

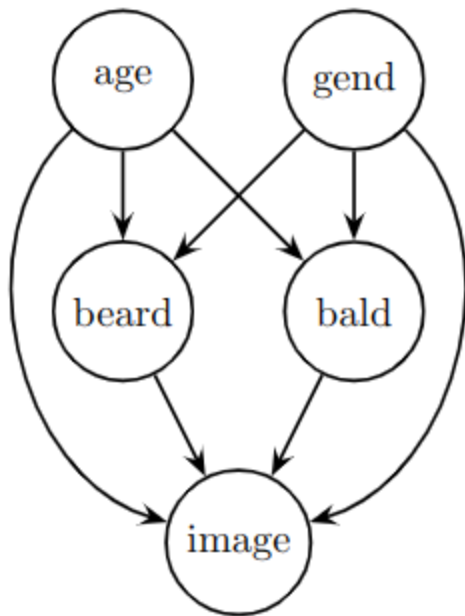
Benchmarking Counterfactual Image Generation

Thomas Melistas

Overview

- Background
 - What is Counterfactual Image Generation?
 - Methods and Models
- Evaluation Metrics
- Benchmarking Setup
- Results

Counterfactual Image Generation



(a) Causal graph



(b) Factual



(c) Causal



(d) Non-causal



Structural Causal Models

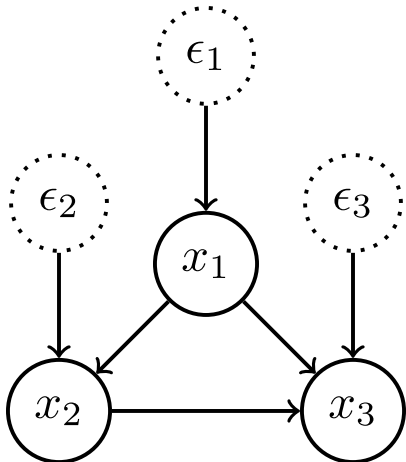
An SCM $\mathcal{G} := (\mathbf{S}, p(\boldsymbol{\epsilon}))$ consists of:

- (i) structural assignments $\mathbf{S} = \{f_i\}_{i=1}^N$, s.t. $x_i := f_i(\epsilon_i, \mathbf{pa}_i)$,
- (ii) a joint distribution $p(\boldsymbol{\epsilon}) = \prod_{i=1}^N p(\epsilon_i)$ over mutually independent noise variables

x_i : an **endogenous** variable (observed)

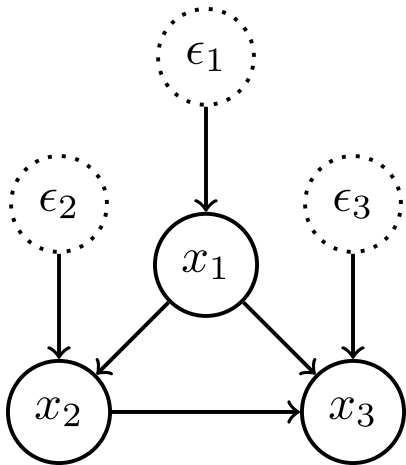
\mathbf{pa}_i : the parents of x_i (its direct *causes*, endogenous)

ϵ_i : an **exogenous** variable (unobserved)



Structural Causal Models

- Causal relations \rightarrow directed acyclic graph (DAG)
- Acyclic \rightarrow solve recursively for x_i and obtain $\mathbf{x} = \mathbf{f}(\boldsymbol{\epsilon})$
- \mathbf{x} : a collection of observable variables, where x_i : image and \mathbf{pa}_i : image attributes



Interventions and Counterfactuals

- *Interventional* distributions $P(x_j | do(x_i = y))$:
 - interventions: $x_i = f_i(\epsilon_i, \mathbf{pa}_i) \rightarrow x_i = y$
 - exogenous noise is sampled from the prior $P(\epsilon)$
- *Counterfactual* distributions $P(x_{j, x_i=y} | \mathbf{x})$:
 - interventions: as above
 - exogenous noise same with the observation $P(\epsilon | \mathbf{x})$

Interventions and Counterfactuals

Layer	Activity	Semantics	Example
(1) Associational $p(y x)$	Seeing 👁️👁️	How would seeing x change my belief in Y ?	What does a symptom tell us about the disease?
(2) Interventional $p(y \text{do}(x), z)$	Doing 💪	What happens to Y if I do x ?	What if I take aspirin, will my headache be cured?
(3) Counterfactual $p(y_{x'} x, y)$	Imagining 🤔	Was it x that caused Y ?	Was it the aspirin that stopped my headache?

Abduction-Action-Prediction

Counterfactuals using SCMs are computed in three steps:

(i) **Abduction**: Infer $P(\epsilon|\mathbf{x})$, the state of the world (exogenous noise) that is compatible with the observation \mathbf{x} .

(ii) **Action**: Replace the structural equations $do(x_i = y)$, resulting in a modified SCM $\tilde{\mathcal{G}} := \mathcal{G}_{\mathbf{x}; do(x_i=y)} = (\tilde{\mathbf{S}}, P(\epsilon|\mathbf{x}))$.

(iii) **Prediction**: Use the modified model to compute $P_{\tilde{\mathcal{G}}}(\mathbf{x})$.

Using Neural Networks for Abduction

Three categories of mechanisms:

- (i) **Invertible, explicit:** Conditional Normalising Flows \rightarrow attributes
- (ii) **Amortised, explicit:** Conditional VAEs or Hierarchical VAEs \rightarrow image
- (iii) **Amortised, implicit:** Conditional GANs \rightarrow image

Invertible, explicit

- Normalising flows perform mappings between probability densities
- A series of invertible transformations
- For an attribute x_i we utilise a conditional NF $f(\epsilon_i; \mathbf{pa}_i)$ which is invertible:
 $\epsilon_i = f^{-1}(x_i; \mathbf{pa}_i)$

$$P(x_i | \mathbf{pa}_i) = p(\epsilon_i) |\det \nabla_{\epsilon_i} f(\epsilon_i; \mathbf{pa}_i)|^{-1}$$

Amortised, explicit (with VAE)

- Encoder: $q_\phi(z|x, pa_x)$ and Decoder: $p_\theta(x|z, pa_x)$, trained with:

$$\text{ELBO}_\beta(\phi, \theta) = \mathbb{E}_{z \sim q_\phi(z|x, pa_x)} [p_\theta(x|z, pa_x)] - \beta D_{KL}[q_\phi(z|x, pa_x) || p(z)]$$

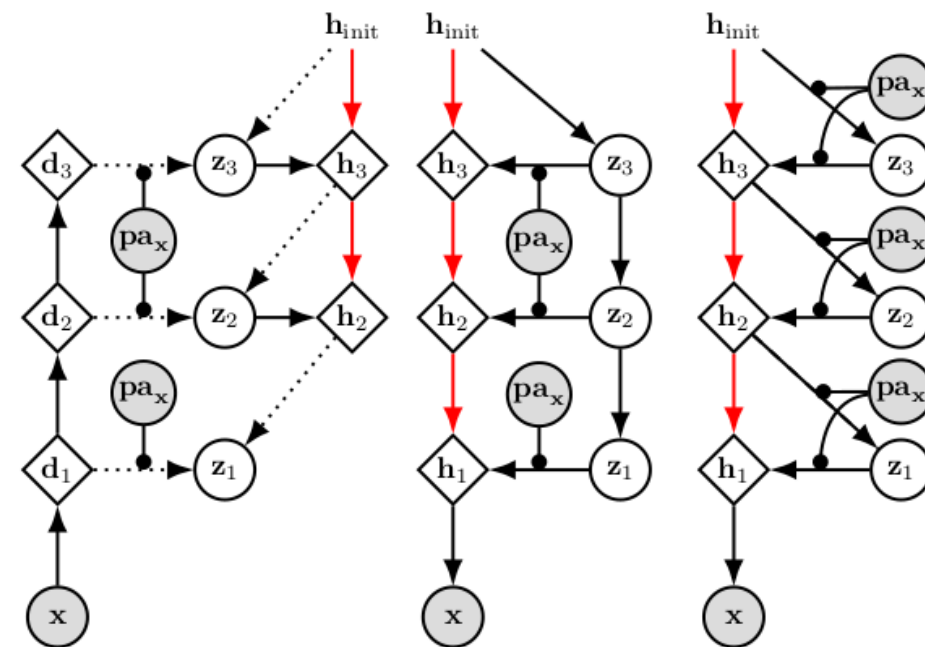
- The above likelihood and posterior are diagonal gaussians, whose μ and σ we predict with neural networks
- Prior $p(z) \sim N(0, I)$.

Amortised, explicit (with VAE)

- Noise ϵ is decomposed into $z \sim q_\phi(z|x, pa_x)$ and $u \sim N(0, I)$
- To perform counterfactual inference:
 - We sample the latent $z = \mu_\phi(x, pa_x) + \sigma_\phi(x, pa_x) * i, i \sim N(0, I)$
 - We sample the counterfactual $x^* = \mu_\theta(z, pa_x^*) + \sigma_\theta(z, pa_x^*) * u$,
where $u = \frac{x - \mu_\theta(z, pa_x)}{\sigma_\theta(z, pa_x)}$

Amortised, explicit (HVAE)

- We have L layers of hierarchical latent variables $\mathbf{z} = \{z_1, z_2, \dots, z_L\}$
- $h_i = h_{i+1} + f_{\omega}^i(z_i, pa_x)$
 $z_i \sim p_{\theta}(z_i | z_{>i})$,
 $p_{\theta}(x | z_{1:L}, pa_x) = N(x | \mu_{\theta}(h_1), \sigma_{\theta}(h_1))$
- $q_{\phi}(z_{1:L} | x, pa_x) = q_{\phi}(z_L | x, pa_x)$
 $q_{\phi}(z_{L-1} | z_L, x, pa_x) \dots q_{\phi}(z_1 | z_{>1}, x, pa_x)$



Amortised, implicit

- Encoder E : $z_x = E(x, pa_x)$
- Generator G : $x' = G(z', pa_x)$, where $z' \sim N(0, I)$ or $z' = z_x$
- Discriminator $D(x', z', pa_x)$: generated \rightarrow fake, data \rightarrow real

$$\min_{E, G} \max_D V(D, G, E) = \mathbb{E}_{q(x)p(pa_x)} [\log(D(x, E(x, pa_x), pa_x))] + \mathbb{E}_{p(z)p(pa_x)} [\log(1 - D(G(z, pa_x), z, pa_x))]$$

Finetuning the encoder:

$$L_x = \mathbb{E}_{x \sim q(x)} \|x - G(E(x, pa_x), pa_x)\|_2$$

$$L_z = \mathbb{E}_{z \sim p(z)} \|z - E(G(z, pa_x), pa_x)\|$$

To produce counterfactuals:

$$x^* = G(E(x, pa_x), pa_x^*)$$

Metrics

We use 4 metrics to evaluate the generated counterfactuals

- **Composition**
- **Effectiveness**
- **Realism (FID)**
- **Minimality (CLD)**

Composition

- Conceptually: the image and its attributes do not change without intervention
- *If we force a variable X to a value x it would have without the intervention, it should have no effect on the other variables*
- Therefore, a *null-intervention* applied m times: f_{\emptyset}^m should change no variable
- $\text{composition}^m(x, pa_x) = d(x, f_{\emptyset}^m(x, pa_x))$,
where $d(\cdot, \cdot)$ is a suitable distance metric

Effectiveness

- Conceptually: how successful was the performed intervention
- *If we force a variable X to have the value x , then X will take on the value x*
- We train an anti-causal predictor g_{θ}^i on observations, for each pa_x^i
- $\text{effectiveness}_i(x^*, pa_x^{*i}) = d(g_{\theta}^i(x^*), pa_x^{*i})$,
where $d(\cdot, \cdot)$: classification metric for categorical variables, regression for continuous

Fréchet Inception Distance (FID)

- Conceptually: Counterfactual image quality
- Measures the similarity of counterfactual images to observational data
- We use Inception-v3 trained on ImageNet to extract features
- Defined as:

$$d^2((m_q, C_q), (m_p, C_p)) = ||m_q - m_p||_2^2 + \text{Tr}(C_q + C_p - 2(C_q C_p)^{\frac{1}{2}})$$

where q : counterfactual features distribution, p : factual features distribution Fréchet distance d is defined as the distance between the Gaussian with mean (m_q, C_q) obtained from q and the Gaussian with mean (m_p, C_p) obtained from p .

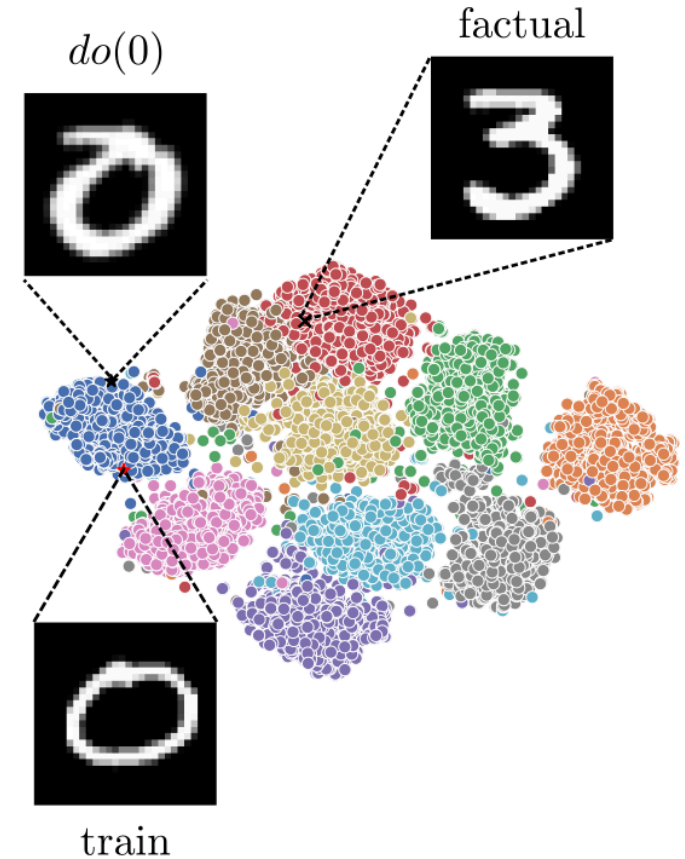
Counterfactual Latent Divergence (CLD)

- Conceptually: Counterfactual only differs in the intervened parent attribute
- Formally:

$$\text{CLD} = \log(w_1 \exp P(S_{x^*} \leq \text{div}) + w_2 \exp P(S_x \geq \text{div}))$$

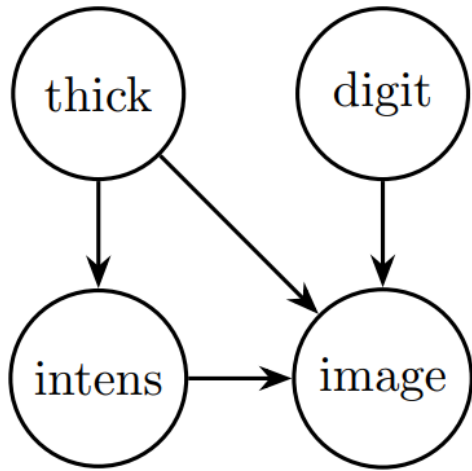
,where: $\text{div} = d(x, x^*)$, $S_x = \{d(x, x') | pa_{x'} = pa_x\}$,
 $S_{x^*} = \{d(x, x') | pa_{x'} = pa_{x^*}\}$

- $d(\cdot, \cdot)$: KL-divergence between the latents given by an unconditional VAE

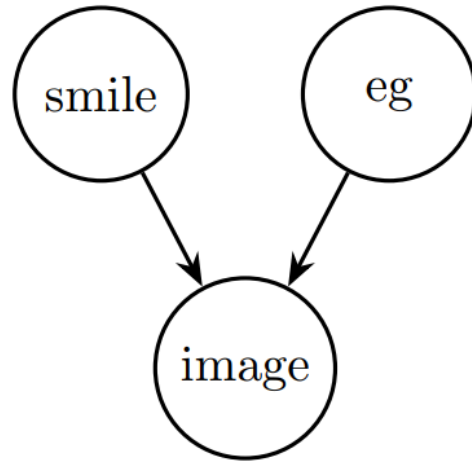


Datasets used for benchmarking

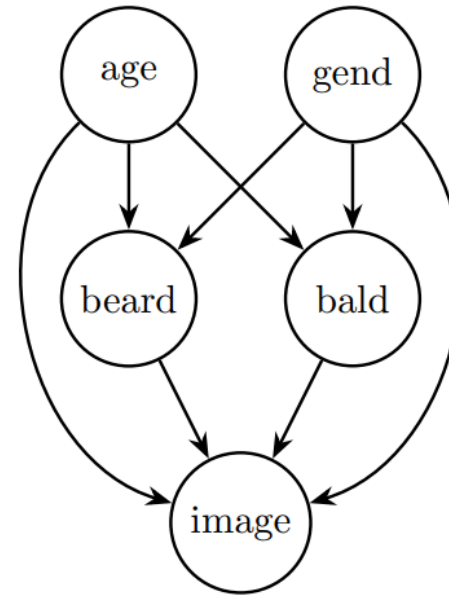
- **MorphoMNIST**: Synthetic dataset of digits (32×32)
- **CelebA**: Human faces (64×64) \rightarrow simple & complex causal graph
- **ADNI**: 2D slices of brain MRIs (192×192)



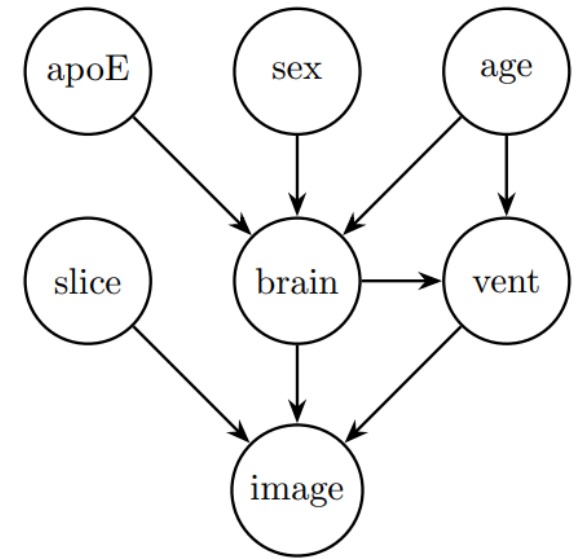
(a) MorphoMNIST



(b) CelebA (simple)

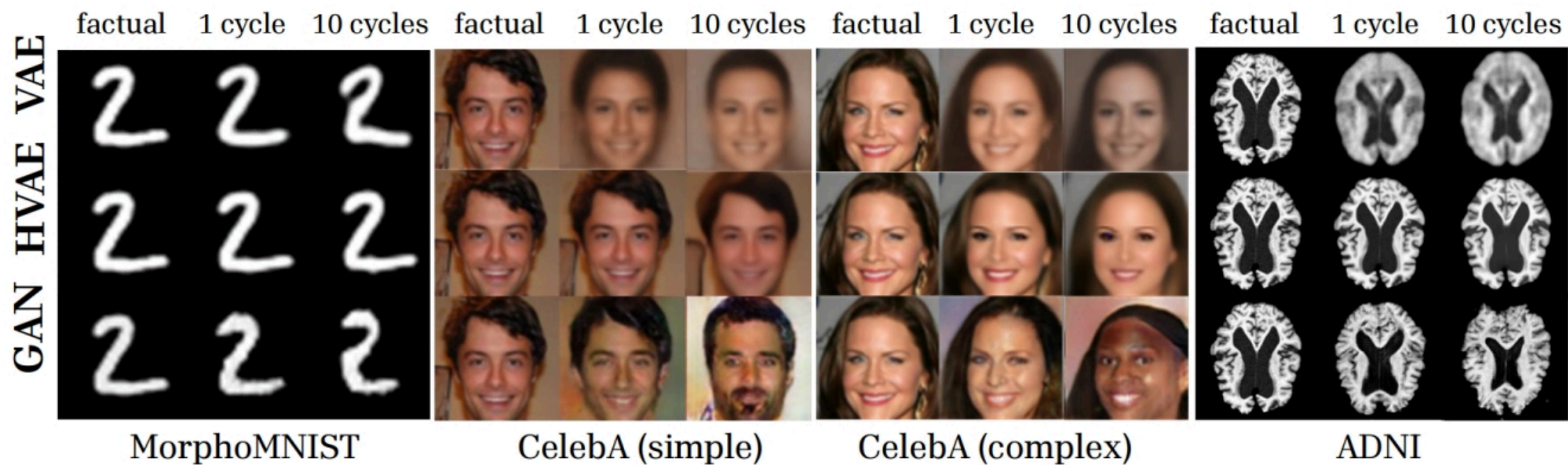


(c) CelebA (complex)



(d) ADNI

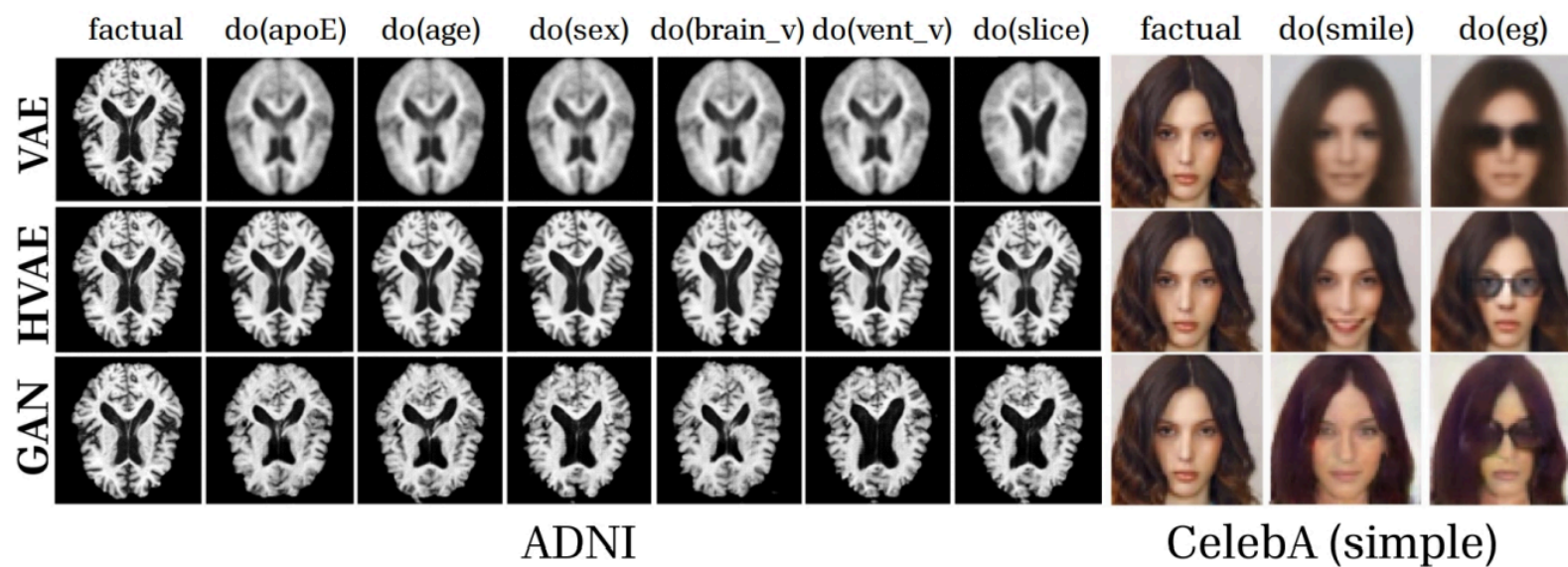
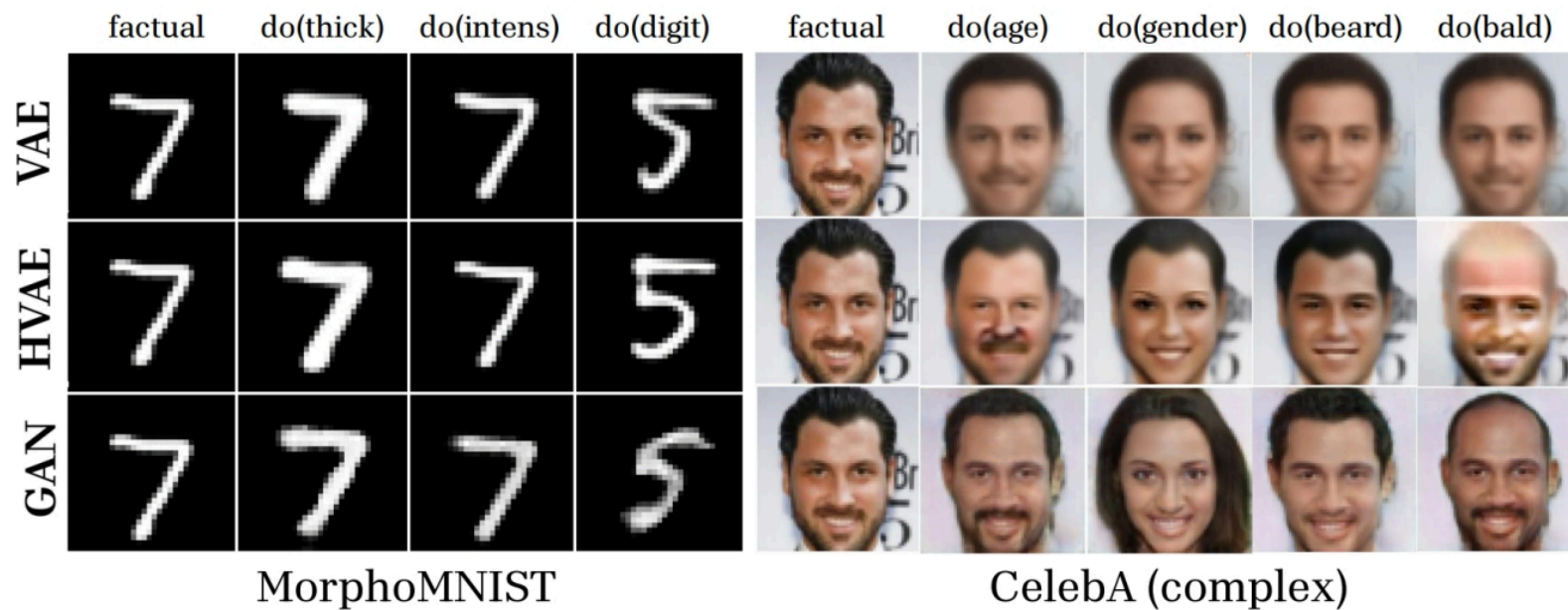
Composition



Composition

Model	l_1 image space ↓		LPIPS ↓	
	1 cycle	10 cycles	1 cycle	10 cycles
MorphoMNIST				
VAE	2.600 _{1.454}	7.698 _{4.199}	0.026 _{0.003}	0.075 _{0.003}
HVAE	0.438 _{0.178}	1.550 _{0.541}	0.006 _{0.001}	0.024 _{0.003}
GAN	3.807 _{1.822}	11.697 _{5.750}	0.049 _{0.004}	0.184 _{0.008}
CelebA (simple/complex)				
VAE	127.8 _{18.1} / 121.2 _{13.2}	123.4 _{22.1} /127.9 _{21.1}	0.295 _{0.004} /0.282 _{0.061}	0.424 _{0.005} /0.412 _{0.091}
HVAE	129.4 _{11.6} /122.7 _{10.4}	138.8 _{23.4} / 124.6 _{16.4}	0.063 _{0.003} / 0.122 _{0.033}	0.200 _{0.008} / 0.240 _{0.053}
GAN	115.3 _{22.0} /127.6 _{17.3}	128.4 _{21.4} /131.8 _{23.3}	0.290 _{0.003} /0.276 _{0.074}	0.462 _{0.005} /0.490 _{0.121}
ADNI				
VAE	18.882 _{1.786}	30.250 _{3.389}	0.306 _{0.008}	0.384 _{0.006}
HVAE	3.384 _{0.367}	7.456 _{0.622}	0.101 _{0.012}	0.156 _{0.014}
GAN	24.261 _{1.821}	32.794 _{3.578}	0.268 _{0.009}	0.323 _{0.007}

Qualitative Results



Effectiveness

Model	Thickness (t) MAE ↓			Intensity (i) MAE ↓			Digit (y) Acc. ↑		
	$do(t)$	$do(i)$	$do(y)$	$do(t)$	$do(i)$	$do(y)$	$do(t)$	$do(i)$	$do(y)$
VAE	0.109 _{0.01}	0.333 _{0.02}	0.139 _{0.01}	3.15 _{0.26}	5.33 _{0.29}	4.64 _{0.35}	0.989 _{0.01}	0.988 _{0.01}	0.775 _{0.02}
HVAE	0.086 _{0.09}	0.224 _{0.03}	0.117 _{0.01}	1.99 _{0.18}	3.52 _{0.26}	2.10 _{0.17}	0.985 _{0.01}	0.935 _{0.02}	0.972 _{0.02}
GAN	0.228 _{0.01}	0.680 _{0.02}	0.393 _{0.02}	9.43 _{0.59}	15.14 _{0.99}	12.39 _{0.57}	0.961 _{0.02}	0.966 _{0.01}	0.451 _{0.024}

CelebA (simple)								
Model	Smiling (s) F1 ↑				Eyeglasses (e) F1 ↑			
	$do(s)$	$do(e)$			$do(s)$	$do(e)$		
VAE	0.897 _{0.02}	0.987 _{0.01}			0.938 _{0.05}	0.810 _{0.02}		
HVAE	0.998 _{0.01}	0.997 _{0.01}			0.883 _{0.06}	0.981 _{0.02}		
GAN	0.819 _{0.02}	0.873 _{0.01}			0.957 _{0.03}	0.891 _{0.01}		

CelebA (complex)								
	Age (a) F1 ↑				Gender (g) F1 ↑			
	$do(a)$	$do(g)$	$do(br)$	$do(bl)$	$do(a)$	$do(g)$	$do(br)$	$do(bl)$
VAE	0.35 _{0.04}	0.782 _{0.02}	0.816 _{0.02}	0.819 _{0.02}	0.977 _{0.01}	0.909 _{0.02}	0.959 _{0.02}	0.973 _{0.01}
HVAE	0.654 _{0.1}	0.893 _{0.04}	0.908 _{0.03}	0.899 _{0.03}	0.988 _{0.02}	0.949 _{0.03}	0.994 _{0.01}	0.95 _{0.03}
GAN	0.413 _{0.04}	0.71 _{0.02}	0.818 _{0.02}	0.799 _{0.01}	0.952 _{0.01}	0.982 _{0.01}	0.92 _{0.01}	0.961 _{0.01}
	Beard (br) F1 ↑				Bald (bl) F1 ↑			
	$do(a)$	$do(g)$	$do(br)$	$do(bl)$	$do(a)$	$do(g)$	$do(br)$	$do(bl)$
VAE	0.944 _{0.01}	0.828 _{0.03}	0.296 _{0.05}	0.945 _{0.02}	0.023 _{0.03}	0.496 _{0.05}	0.045 _{0.04}	0.412 _{0.03}
HVAE	0.952 _{0.03}	0.951 _{0.03}	0.441 _{0.11}	0.916 _{0.04}	0.02 _{0.05}	0.86 _{0.05}	0.045 _{0.07}	0.611 _{0.04}
GAN	0.908 _{0.01}	0.838 _{0.02}	0.233 _{0.03}	0.907 _{0.01}	0.021 _{0.02}	0.82 _{0.02}	0.055 _{0.02}	0.492 _{0.02}

Model	Brain volume (b) MAE ↓			Ventricular volume (v) MAE ↓			Slice (s) F1 ↑		
	$do(b)$	$do(v)$	$do(s)$	$do(b)$	$do(v)$	$do(s)$	$do(b)$	$do(v)$	$do(s)$
VAE	0.17 _{0.03}	0.15 _{0.06}	0.15 _{0.06}	0.08 _{0.05}	0.20 _{0.04}	0.08 _{0.05}	0.52 _{0.15}	0.48 _{0.15}	0.46 _{0.10}
HVAE	0.09 _{0.03}	0.12 _{0.06}	0.13 _{0.06}	0.06 _{0.04}	0.04 _{0.01}	0.06 _{0.04}	0.38 _{0.15}	0.41 _{0.16}	0.41 _{0.11}
GAN	0.17 _{0.02}	0.16 _{0.07}	0.16 _{0.06}	0.12 _{0.02}	0.22 _{0.03}	0.12 _{0.03}	0.14 _{0.03}	0.16 _{0.03}	0.05 _{0.02}

Realism (FID) & Minimality (CLD)

Model	MorphoMNIST		CelebA (simple/complex)		ADNI	
	FID ↓	CLD ↓	FID ↓	CLD ↓	FID ↓	CLD ↓
VAE	10.124	0.268	66.412/59.393	0.301/ 0.299	278.245	0.352
HVAE	5.362	0.272	22.047 /35.712	0.295 /0.305	74.696	0.347
GAN	35.568	0.286	31.560/ 27.861	0.38/0.304	113.749	0.353

Conclusions

- HVAE outperforms other models across metrics and datasets
- GAN counterfactuals more realistic than VAE but far from factials for complex datasets
- Amortised implicit mechanisms → better abduction, but hierarchical latents are important

Future Work

- Extend to other
 - Generative models for the image mechanism (i.e. Diffusion Models)
 - Counterfactual paradigms (i.e. Deep Twin Networks, Backtracking Counterfactuals)
- Metrics
 - Limit bias of used model-dependent metrics (i.e. effectiveness, CLD)
 - Come up with new metrics

Sanchez & Tsafaris: Diffusion Causal Models for Counterfactual Estimation. 2022

Vlontzos et al. Estimating categorical counterfactuals via deep twin networks. 2023

Kladny et al. Deep backtracking counterfactuals for causally compliant explanations. 2024

Thank you for your attention!!!