

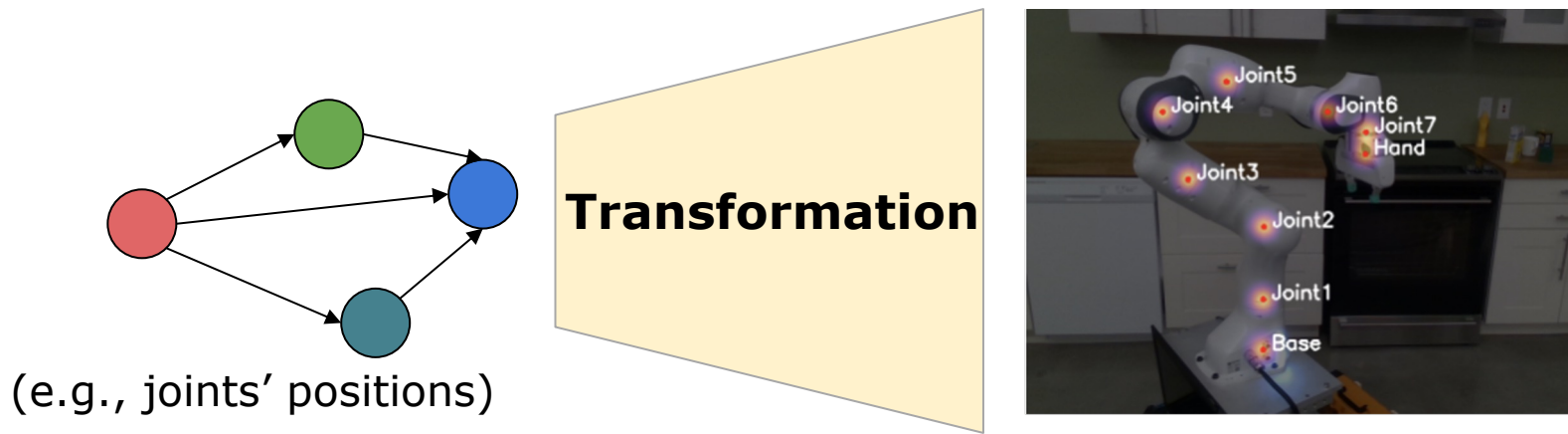


find the paper

Burak Varıcı¹ Emre Acartürk¹ Karthikeyan Shanmugam² Ali Tajer¹

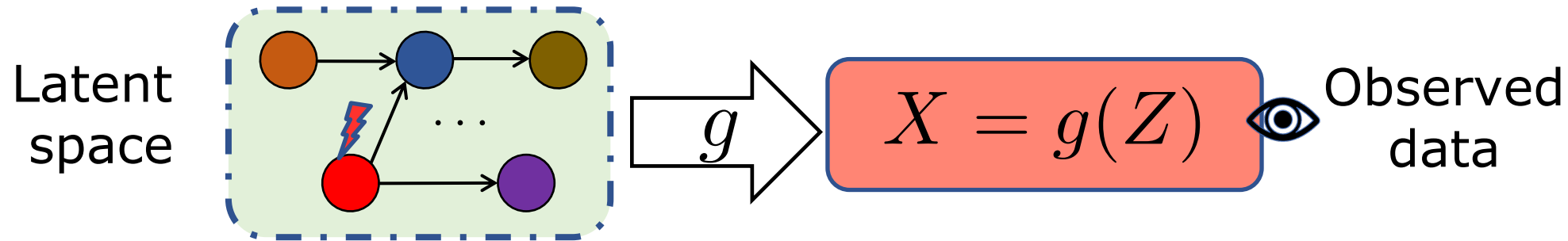
¹Rensselaer Polytechnic Institute ²Google Research India

CRL from Interventions



".. learn a representation (partially) exposing the unknown causal structure, e.g., which variables describe the system, and their relations .." Schölkopf et al., 2021

Generic goal: Invert the unknown transformation to recover **1) latent representation** and **2) the latent causal structure**



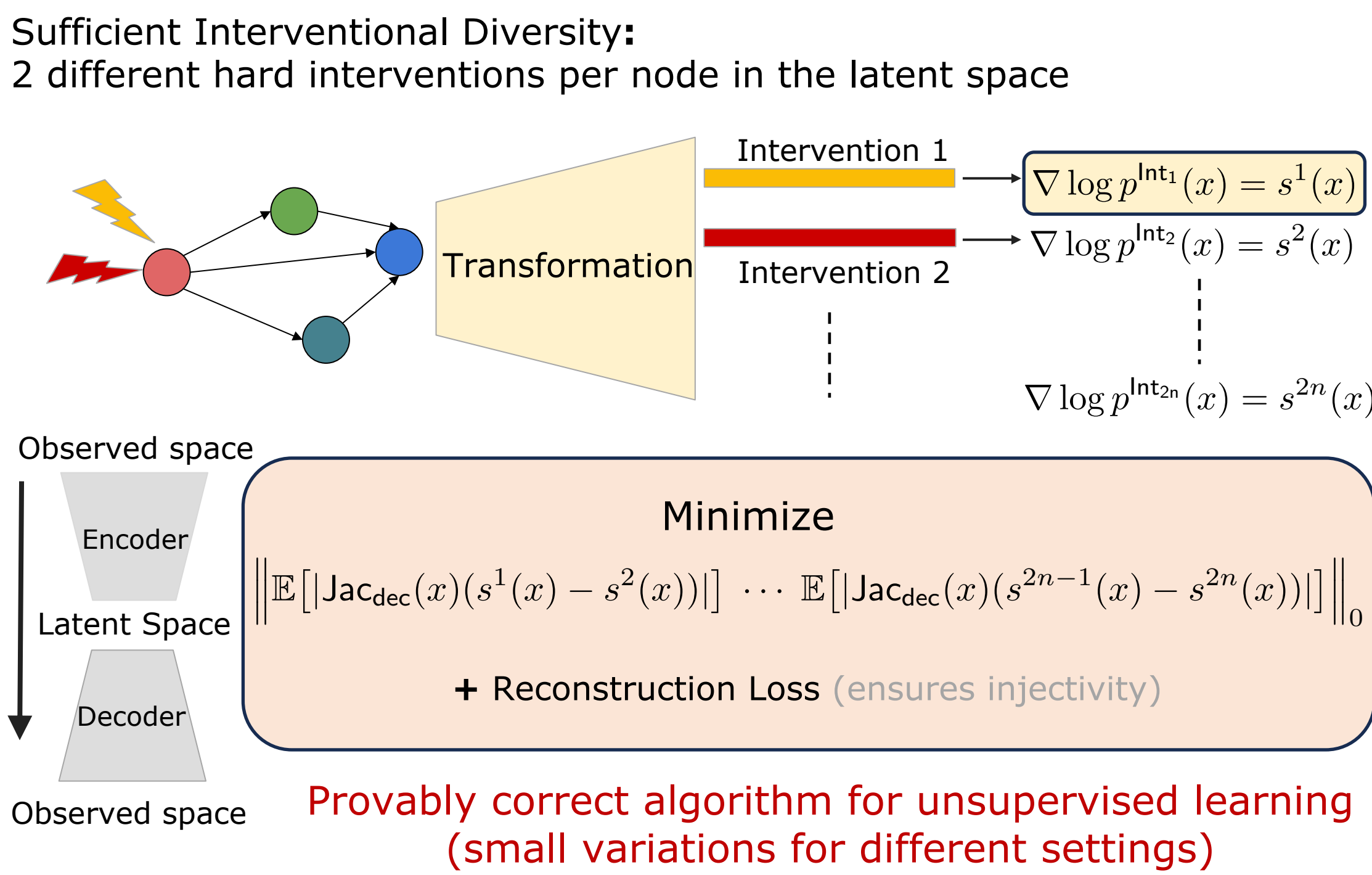
- Identifiability:** Conditions for uniquely recovering Z and g_Z
- Achievability:** Provably correct algorithms to recover Z and g_Z

Our contributions

Latent model	Transform	Interv. / node	Main results
Nonparametric + Nonparametric	Linear	2 hard	= perfect ID
Sufficiently nonlinear	Linear	+ 1 hard (soft)	= perfect ID (true DAG + Markov)
Nonparametric	Linear	+ 1 hard (soft)	= perfect ID (ID up to ancestors)

provably correct algorithms for all settings

Algorithm Overview



Experiments

Non-linear latent model: $Z_i = \sqrt{Z_{\text{pa}(i)}^\top A_{p,i} Z_{\text{pa}(i)}} + N_{p,i}$ $n=8$ latent variables

Input score differences ($s_X - s_X^m$): Perfect score oracle or Sliced Score Matching

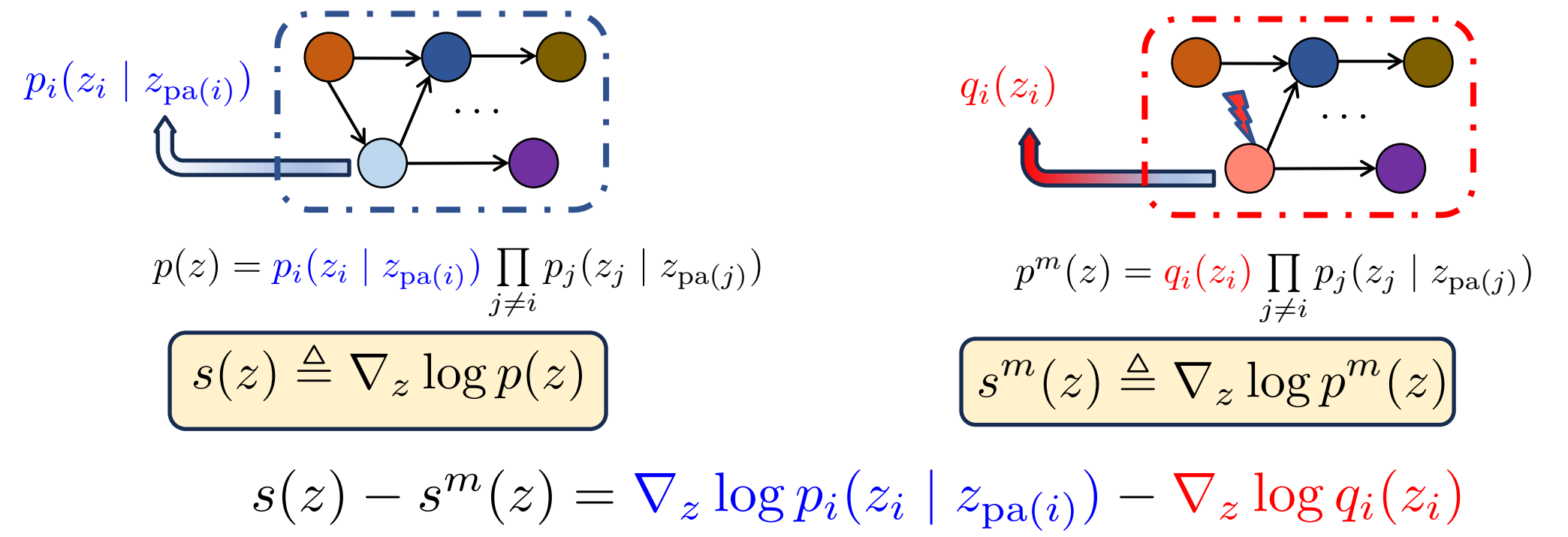
Non-linear transform: $X = \tanh(T \cdot Z)$

Linear transform: $X = T \cdot Z$

Two hard / node					One hard / node				
Obs. dim	Norm. Z error	DAG error (SHD)	Norm. Z error	DAG error (SHD)	Obs. dim	Norm. Z error	DAG error (SHD)	Norm. Z error	DAG error (SHD)
8	0.16	1.56	0.70	11.9	8	0.50	5.4	0.75	10.3
25	0.20	1.55	0.68	10.5	25	0.51	6.0	0.78	8.9
40	0.21	1.14	0.71	11.8	40	0.50	5.3	0.61	11.9

score oracle noisy scores score oracle noisy scores

Why score functions?



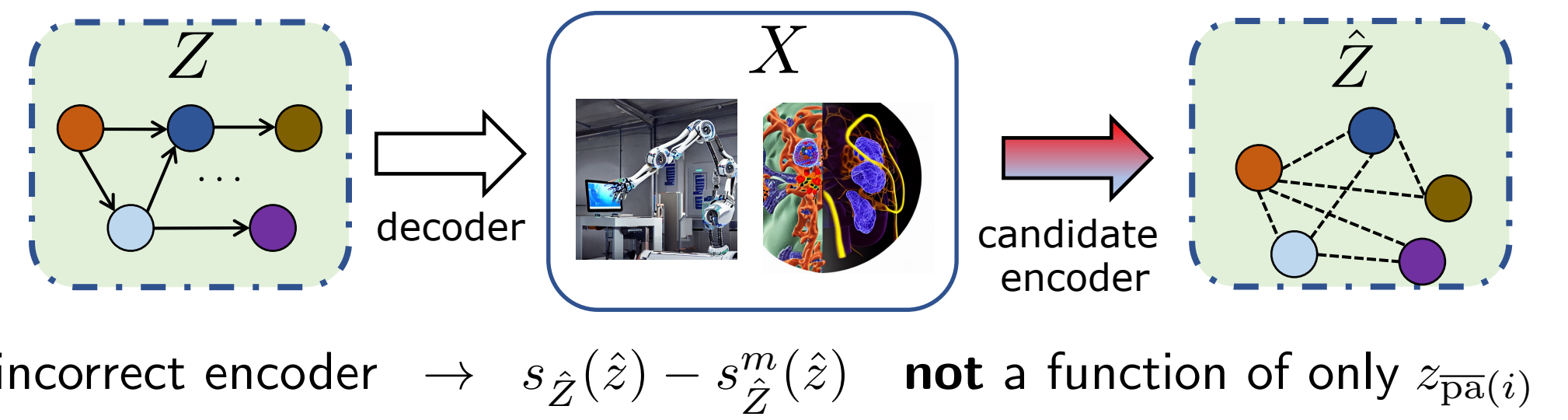
Score functions contain all information about latent DAGs

node i intervened: $s(z) - s^m(z)$ becomes a function of only $z_{\overline{\text{pa}(i)}}$

$$s(z) - s^m(z) = [0 \ 0 \ \times \ 0 \ \times \ 0]^\top$$

coordinates of parents of node i

Methodology



estimated score differences cannot be sparser than true score differences

Min. score variations over environment pairs = correct encoder

$$s_{\hat{Z}}(\hat{z}) - s_{\hat{Z}}^m(\hat{z}) = [J_{\text{decoder}}(\hat{z})]^\top (s_X(x) - s_X^m(x))$$

Results

Nonparametric transform

Interventional discrepancy: $\frac{\partial q_i(z_i)}{\partial z_i} \neq \frac{\partial \tilde{q}_i(z_i)}{\partial z_i}$ almost everywhere

Theorem : Observational data and **two hard** interventions/node **Perfect ID**

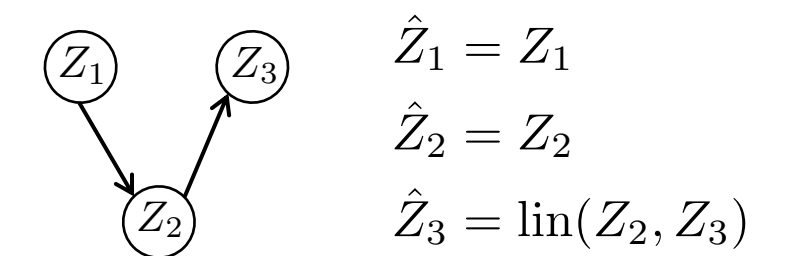
von Kügelgen et al.(2023): **Coupled** two hard + **faithfulness** = Perfect ID

Linear transform + nonlinear latents

A1 (nonlinearity): $\text{rank}(\text{im}(s - s^m)) = |\overline{\text{pa}(i)}|$ e.g., 2-layer NN with additive noise

Theorem : Linear transform + **one** intervention/node + **A1 hard: Perfect ID ; soft: Perfect DAG + Markov Property**

Going beyond 'ID ancestors' for soft (nonlinearity = up to ancestors in Zhang'23)



Linear transform + any latents

A2 (mild): $\forall j \in \text{pa}(I^m), \frac{[s - s^m]_j}{[s - s^m]_{I^m}} \neq \text{constant}$ e.g., weights change in linear model

Theorem : Linear transform + **one** intervention/node + **A2 hard : Perfect ID ; soft: ID up to ancestors**

No parametric restrictions on latents (linear models on Squires'23, Buchholz'23)

reach out at: burakvarici@gmail.com