

# Causal Diffusion Autoencoders: Toward Representation-Enabled Counterfactual Generation via Diffusion Probabilistic Models

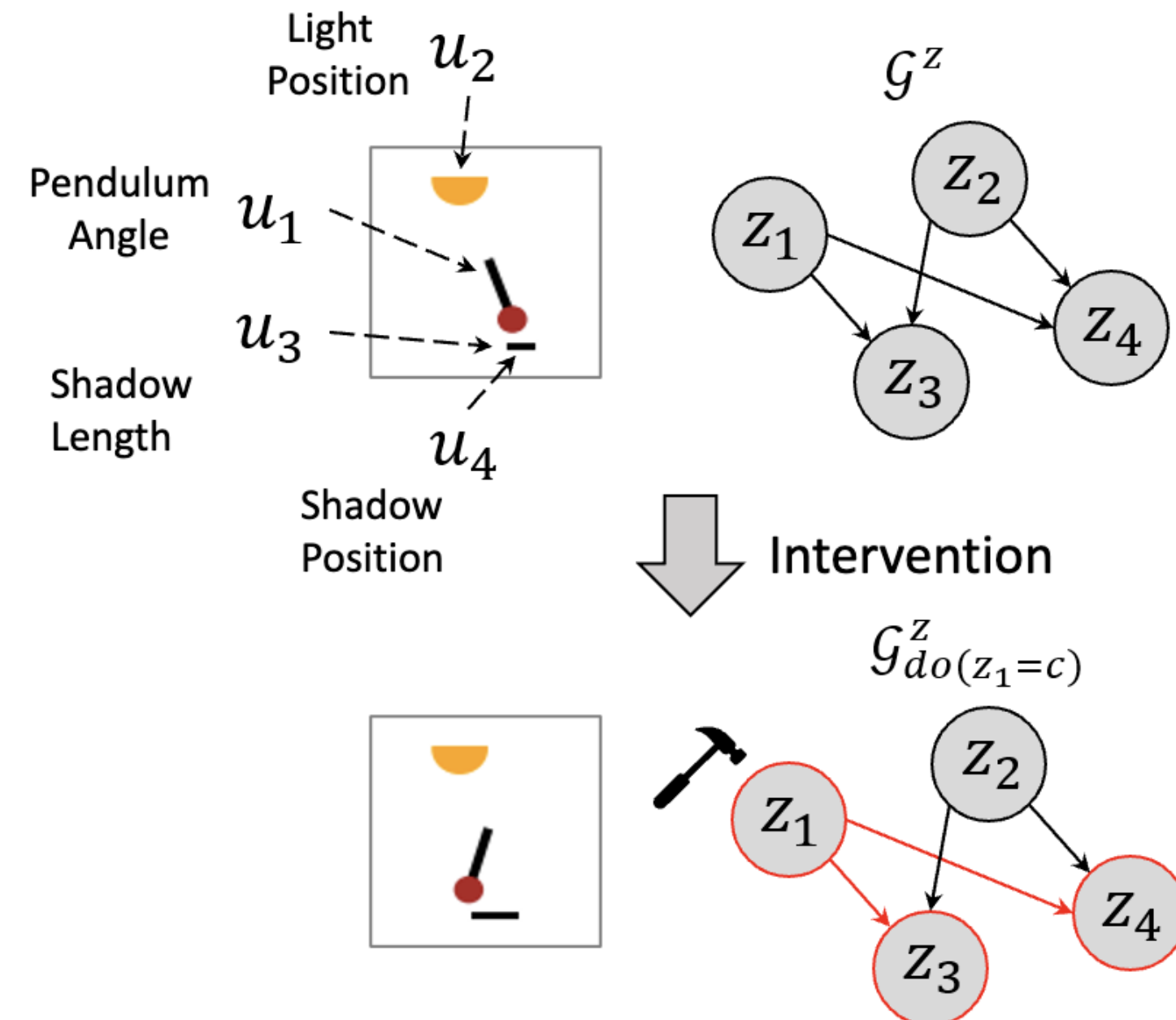
Aneesh Komanduri<sup>1</sup>, Chen Zhao<sup>2</sup>, Feng Chen<sup>3</sup>, and Xintao Wu<sup>1</sup>

<sup>1</sup>University of Arkansas, <sup>2</sup>Baylor University, <sup>3</sup>University of Texas at Dallas



## Motivation

- Diffusion models have become state-of-the-art in image generation, but do not model complex causal dependencies in latent space
- Representation-enabled **counterfactual generation** [1] is an underexplored area in diffusion models
- **Applications:** Healthcare, medicine, biology, etc.



## Contributions

- We propose **CausalDiffAE**, a causal representation learning diffusion-based framework to achieve counterfactual generation. We
- learn a causal representation via stochastic encoder and model causal mechanisms via neural networks in the latent space
- formulate a variational objective with an alignment prior to ensure disentangled representations
- propose DDIM for counterfactual generation subject to interventions
- explore a weak supervision scenario with limited label supervision to enable granular control over generated counterfactuals

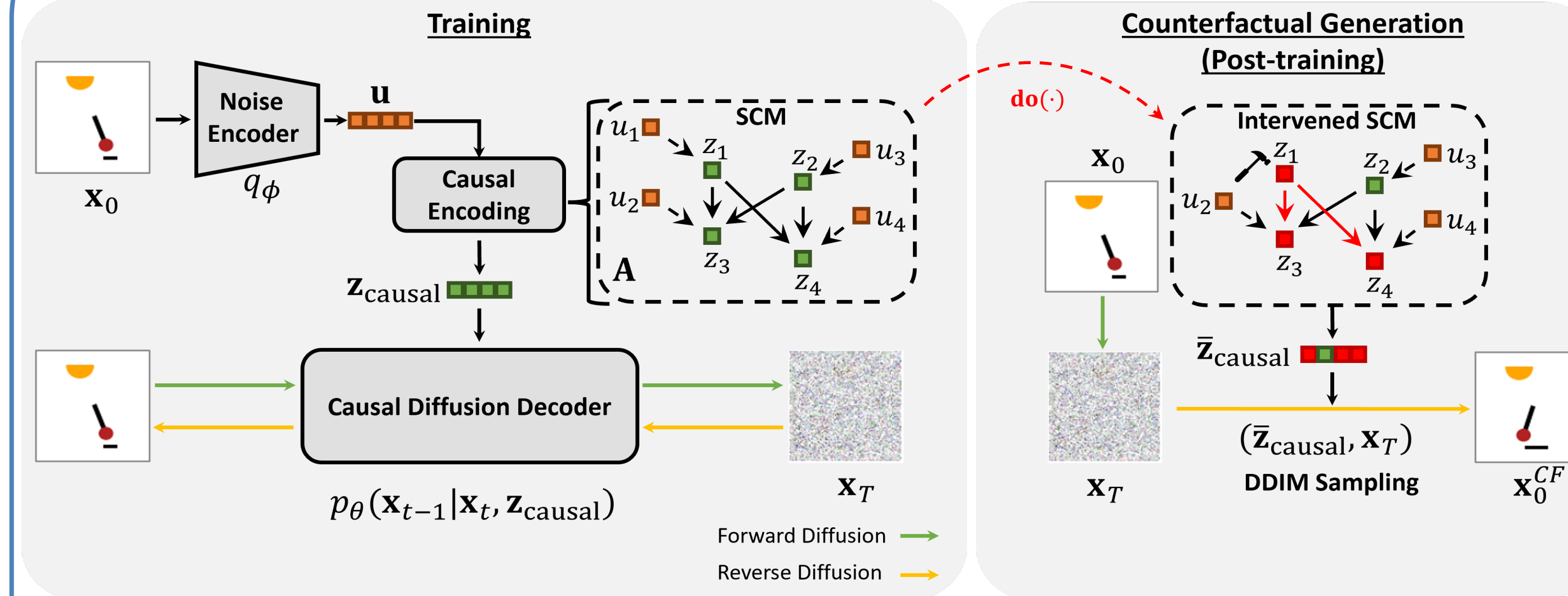
## Acknowledgements

This work is supported in part by NSF 1910284, 1946391, 2147375, and NIH P20GM139768

## References

[1] A. Komanduri et al. From Identifiable Causal Representations to Controllable Counterfactual Generation: A Survey on Causal Generative Modeling. TMLR. 2024.

## Causal Diffusion Autoencoder



$$\text{Loss} = \sum_{t=1}^T \mathbb{E}_{\mathbf{x}_0, \epsilon_t} [\|\epsilon_\theta(\mathbf{x}_t, t, \mathbf{z}_{\text{causal}}) - \epsilon_t\|_2^2] + \gamma \left\{ \mathcal{D}_{KL}(q_\phi(\mathbf{z}_{\text{causal}} | \mathbf{x}_0, \mathbf{y}) \| p(\mathbf{z}_{\text{causal}} | \mathbf{y})) + \mathcal{D}_{KL}(q_\phi(\mathbf{u} | \mathbf{x}_0) \| \mathcal{N}(0, \mathbf{I})) \right\}$$

$$\text{Causal Model} \quad \begin{aligned} z_i &= f_i(z_{\text{pa}_i}, u_i) \\ \mathbf{z} &= (\mathbf{I} - \mathbf{A}^T)^{-1} \mathbf{u} \\ z_i &= f_i(\mathbf{A}_i \odot \mathbf{z}; \nu_i) + u_i \end{aligned}$$

$$\text{Prior} \quad p(\mathbf{z}_{\text{causal}} | \mathbf{y}) = \prod_{i=1}^n p(z_i | y_i) = \prod_{i=1}^n \mathcal{N}(z_i; \mu_\nu(y_i), \sigma_\nu^2(y_i) \mathbf{I})$$

## Counterfactual DDIM

$$\mathbf{x}_{t-1}^{CF} = \sqrt{\alpha_{t-1}} \left( \frac{\mathbf{x}_t^{CF} - \sqrt{1 - \alpha_t} \epsilon_\theta(\mathbf{x}_t^{CF}, t, \mathbf{z}_{\text{causal}})}{\sqrt{\alpha_t}} \right) + \sqrt{1 - \alpha_{t-1}} \epsilon_\theta(\mathbf{x}_t^{CF}, t, \mathbf{z}_{\text{causal}})$$

## Counterfactual Generation Algorithm

**Input:** Factual sample  $\mathbf{x}_0$ , intervention target set  $\mathcal{I}$  with intervention values  $c$ , noise predictor  $\epsilon_\theta$ , encoder  $\phi$   
**Output:** Counterfactual sample  $\mathbf{x}_0^{CF}$

```

1:  $\mathbf{u} \sim q_\phi(\mathbf{u} | \mathbf{x}_0)$  ▷ Noise encoding
2: for  $i = 1$  to  $n$  do ▷ in topological order
3:   if  $i \in \mathcal{I}$  then
4:      $z_i = c_i$ 
5:   else
6:      $z_i = f_i(u_i, z_{\text{pa}_i})$ 
7:   end if
8: end for
9:  $\mathbf{z}_{\text{causal}} = \{z_1, \dots, z_n\}$  ▷ Intervened representation
10:  $\mathbf{x}_T \sim \mathcal{N}(\sqrt{\alpha_T} \mathbf{x}_0, (1 - \alpha_T) \mathbf{I})$ 
11:  $\mathbf{x}_T^{CF} = \mathbf{x}_T$ 
12: for  $t = T, \dots, 1$  do ▷ DDIM sampling
13:    $\mathbf{x}_{t-1}^{CF} = \sqrt{\alpha_{t-1}} \left( \frac{\mathbf{x}_t^{CF} - \sqrt{1 - \alpha_t} \epsilon_\theta(\mathbf{x}_t^{CF}, t, \mathbf{z}_{\text{causal}})}{\sqrt{\alpha_t}} \right) + \sqrt{1 - \alpha_{t-1}} \epsilon_\theta(\mathbf{x}_t^{CF}, t, \mathbf{z}_{\text{causal}})$ 
14: end for
15: return  $\mathbf{x}_0^{CF}$ 

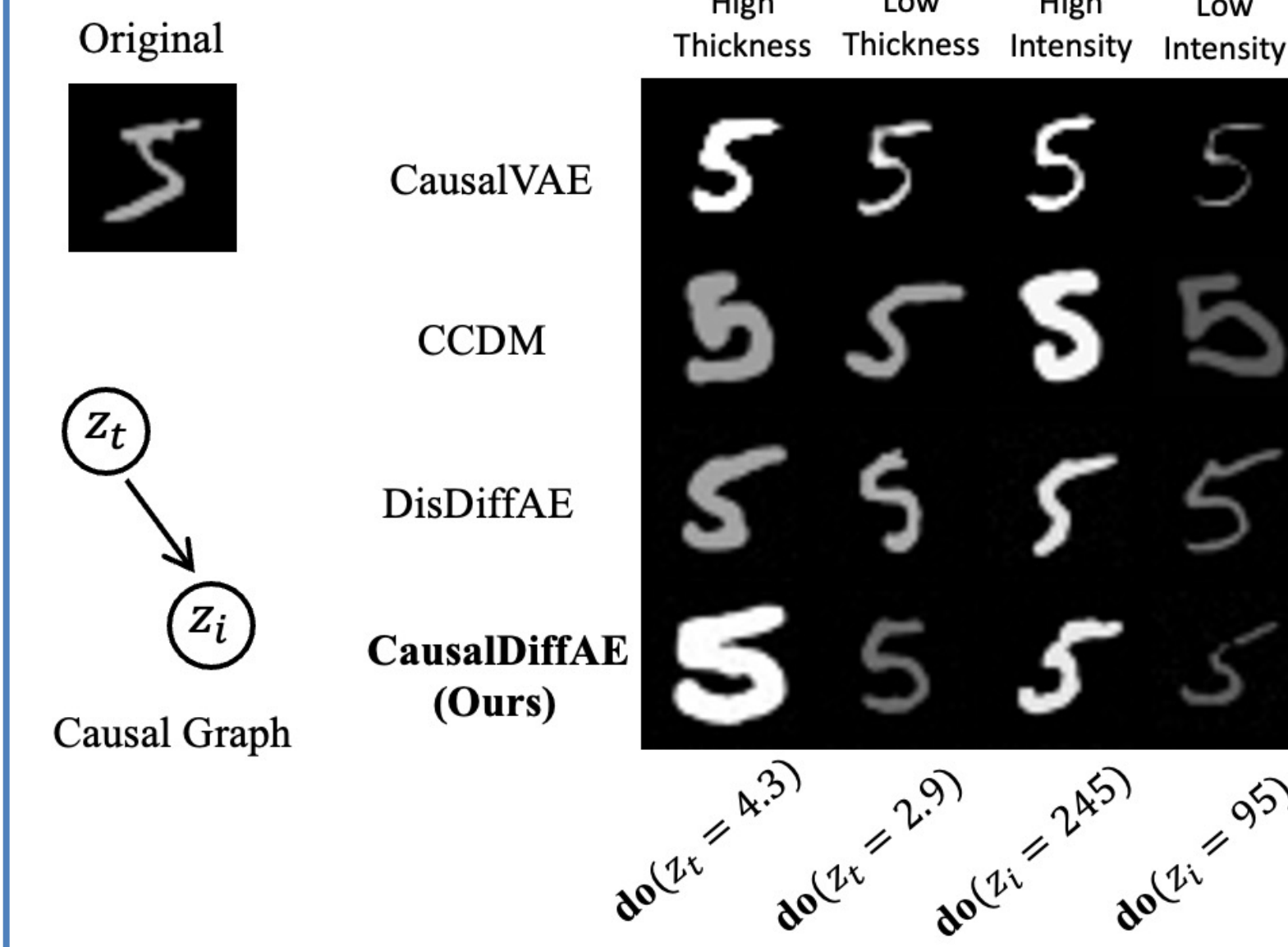
```

## General Procedure:

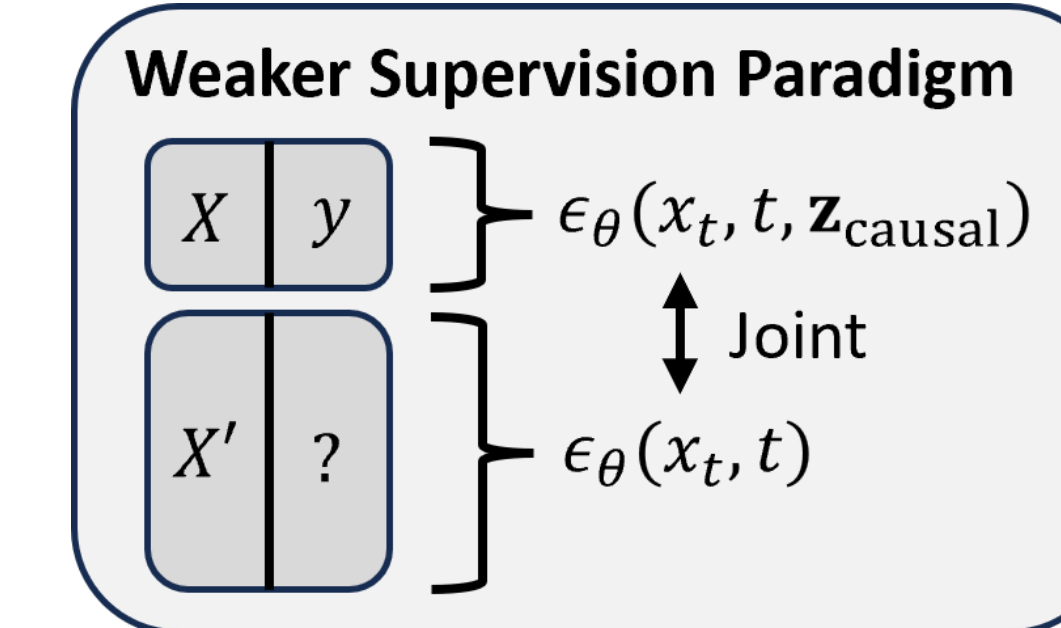
1. Noise abduction from factual
2. Intervene on causal latents
3. Propagate effects
4. Decode using DDIM

## Experimental Evaluation

- **Datasets:** MorphoMNIST and Pendulum, **Metrics:** DCI, Effectiveness, Visual Inspection
- **Generated Counterfactuals**



- **Weaker Supervision Training**



Limited label scenario  
 ➔ Jointly train conditional & unconditional model

- **Quantitative Metrics**

## High Disentanglement

Dataset	Model	DCI ↑
MorphoMNIST	CausalVAE	0.7838 <sub>0.01</sub>
	DiffAE	0.3578 <sub>0.01</sub>
	CausalDiffAE (Ours)	<b>0.9934<sub>0.01</sub></b>
Pendulum	CausalVAE	0.8850 <sub>0.01</sub>
	DiffAE	0.3525 <sub>0.01</sub>
	CausalDiffAE (Ours)	<b>0.9995<sub>0.01</sub></b>

## High Effectiveness

MorphoMNIST	Thickness MAE ↓		Intensity MAE ↓	
	do(t)	do(i)	do(t)	do(i)
CausalVAE	3.763 <sub>0.01</sub>	4.645 <sub>0.01</sub>	13.233 <sub>0.01</sub>	15.087 <sub>0.01</sub>
DisDiffAE	<b>0.377<sub>0.02</sub></b>	0.326 <sub>0.02</sub>	0.794 <sub>0.02</sub>	0.262 <sub>0.02</sub>
CausalDiffAE	0.392 <sub>0.02</sub>	<b>0.309<sub>0.02</sub></b>	<b>0.503<sub>0.01</sub></b>	<b>0.256<sub>0.01</sub></b>

Pendulum		Angle MAE ↓			
Model	do(a)	do(l)	do(sl)	do(sp)	
CausalVAE	24.860	23.030	20.470	11.580	
DisDiffAE	0.668	0.648	0.647	0.647	
CausalDiffAE	<b>0.297</b>	<b>0.132</b>	<b>0.031</b>	<b>0.034</b>	

\*for more results, see full paper