isomorphism

$$\mathfrak{H}(\text{img}, \text{person}) = \mathfrak{H}(\text{img}, \text{Iso}(\text{person}))$$

$$P_{\text{person}} = \text{Iso}_{\#} P_{\text{person}}$$

Recursive SCMs $\mathcal{M}^{(1)}$ and $\mathcal{M}^{(2)}$ are said to be **exogenously isomorphic**, denoted $\mathcal{M}^{(1)} \sim_{\text{EI}} \mathcal{M}^{(2)}$, if there exists a shared causal ordering \leq and function $\mathbf{h} : \Omega_{\mathbf{U}}^{(1)} \rightarrow \Omega_{\mathbf{U}}^{(2)}$ satisfying:

- Component-wise bijection $\mathbf{h} = (h_i)_{i \in \mathcal{I}}$, where each $h_i : \Omega_{U_i}^{(1)} \rightarrow \Omega_{U_i}^{(2)}$ is a bijection
- Exogenous distribution isomorphism $P_{\mathbf{U}}^{(2)} = \mathbf{h}_{\#} P_{\mathbf{U}}^{(1)}$
- Causal mechanism isomorphism $f^{(1)}(\mathbf{v}_{\text{pa}(1)}, u^{(1)}) = f^{(2)}(\mathbf{v}_{\text{pa}(2)}, h(u^{(1)}))$

- Exogenous isomorphism implies counterfactual consistency

$$\mathcal{M}^{(1)} \sim_{\text{EI}} \mathcal{M}^{(2)} \implies \mathcal{M}^{(1)} \sim_{\mathcal{L}_3} \mathcal{M}^{(2)}.$$

$\sim_{\mathcal{L}_3}$ is an equivalence relation over Structural Causal Models (SCMs), indicating that the models yield the same answers to any counterfactual statement.

- This includes counterfactual outcomes, counterfactual effects, joint counterfactuals, and nested counterfactuals.
- From the perspective of the Pearl Causal Hierarchy (PCH), $\sim_{\mathcal{L}_3}$ implies that the complete counterfactual layer is identical, meaning the models are indistinguishable across the entire hierarchy.

To Achieve Exogenous Isomorphism

- EI-identifiability problem

$$\begin{aligned}\mathcal{A} \models \quad & \mathcal{M}^{(1)} \sim_{\square} \mathcal{M}^{(2)} \\ \implies & \mathcal{M}^{(1)} \sim_{\text{EI}} \mathcal{M}^{(2)} \\ \implies & \mathcal{M}^{(1)} \sim_{\mathcal{L}_3} \mathcal{M}^{(2)}\end{aligned}$$

Under what assumptions \mathcal{A} does the model equivalence class \sim_{\square} imply \sim_{EI} ?

Bijjective SCMs

Assumption (a). The causal mechanisms are **bijjective** with respect to the exogenous variables.

Assumption (b). The structural causal model is **Markovian**.

Assumption (c). The **causal order** is given.

Assumption (d). The **observational distribution** is available.

Assumption (e). The model induces identical **counterfactual transport** (e.g., Knothe–Rosenblatt transport).

Triangular Monotonic SCMs

Assumption (a). The causal mechanisms are **monotonic** with respect to the exogenous variables.

Assumption (b). The structural causal model is **Markovian**.

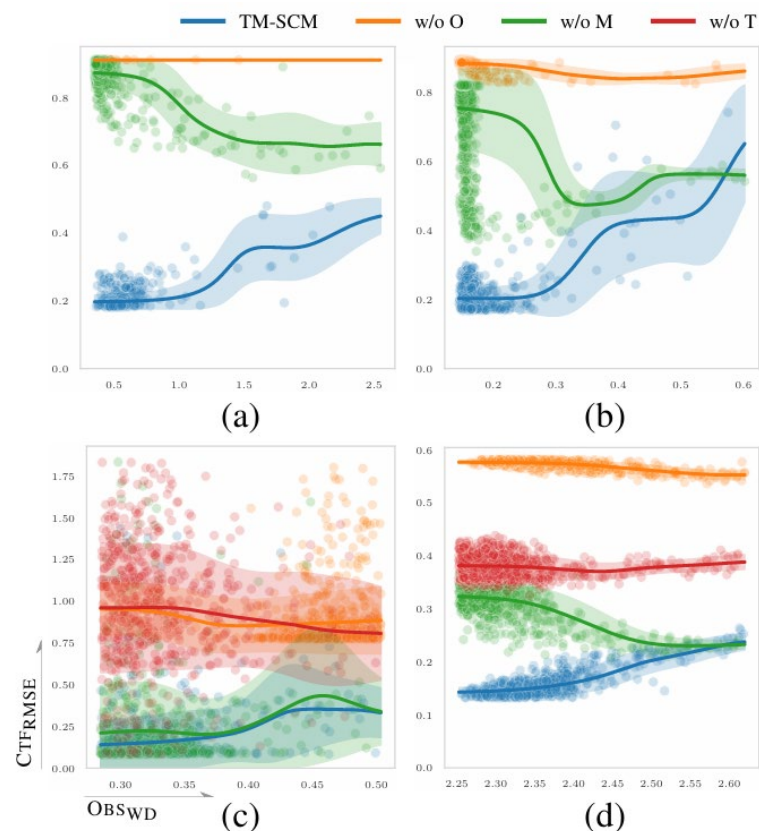
Assumption (c). The **causal order** is given.

Assumption (d). The **observational distribution** is available.

- The model equivalence classes that satisfy the above assumptions entails EI-identifiability.
- As a result, counterfactual identifiability are ensured within these classes.

Ablation Study

- Ablation experiments on synthetic datasets with neural network-based Triangular Monotonic SCMs provide empirical support for the theoretical results.



METHOD		ER-DIAG-50	ER-TRIL-50
DNME	-	0.53 ± 0.05	0.51 ± 0.12
	w/o O	0.78 ± 0.05	0.89 ± 0.10
	w/o M	0.62 ± 0.04	0.58 ± 0.10
TNME	-	0.47 ± 0.05	0.55 ± 0.12
	w/o O	11.24 ± 20.98	6.41 ± 9.84
	w/o M	0.62 ± 0.04	0.73 ± 0.21
CMSM	-	0.37 ± 0.05	0.42 ± 0.12
	w/o O	2.64 ± 3.72	2.12 ± 2.49
	w/o M	1.69 ± 2.60	0.75 ± 0.49
	w/o T	0.64 ± 0.05	1.25 ± 1.29
TVSM	-	0.46 ± 0.05	0.50 ± 0.12
	w/o O	0.79 ± 0.04	0.88 ± 0.10
	w/o M	0.53 ± 0.05	0.53 ± 0.11
	w/o T	0.67 ± 0.05	0.78 ± 0.12

Summary

- We introduced exogenous isomorphism ensuring consistency across all counterfactual queries.
- We identified and proved that two specific classes of models satisfy EI-identifiability.
- Empirical results on neural SCMs and synthetic datasets support our theoretical claims.
- This work provides a principled foundation for learning counterfactually reliable models.

THANKS

Full paper is available at:
<https://arxiv.org/abs/2505.02212>



Code is available at:
<https://github.com/cyisk/tmscm>

