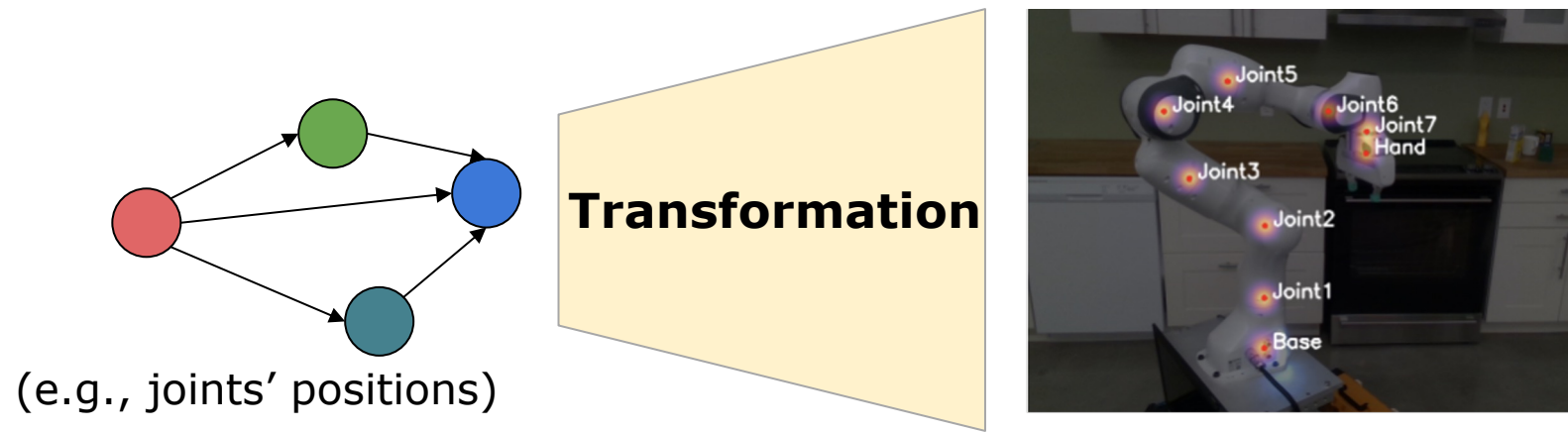


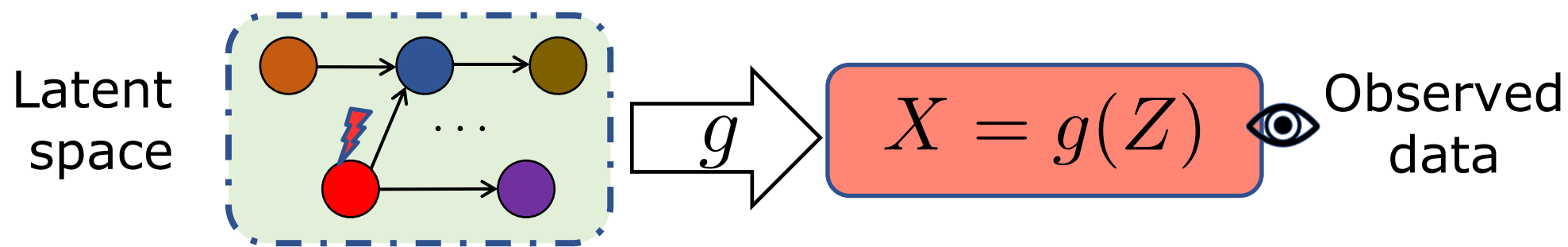
# General Identifiability and Achievability for Causal Representation Learning

## 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**

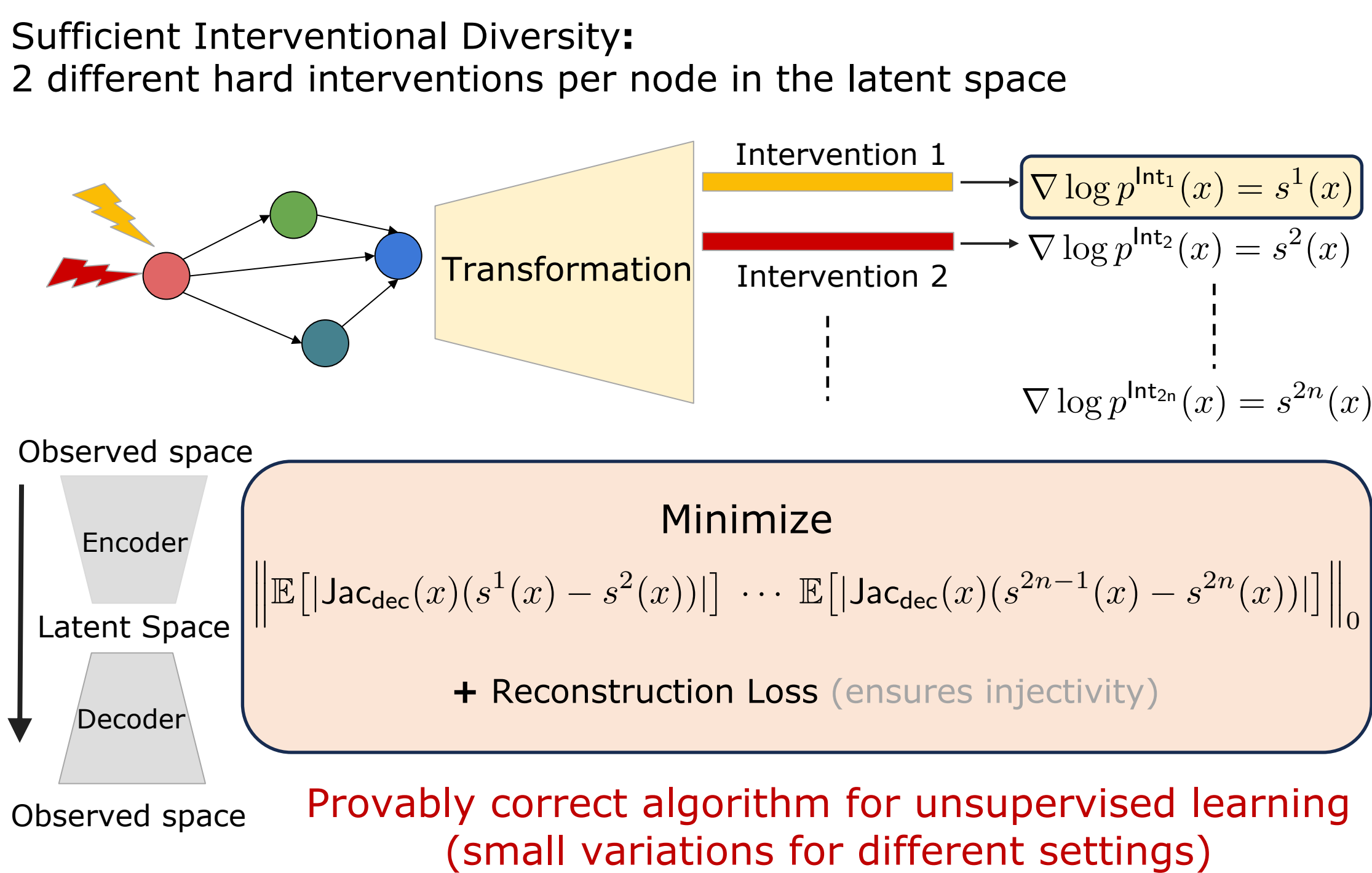


- Identifiability:** Conditions for uniquely recovering  $Z$  and  $g_Z$
- Achievability:** Provably correct algorithms to recover  $Z$  and  $g_Z$

## Our contributions

Related work for perfect ID	Transform	Requirements	Provable Algorithm
Varıcı et al. 2024	Linear	1 int/node	✓
von Kügelgen et al. 2023	General	2 (coupled) int/node + faithfulness	✗
<b>This work</b>	<b>General</b>	<b>2 (uncoupled) int/node</b>	<b>✓</b>

## 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

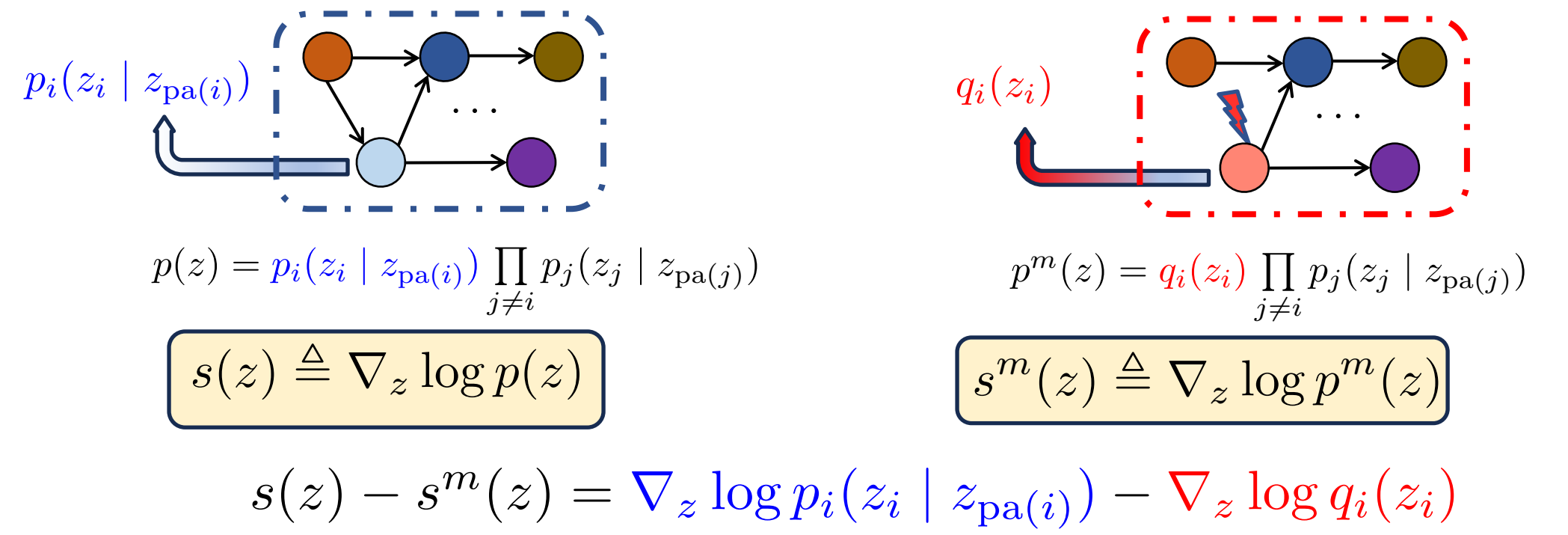
Input score differences ( $s_X - s_X^m$ ): Perfect score oracle or Sliced Score Matching

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

Obs. dim	Norm. Z error	DAG error (SHD)	Norm. Z error	DAG error (SHD)
8	0.16	1.56	0.70	11.9
25	0.20	1.55	0.68	10.5
40	0.21	1.14	0.71	11.8

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  $i$

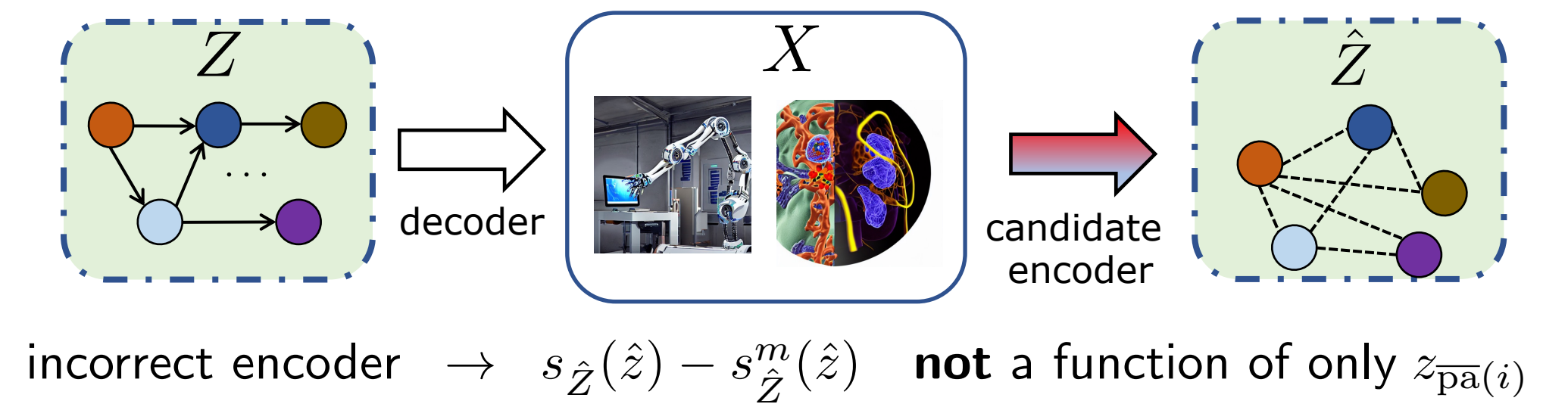
Coupled hard interventions have sparse score differences

node  $i$  intervened twice:  $s^m(z) - \tilde{s}^m(z)$  becomes a function of only  $z_i$

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

coordinate of node  $i$

## Methodology



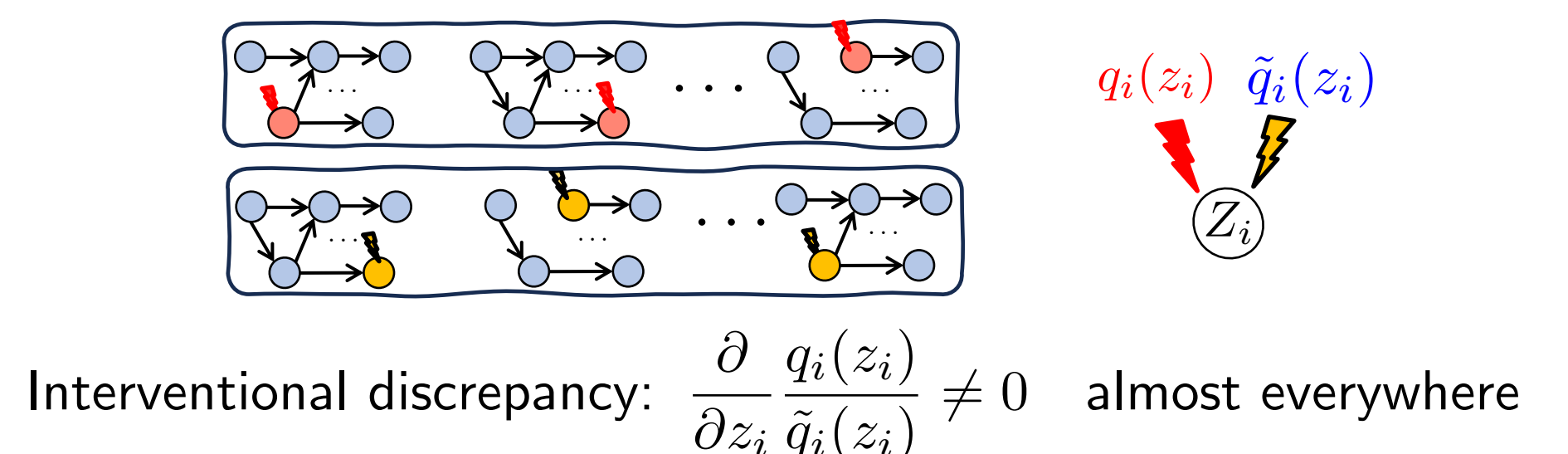
estimated score differences cannot be sparser than true score differences

**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



**Theorem :** Observational data and **two hard** interventions/node suffice for **perfect ID:**

- Latent graph recovery up to isomorphism
- Latent variables recovery up to elementwise transform

- Achievable algorithm
- No faithfulness assumption for identifiability
- Uncoupled two hard interventions per node

