

# Causality & Generative Models



Fabio De Sousa Ribeiro

fdesousa@ic.ac.uk



Miguel Monteiro



Nick Pawlowski



Daniel C. Castro



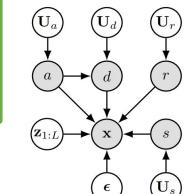
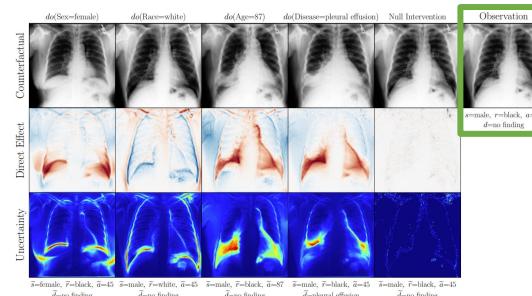
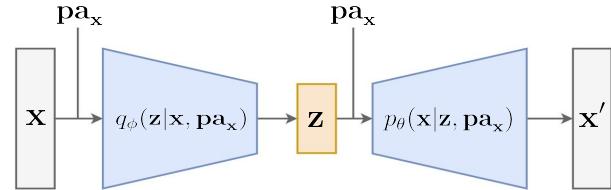
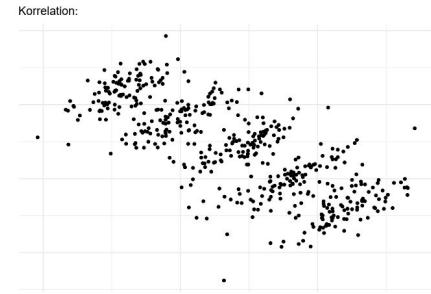
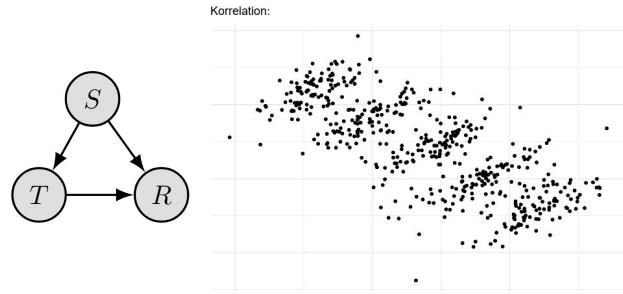
Tian Xia



Ben Glocker

# Outline

- Introduction & Motivation
  - ◆ Causality: The Ladder of Causation
  - ◆ Structural Causal Models
- Image Counterfactuals
  - ◆ Deep Causal Mechanisms for Structured Variables
  - ◆ Causal Mediation Analysis
- Case Studies: Medical Imaging
  - ◆ Brain Imaging
  - ◆ Chest X-ray Imaging
- Conclusion & Outlook



(a) Deep SCM for MIMIC-CXR. The variables in the causal graph are: age ( $a$ ), sex ( $s$ ), race ( $r$ ), disease ( $d$ ) and chest x-ray ( $x$ ). The disease  $d$  is pleural effusion.

# What is Causality?

- Scientific inquiry is often motivated by causal questions:
  - I. How effective is a **treatment** in preventing a **disease**?
  - II. If I take this **pill**, will my **headache** be gone?
  - III. Would my **grades** have been better had I **studied** more?

Causality is the relationship between **cause** and **effect**

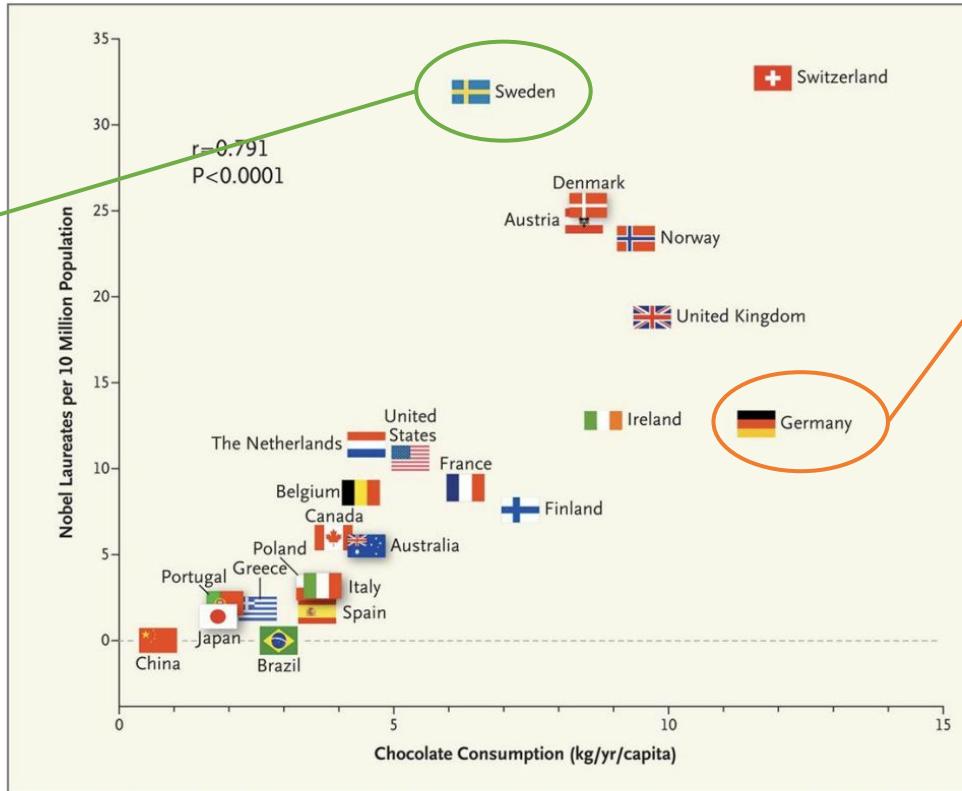
# Example: chocolate

“Correlation does not imply causation.”

Nobel Prize Correlates with Chocolate Consumption

Secretly the best chocolate?

Least effective chocolate?



# Reichenbach's Common Cause Principle



Hans Reichenbach

- I. Chocolate consumption causes Nobel prize win

$$C \rightarrow N$$

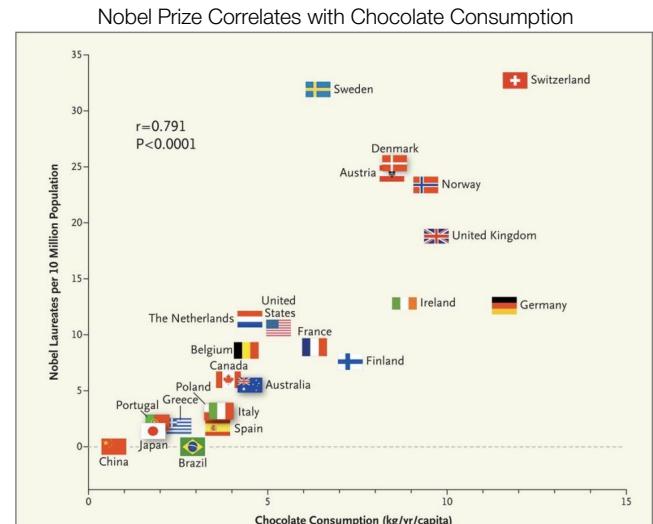
- II. Nobel prize win causes chocolate consumption

$$C \leftarrow N$$

- III. Both are caused by an unknown factor  $U$

$$C \leftarrow U \rightarrow N$$

Could  $U$  be GDP per capita or wealth per adult?



F. H. Messerli: Chocolate Consumption, Cognitive Function, and Nobel Laureates, N Engl J Med 2010

# Simpson's Paradox



Edward Simpson

Kidney Stone Size	Treatment A	Treatment B
Small	<b>93%</b> (81/87)	87% (234/270)
Large	<b>73%</b> (192/263)	69% (55/80)
Both	78% (273/350)	<b>83%</b> (289/350)

# Simpson's Paradox



Edward Simpson

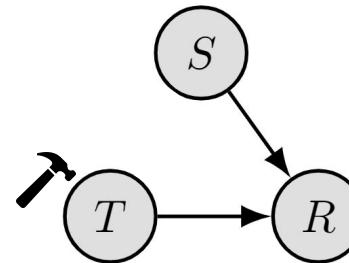
Kidney Stone Size	Treatment A	Treatment B
Small	<b>93% (81/87)</b>	87% (234/270)
Large	<b>73% (192/263)</b>	69% (55/80)
Both	78% (273/350)	<b>83% (289/350)</b>

# Simpson's Paradox



Edward Simpson

- Stone size ( $S$ ) is a **confounder**
- Make treatment ( $T$ ) independent of size:
  - ◆ What's the recovery ( $R$ ) rate if **all** subjects receive treatment **A** vs **B**?



Kidney Stone Size	Treatment A	Treatment B
Small	93% (81/87)	87% (234/270)
Large	73% (192/263)	69% (55/80)
Both	78% (273/350)	83% (289/350)

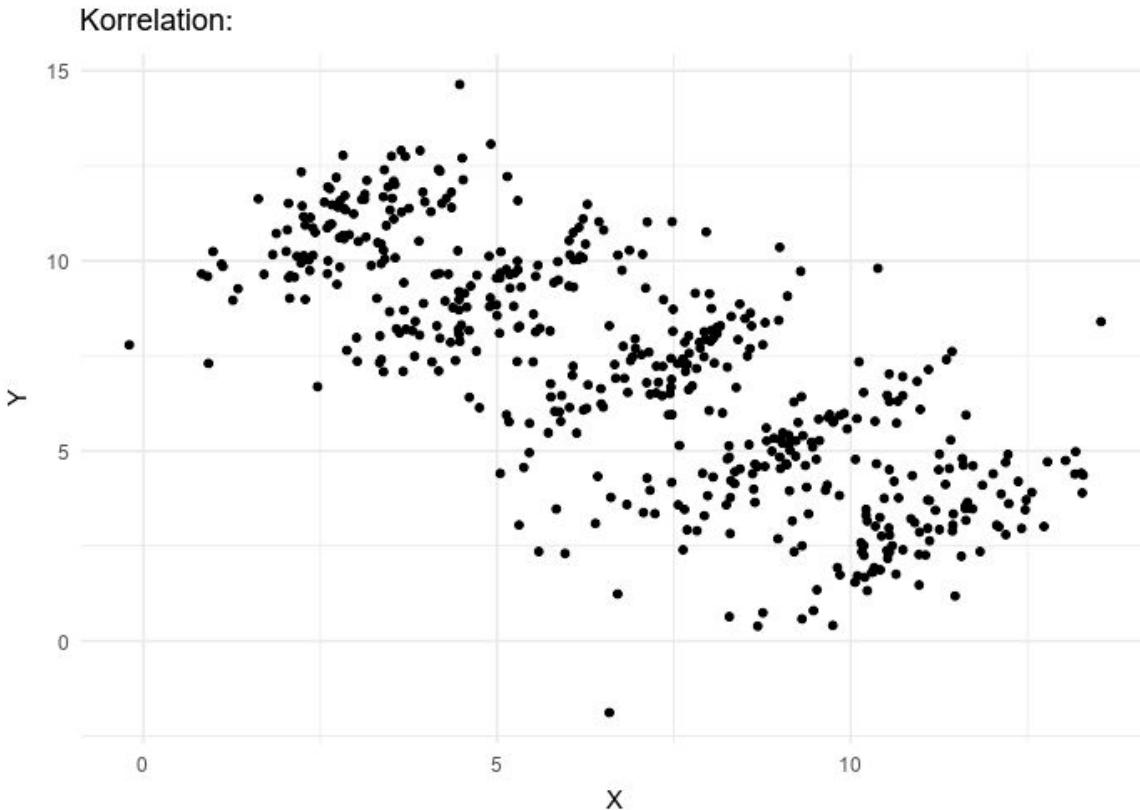
??

Patients with **large stones** received the **better treatment (A)**, and those with **small stones** received the **inferior treatment (B)**

# Simpson's Paradox



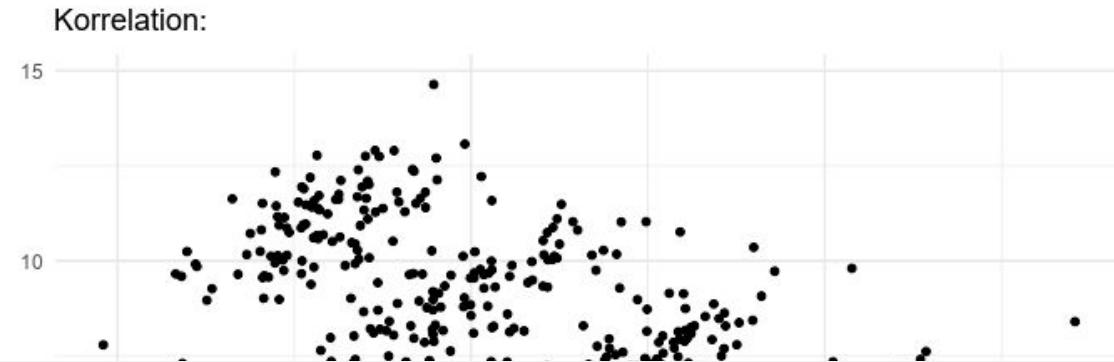
Edward Simpson



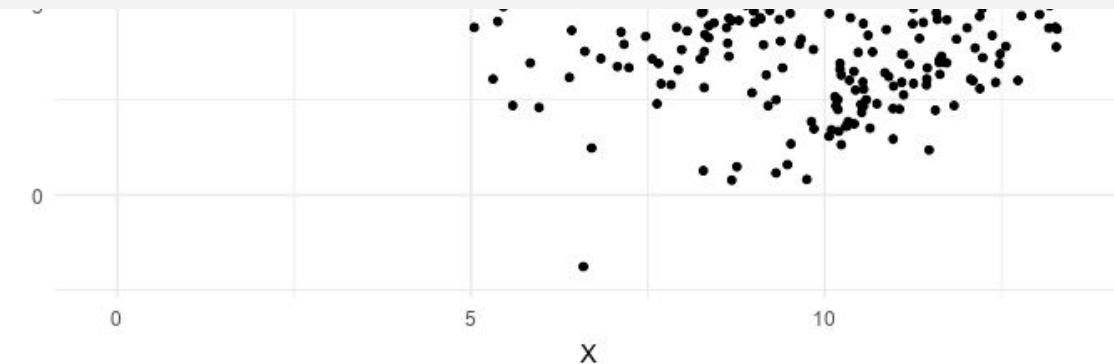
# Simpson's Paradox



Edward Simpson



Correlations may **reverse** depending on how we aggregate or filter data and its subpopulations



# Predictive Modelling

Given an image  $X$ , train a model to predict some label  $Y$

$$P(Y|X)$$

- Assumptions:
  - ◆ Sufficient training data ( $X, Y$ ) is available
  - ◆ Training and test data come from the same distribution

# A Causal Perspective: Predictive Modelling

What is the causal relationship between image  $X$  and label  $Y$ ?

$$P(Y|X)$$

$$X \rightarrow Y$$

**causal**  
(predict effect from cause)

$$Y \rightarrow X$$

**anti-causal**  
(predict cause from effect)

# Example: Skin Lesion Classification

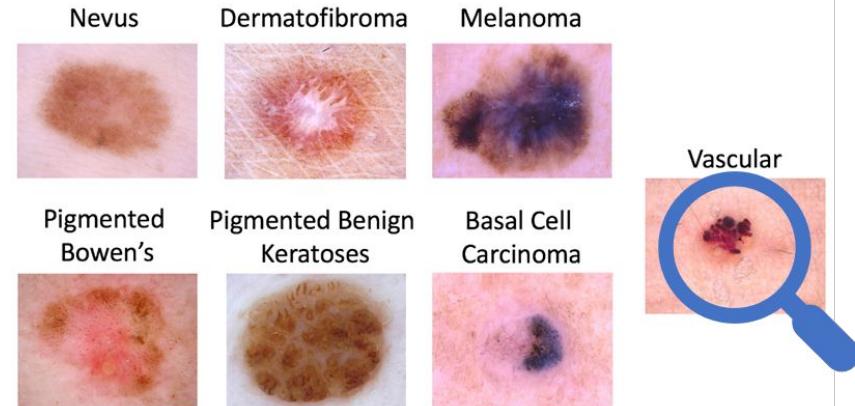
$X$  – dermascopic image

$Y$  – biopsy-derived diagnosis

$$X \xleftrightarrow{?} Y$$

causal or anti-causal?

$$P(Y|X)$$



# Example: Skin Lesion Classification

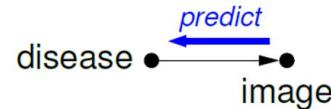
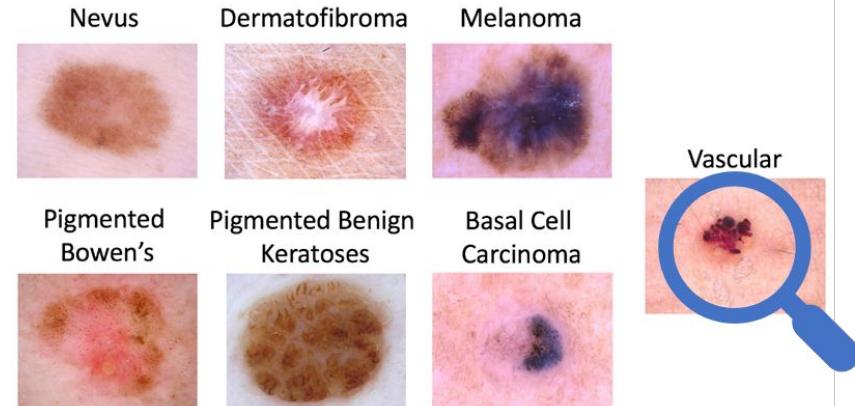
$X$  – dermascopic image

$Y$  – biopsy-derived diagnosis

$$X \leftarrow Y$$

**anti-causal**  
(predict **cause** from **effect**)

$$P(Y|X)$$



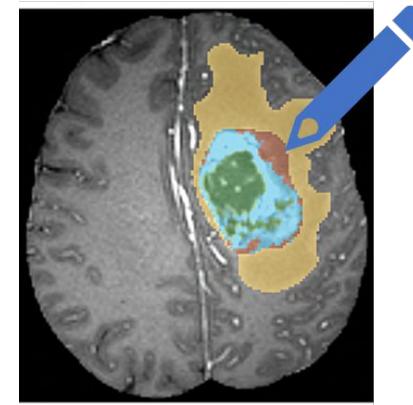
# Example: Brain Tumour Segmentation

$X$  – structural brain MRI

$Y$  – manually drawn contour

$$X \xleftrightarrow{?} Y$$

causal or anti-causal?



$$P(Y|X)$$

# Example: Brain Tumour Segmentation

$X$  – structural brain MRI

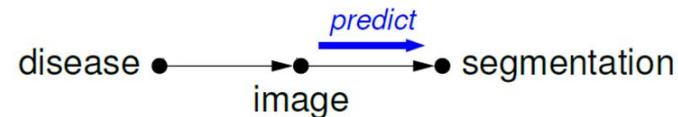
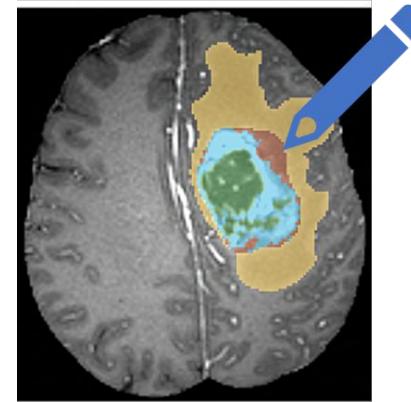
$Y$  – manually drawn contour

$$X \rightarrow Y$$

**causal**

(predict **effect** from **cause**)

$$P(Y|X)$$



# Example: Radiology Reports

$X$  – chest X-ray

$Y$  – diagnosis extracted from report

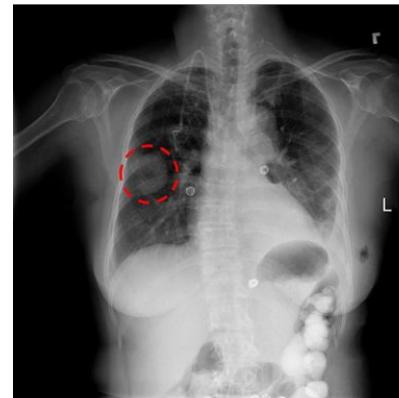
$$X \xleftrightarrow{?} Y$$

causal or anti-causal?

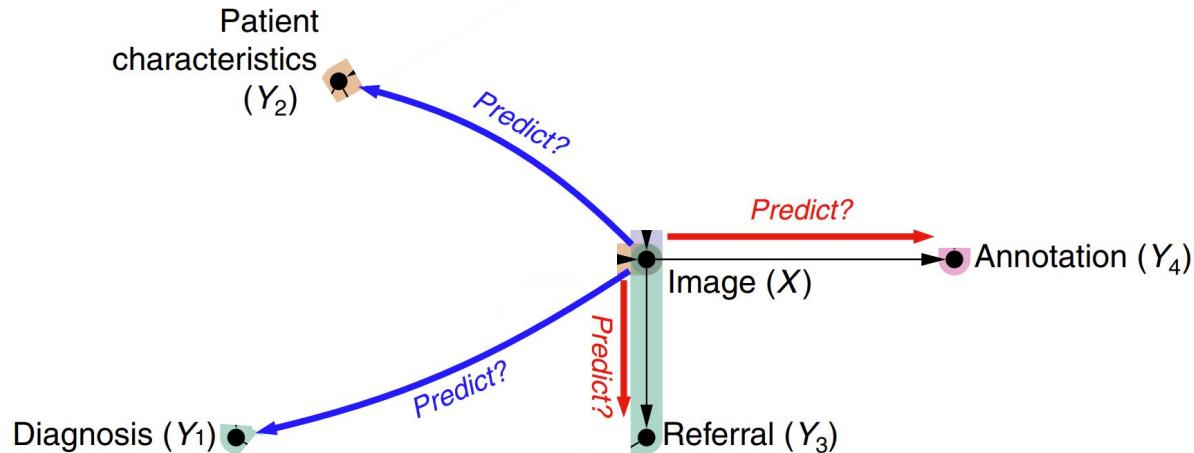
Not clear whether the diagnosis was derived from the image, or from another lab result...

$$P(Y|X)$$

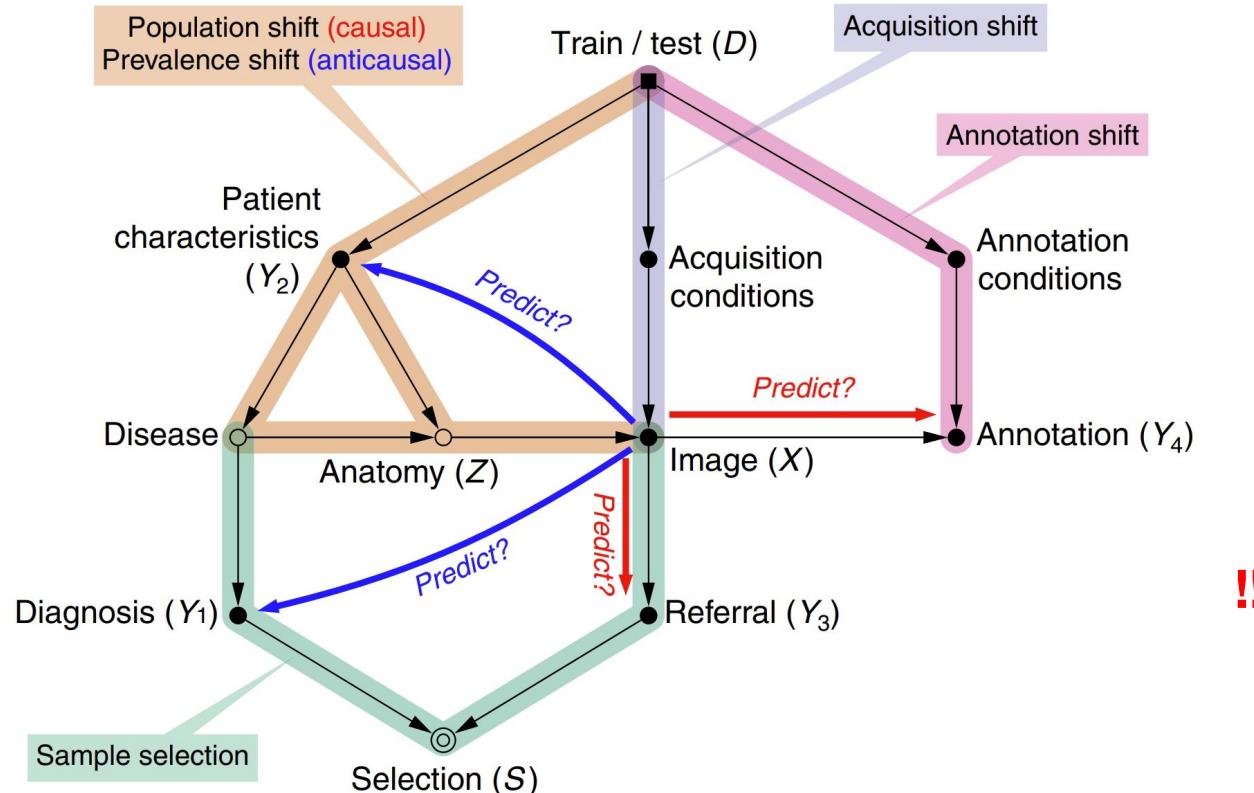
Meta-information is needed to establish causal relationships



# Data Generating Process



# Data Generating Process



!!

# Ladder of Causation



## Counterfactuals

**What if  $X$  had not occurred?**

- I. Counterfactual reasoning
- II. Deduce causes for observed events

3.

## Intervention

**What will  $Y$  be if i **do**  $X$ ?**

- I. Use experiments to identify causal effects
- II. Crucial for planning and policy making

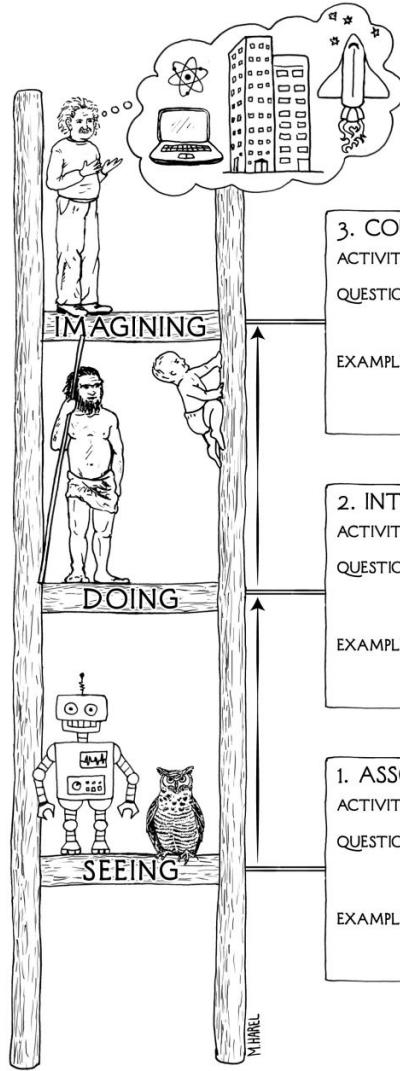
2.

## Association

**How would **seeing**  $X$  change my belief in  $Y$ ?**

- I. Modelling correlations  $P(Y|X)$
- II. Modern ML excels at this task

1.



### 3. COUNTERFACTUALS

ACTIVITY: Imagining, Retrospection, Understanding

QUESTIONS: *What if I had done ...? Why?*  
(Was it  $X$  that caused  $Y$ ? What if  $X$  had not occurred? What if I had acted differently?)

EXAMPLES: Was it the aspirin that stopped my headache?  
Would Kennedy be alive if Oswald had not killed him? What if I had not smoked for the last 2 years?

### 2. INTERVENTION

ACTIVITY: Doing, Intervening

QUESTIONS: *What if I do ...? How?*  
(What would  $Y$  be if I do  $X$ ?  
How can I make  $Y$  happen?)

EXAMPLES: If I take aspirin, will my headache be cured?  
What if we ban cigarettes?

### 1. ASSOCIATION

ACTIVITY: Seeing, Observing

QUESTIONS: *What if I see ...?*  
(How are the variables related?  
How would seeing  $X$  change my belief in  $Y$ ?)

EXAMPLES: What does a symptom tell me about a disease?  
What does a survey tell us about the election results?

# Structural Causal Models

→ A Structural Causal Model (SCM) is a triple:  $\mathcal{M} := \langle X, U, F \rangle$

- ◆ Two sets of variables:

$$X = \{x_1, \dots, x_N\} \qquad \qquad U = \{u_1, \dots, u_N\}$$

- ◆ A set of functions known as **causal mechanisms**:

$$F = \{f_1, \dots, f_N\}$$

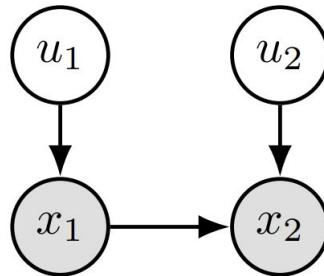
- ◆ The value of each variable is a function of its **parents** (direct causes):

$$x_k := f_k(\mathbf{pa}_k, u_k), \qquad k = 1, \dots, N$$

# Structural Causal Models

→ **Example:** A simple SCM:

- ◆  $x_1, x_2$  are **endogenous** whereas  $u_1, u_2$  are **exogenous**

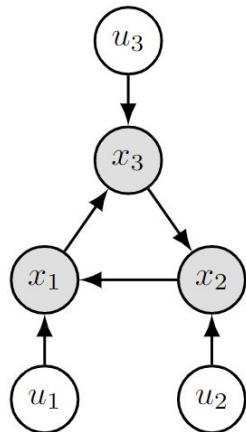


$$x_1 := f_1(u_1), \quad u_1 \sim \mathcal{N}(0, 1)$$

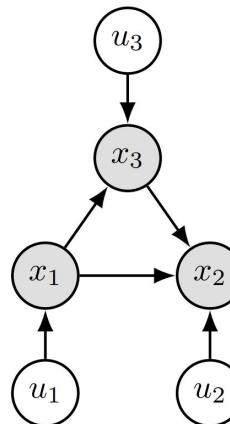
$$x_2 := f_2(x_1, u_2), \quad u_2 \sim \mathcal{N}(0, 1)$$

# Structural Causal Models

- Acyclic SCMs can be represented by a Directed Acyclic Graphs (DAGs), with edges pointing from **causes** to **effects**



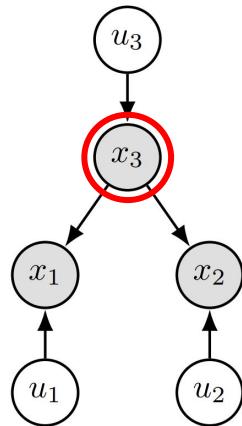
(a) Cyclic



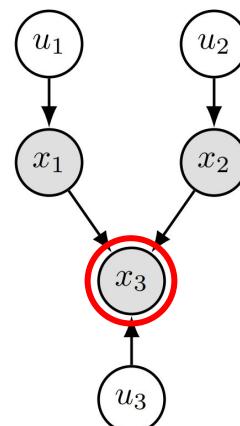
(b) Acyclic

# Structural Causal Models

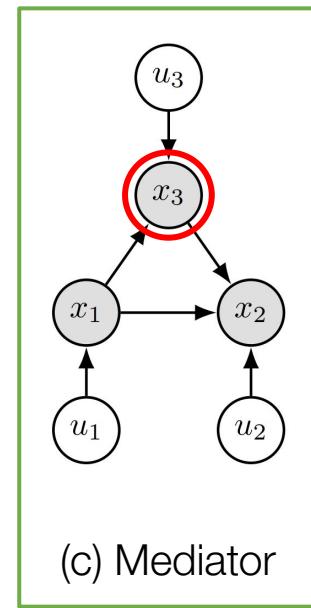
- Acyclic SCMs can be represented by a **Directed Acyclic Graphs (DAGs)**, with edges pointing from **causes** to **effects**



(a) Confounder



(b) Collider



(c) Mediator

# Observational Distribution

(rung 1.)

- SCMs with jointly independent exogenous noises are called **Markovian**:

$$P(u_1, \dots, u_N) = \prod_{k=1}^N P(u_k)$$

- Markovian SCMs induce a unique joint **observational distribution** over the endogenous variables:

$$P_{\mathcal{M}}(x_1, \dots, x_N) = \prod_{k=1}^N P_{\mathcal{M}}(x_k \mid \mathbf{pa}_k)$$

- Each variable is **independent** of its non-descendants given its direct causes (*causal Markov condition*)

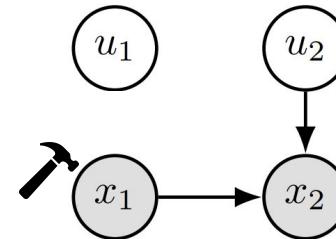
# Interventional Distribution

(rung 2.)

→ SCMs predict the causal effects of actions via interventions

→ Interventions answer causal questions like:

- ◆ E.g. what would  $x_2$  be if we set  $x_1 := c$ ?



→ Interventions replace one or more of the structural assignments and are denoted with the *do* operator:  $do(x_k := c)$

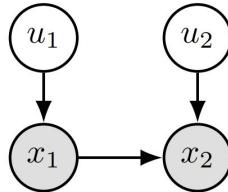
→ This induces a submodel  $\mathcal{M}_c$  and its entailed distribution is known as the interventional distribution:  $P_{\mathcal{M}_c}(X \mid do(c))$

# Counterfactuals

(rung 3.)

- SCMs can consider **hypothetical** scenarios:
  - ◆ Given that we observed  $(x_1, x_2)$ , what would  $x_2$  have been had  $x_1$  been  $c$ ?
  - ◆ All else being equal, would I have been **late** had I not **missed the train**?
- Counterfactual inference involves three steps:
  1. **Abduction:** Update  $P(U)$  given observed evidence, i.e. infer posterior  $P(U | X)$
  2. **Action:** Perform an intervention e.g.  $do(\tilde{x}_k := c)$  and obtain the submodel  $\mathcal{M}_c$
  3. **Prediction:** Use the modified model  $\langle \mathcal{M}_c, P(U | X) \rangle$  to compute counterfactuals

# Example: Computing Counterfactuals



$$x_1 := f_1(u_1) = 1 + u_1$$

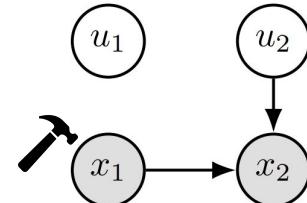
$$x_2 := f_2(x_1, u_2) = 3x_1 + u_2$$

**Q:** Given we observed  $\{x_1=2, x_2=4\}$ , what would  $x_2$  have been had  $x_1$  been 5?

1. Abduction:  $u_1 = f_1^{-1}(x_1) = x_1 - 1 \implies u_1 = 1$

$$u_2 = f_2^{-1}(x_1, x_2) = x_2 - 3x_1 \implies u_2 = -2$$

2. Action:  $\tilde{x}_1 := 5$



3. Prediction:  $\tilde{x}_2 = 3\tilde{x}_1 + u_2 = 3 \cdot 5 - 2 = 13$

# Deep Structural Causal Models<sup>1</sup>

- Leverage **deep generative models** to learn SCM mechanisms:

$$x_k := f_k(\mathbf{pa}_k, u_k)$$

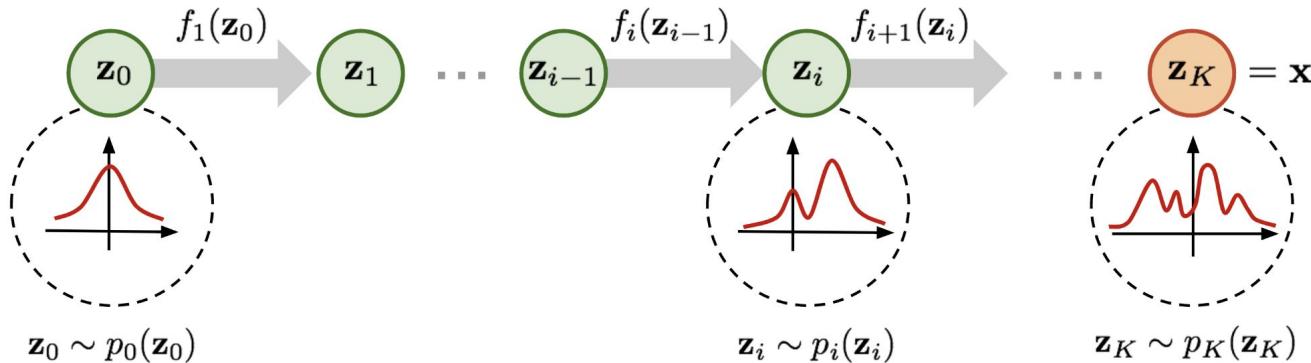
- Tractably estimate causal effects of interventions and perform counterfactual inference, i.e. answer “what if...?” type questions
- Abduction is challenging in complex problems, e.g. medical imaging

**Research Question:** Can we generate plausible high-fidelity image counterfactuals of real-world data, and if so, how do we evaluate them?

<sup>1</sup> Pawlowski, Castro, Glocker. Deep Structural Causal Models for Tractable Counterfactual Inference. NeurIPS 2020

# Deep Mechanisms: Normalizing Flows

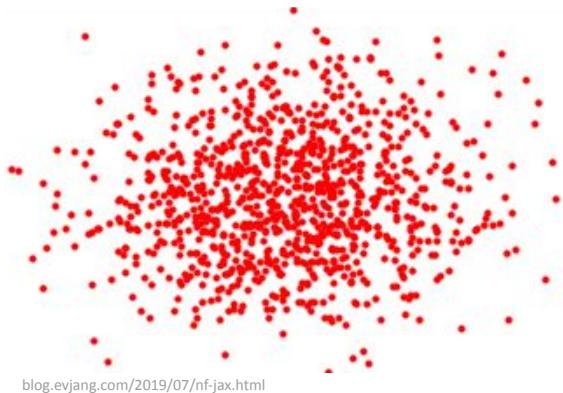
- Normalizing Flows (NFs) build complex probability distributions via successive (invertible) transformations of simple distributions



[lilianweng.github.io/posts/2018-10-13-flow-models/](https://lilianweng.github.io/posts/2018-10-13-flow-models/)

# Deep Mechanisms: Normalizing Flows

- Normalizing Flows (NFs) build complex probability distributions via successive (invertible) transformations of simple distributions



[blog.evjang.com/2019/07/nf-jax.html](http://blog.evjang.com/2019/07/nf-jax.html)

- **TLDR:** NFs enable invertible mechanisms and deterministic abduction:

$$x = f_{\theta}(\mathbf{pa}_x, u_x) \quad u_x = f_{\theta}^{-1}(\mathbf{pa}_x, x)$$

# Intuition: Normalizing Flows

→ Change-of-variables formula:

$$p(x) = p(u) \left| \frac{du}{dx} \right|$$

volume correction

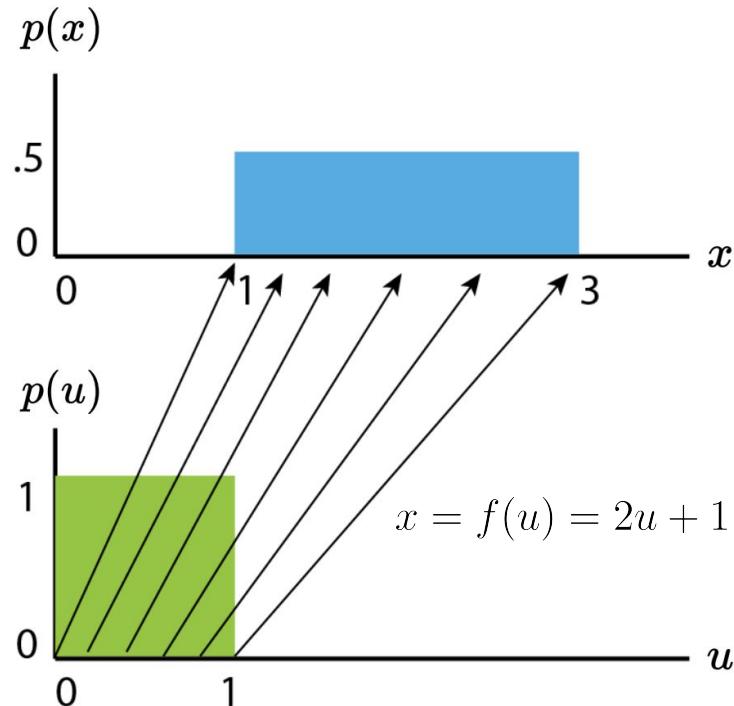
$$u = f^{-1}(x) = \frac{x - 1}{2}$$

$$\frac{du}{dx} = \frac{df^{-1}(x)}{dx} = \frac{d}{dx} \left( \frac{x - 1}{2} \right) = \frac{1}{2}$$

$$p(u) = p_U(f^{-1}(x)) = \frac{1}{1 - 0} = 1$$

PDF of  $\mathcal{U}[a, b]$  is  $\frac{1}{b - a}$

$$X = f(U)$$



# Training Objective: Normalizing Flows

- Maximum (log) likelihood training objective:

$$\log p(x) = \log p(f_\theta^{-1}(x)) + \log \left| \frac{df_\theta^{-1}(x)}{dx} \right|$$

- Can condition on parents via a learned parameterised function:

$$\log p(x \mid \mathbf{pa}_x) = \log p(f_\theta^{-1}(\mathbf{pa}_x, x)) + \log \left| \det \left( \frac{\partial f_\theta^{-1}(\mathbf{pa}_x, x)}{\partial x} \right) \right|$$

- In multivariate settings, we need to compute the determinant of a **Jacobian** matrix, which can become pretty expensive!

# Deep Mechanisms: Variational Autoencoders

→ Well suited for modelling structured variable mechanisms, e.g. for images

→ Caveat: **abduction** is non-deterministic

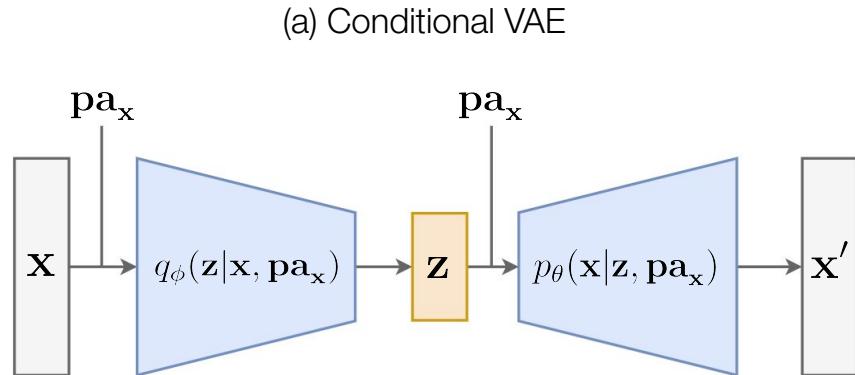
→ Causal mechanism:

invertible                   non-invertible

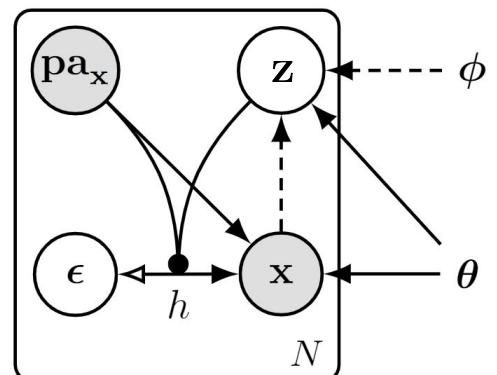
$$\begin{aligned} \mathbf{x} &:= f_{\theta}(\mathbf{pa}_x, \mathbf{u}_x) = h(\epsilon; g_{\theta}(z, \mathbf{pa}_x)) \\ &= \mu(z, \mathbf{pa}_x) + \sigma(z, \mathbf{pa}_x) \odot \epsilon, \quad \epsilon \sim \mathcal{N}(0, \mathbf{I}) \end{aligned}$$

→ Factored exogenous noise:  $p(\mathbf{u}_x) = p_{\theta}(\mathbf{z})p(\epsilon)$

$$p(\mathbf{u}_x | \mathbf{x}, \mathbf{pa}_x) \approx q_{\phi}(\mathbf{z} | \mathbf{x}, \mathbf{pa}_x) \delta(\epsilon - h^{-1}(\mathbf{x}; g_{\theta}(\mathbf{z}, \mathbf{pa}_x)))$$



Causal mechanism for  $\mathbf{x}$



# Deep Mechanisms: Variational Autoencoders

→ **Example:** computing counterfactuals:

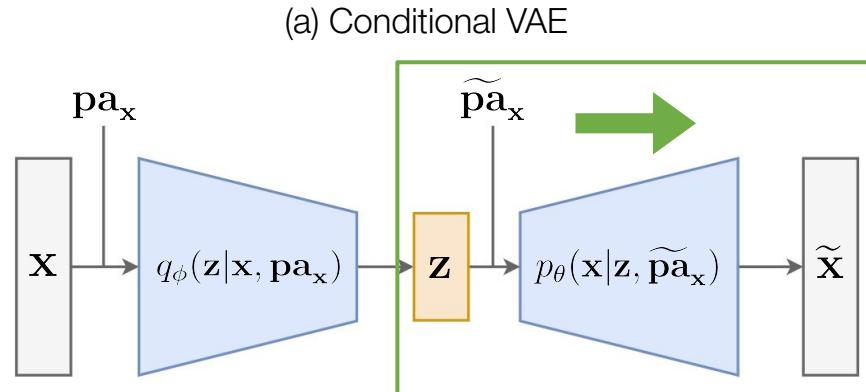
1. **Abduction:**  $\mathbf{z} \sim q_\phi(\mathbf{z} \mid \mathbf{x}, \mathbf{pa}_x)$

$$\epsilon = h^{-1}(\mathbf{x}; g_\theta(\mathbf{z}, \mathbf{pa}_x)) = \frac{\mathbf{x} - \mu(\mathbf{z}, \mathbf{pa}_x)}{\sigma(\mathbf{z}, \mathbf{pa}_x)}$$

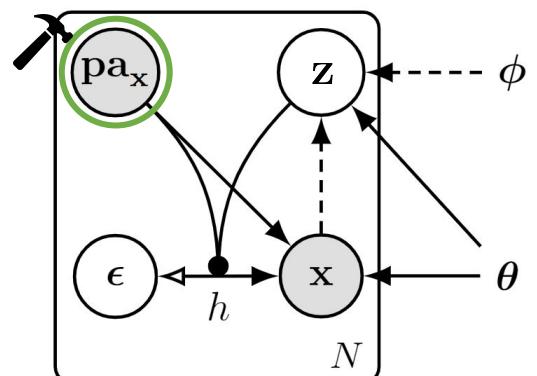
2. **Action:**  $do(\mathbf{pa}_x := \widetilde{\mathbf{pa}}_x)$

3. **Prediction:**  $\widetilde{\mathbf{x}} \sim p_\theta(\widetilde{\mathbf{x}} \mid \mathbf{z}, \widetilde{\mathbf{pa}}_x)$

$$\widetilde{\mathbf{x}} := h(\epsilon; g_\theta(\mathbf{z}, \widetilde{\mathbf{pa}}_x)) = \mu(\mathbf{z}, \widetilde{\mathbf{pa}}_x) + \sigma(\mathbf{z}, \widetilde{\mathbf{pa}}_x) \odot \epsilon$$

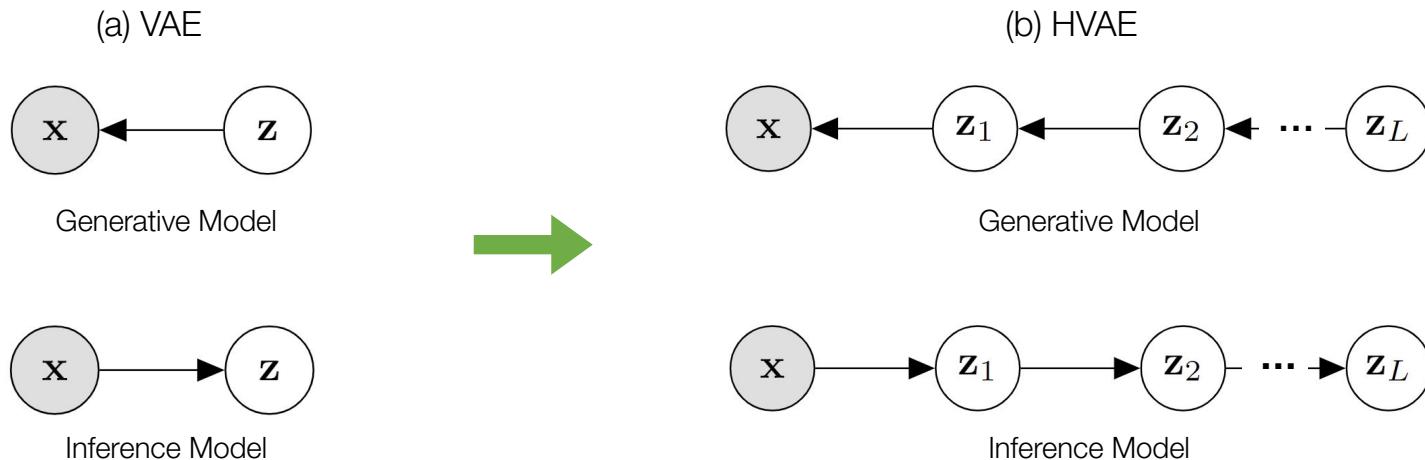


Causal mechanism for  $\mathbf{x}$



# Deep Mechanisms: Hierarchical VAEs

- To produce high-fidelity **image counterfactuals**, we require a powerful generative model that is amenable to principled **abduction**
- Hierarchical VAEs extend VAEs to multiple layers of latent variables:



# Deep Mechanisms: Hierarchical VAEs

- To produce high-fidelity **image counterfactuals**, we require a powerful generative model that is amenable to principled **abduction**
- Hierarchical VAEs extend VAEs to multiple layers of latent variables:

$$p(\mathbf{x}, \mathbf{z}_{1:L}) = \underbrace{p(\mathbf{x} \mid \mathbf{z}_{1:L}) p(\mathbf{z}_L)}_{\text{Generative Model}} \prod_{i=1}^{L-1} p(\mathbf{z}_i \mid \mathbf{z}_{>i})$$

- **Goal:** Optimize our model  $p$  to be close to a given data distribution  $p_{\text{data}}$

$$\log p(\mathbf{x}) \geq \mathbb{E}_{q(\mathbf{z}_{1:L} \mid \mathbf{x})} [\log p(\mathbf{x} \mid \mathbf{z}_{1:L})] - D_{\text{KL}}(q(\mathbf{z}_{1:L} \mid \mathbf{x}) \parallel p(\mathbf{z}_{1:L})) =: \text{ELBO}(\mathbf{x})$$

## NVAE: A Deep Hiera

Ar  
{avahd

Casper Kaae Sønderby\*  
casperkaae@gmail.com

Normalizing flows, autoregressive generative models, or deep generative learning, are tractable sampling and easily tractably outperformed by other neural models. While the majority of challenges, we explore the architectures for hierarchical VAEs, which VAE built for image generation, batch normalization. NVAE is a distribution and its training is NVAE achieves state-of-the-art models on the MNIST, CIFAR provides a strong baseline on F images on CelebA HQ as shown the first successful VAE application code is available at <http://>



Søren Kaae Sønd  
skaaesonderby@gn

Variational autoencoders are powerful generative models with several layers. However, they often train which limits the improvement of the model. We propose a new inference method that recursively corrects the generated likelihood in a process resembling backpropagation. We show that this model provides a log-likelihood lower bound competitive with state-of-the-art Variational Autoencoders and that the analysis of the learned hierarchical inference model is qualitative. Finally, we observe that batch normalization and deterministic warm-up (gradually turning on the KL-term) are crucial for training variational models with many stochastic layers.

casperkaae@gmail.com

Diederik P. Kingma  
dpkingma@openai.com

Tim Salimans  
tim@openai.com

Rafal Jozefowicz  
rafał@openai.com

Xi Chen  
peter@openai.com

Ilya Sutskever  
ilya@openai.com

Max Welling\*  
M.Welling@uva.nl

## Abstract

The framework of normalizing flows provides a general strategy for flexible variational inference of posteriors over latent variables. We propose a new type of normalizing flow, inverse autoregressive flow (IAF), that, in contrast to earlier published flows, scales well to high-dimensional latent spaces. The proposed flow consists of a chain of invertible transformations, where each transformation is based on an autoregressive neural network. In experiments, we show that IAF significantly improves upon diagonal Gaussian approximate posteriors. In addition, we demonstrate that a novel type of variational autoencoder, coupled with IAF, is competitive with neural autoregressive models in terms of attained log-likelihood on natural images, while allowing significantly faster synthesis.

ZE AUTOREGRESSIVE  
M THEM ON IMAGES

ime, generates samples quickly on all natural image benchmarks. It can also represent autoregressive distributions, when made sufficiently deep. IAF has been shown to significantly outperform VAEs in log-likelihood by scaling a VAE to greater resolution. Using it on CIFAR-10, ImageNet, and other datasets, very deep VAEs achieve higher log-likelihoods thousands of times faster, and more accurately. Qualitative studies suggest this is due to better visual representations. We release the source code at [http://openai.com/vdvae](http://).



# Quick Detour: Connection to Diffusion Models

→ A Diffusion Model is a type of Hierarchical VAE where:

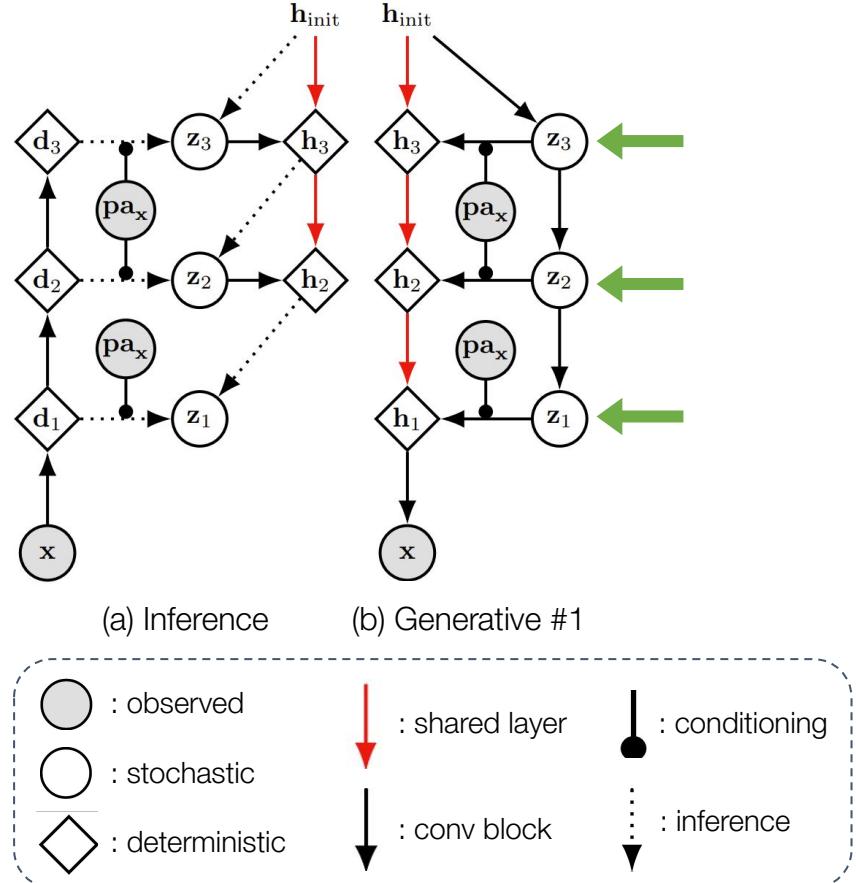
- I. The encoder  $q(\mathbf{z}_{1:T} \mid \mathbf{x})$  is fixed rather than learned
- II. Each latent variable  $\mathbf{z}_t$  has the same dimensionality as  $\mathbf{x}$
- III. A single (denoising) model is shared amongst all layers in the hierarchy

$$\begin{aligned} -\text{ELBO}(\mathbf{x}) &= \mathbb{E}_{q(\mathbf{z}_{1:T} \mid \mathbf{x})} [-\log p(\mathbf{x} \mid \mathbf{z}_{1:T})] + D_{\text{KL}}(q(\mathbf{z}_{1:T} \mid \mathbf{x}) \parallel p(\mathbf{z}_{1:T})) \\ &= \mathbb{E}_{q(\mathbf{z}_1 \mid \mathbf{x})} [-\log p(\mathbf{x} \mid \mathbf{z}_1)] + D_{\text{KL}}(q(\mathbf{z}_T \mid \mathbf{x}) \parallel p(\mathbf{z}_T)) + \mathcal{L}_T(\mathbf{x}) \quad \text{diffusion loss} \end{aligned}$$

$$\mathcal{L}_T(\mathbf{x}) = \sum_{t=2}^T \mathbb{E}_{q(\mathbf{z}_t \mid \mathbf{x})} [D_{\text{KL}}(q(\mathbf{z}_{t-1} \mid \mathbf{z}_t, \mathbf{x}) \parallel p(\mathbf{z}_{t-1} \mid \mathbf{z}_t))] = \boxed{\frac{T}{2} \mathbb{E}_{\boldsymbol{\epsilon} \sim \mathcal{N}(0, \mathbf{I}), t \sim U\{1, T\}} [w(t) \|\boldsymbol{\epsilon} - \hat{\boldsymbol{\epsilon}}_\theta(\mathbf{z}_t; t)\|_2^2]}$$

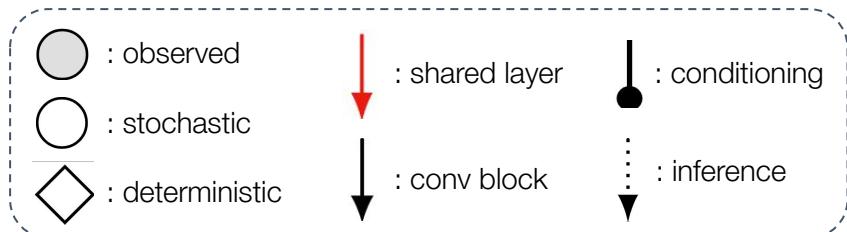
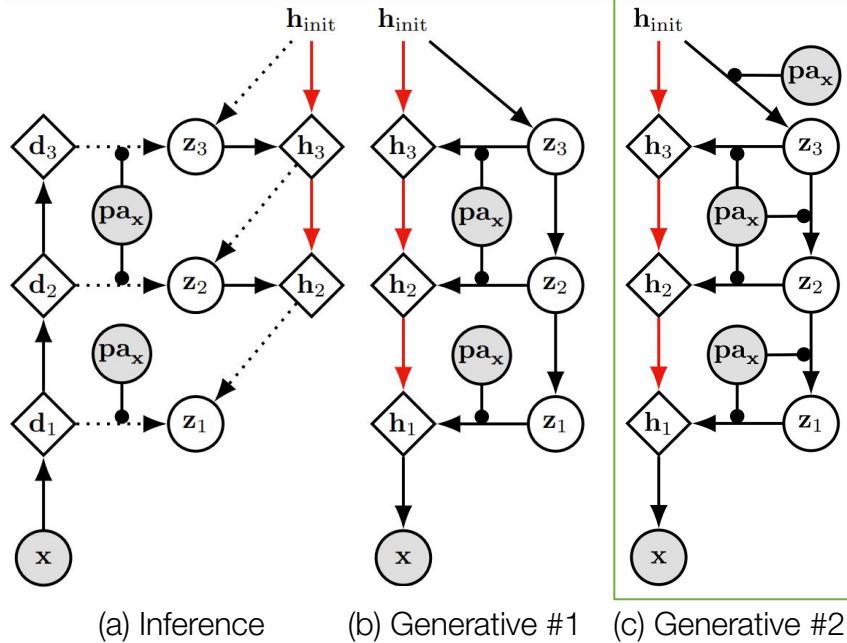
# Deep Mechanisms: Conditional HVAEs

- Generative model structures:
  1. Exogenous Prior:  $p_\theta(\mathbf{z}_{1:L})$



# Deep Mechanisms: Conditional HVAEs

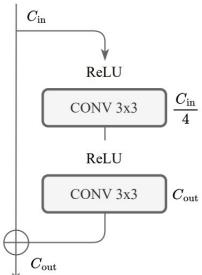
- Generative model structures:
1. Exogenous Prior:  $p_\theta(\mathbf{z}_{1:L})$
  2. Conditional Prior:  $p_\theta(\mathbf{z}_{1:L} \mid \mathbf{pa}_x)$



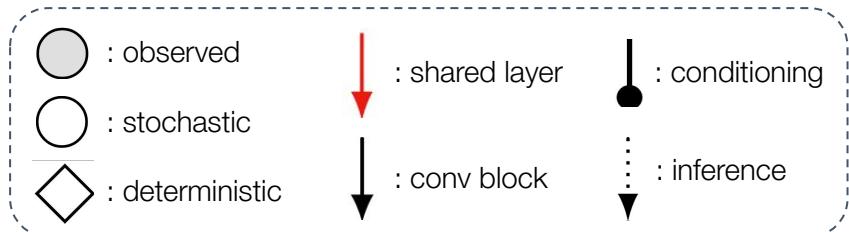
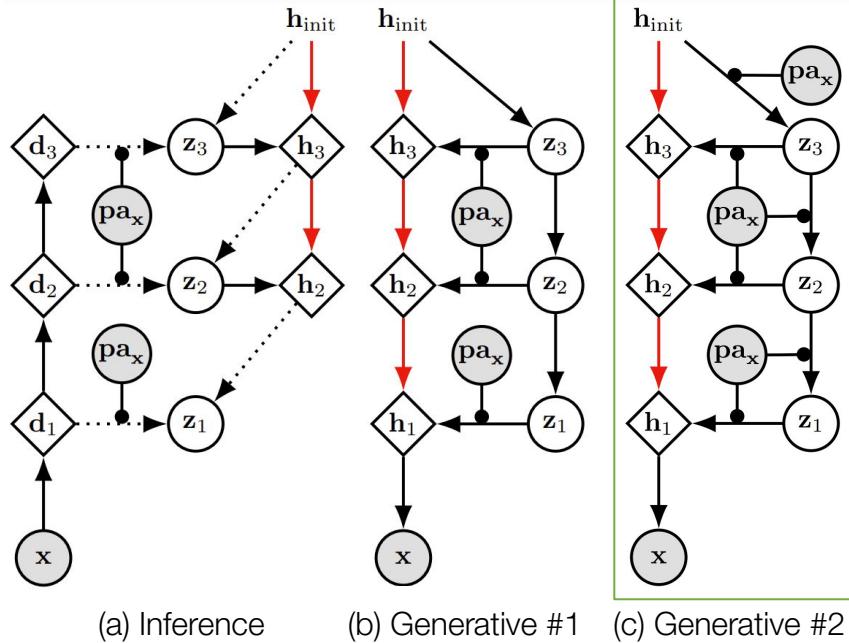
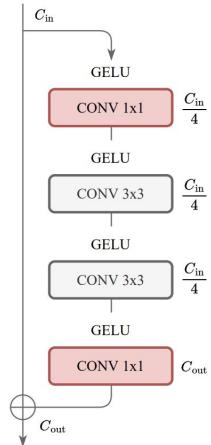
# Deep Mechanisms: Conditional HVAEs

- Generative model structures:
  1. Exogenous Prior:  $p_\theta(\mathbf{z}_{1:L})$
  2. Conditional Prior:  $p_\theta(\mathbf{z}_{1:L} \mid \mathbf{pa}_x)$

- Standard residual blocks:



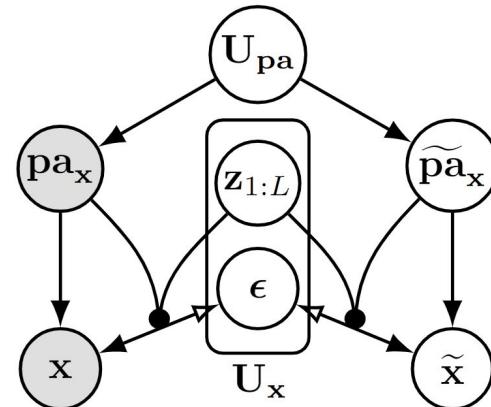
Option 2:



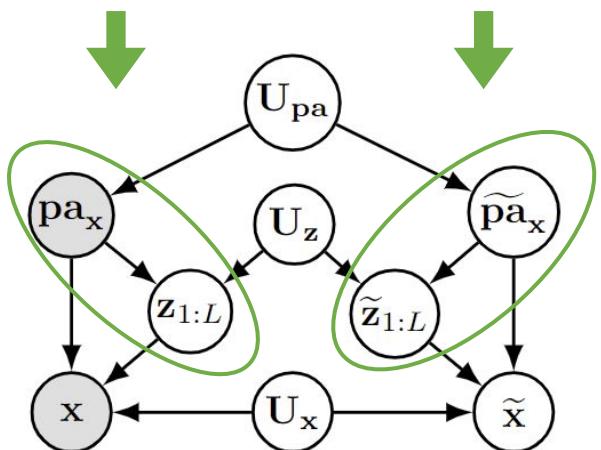
# Latent Mediator Model

- The conditional prior  $p_\theta(\mathbf{z}_{1:L} | \mathbf{pa}_x)$  induces a **latent mediator**, as  $\mathbf{z}_{1:L}$  is no longer exogenous
- Nonetheless, the underlying SCM has a Markovian interpretation:

$$p(\mathbf{U}) = p(\mathbf{U}_x) \left( \prod_{k=1}^K p(\mathbf{U}_{\mathbf{pa}_k}) \right) \left( \prod_{i=1}^L p(\mathbf{U}_{\mathbf{z}_i}) \right)$$



(a) Exogenous Prior:  $p_\theta(\mathbf{z}_{1:L})$



(b) Latent Mediator:  $p_\theta(\mathbf{z}_{1:L} | \mathbf{pa}_x)$

# Latent Mediator Model

## Causal Mediation Analysis

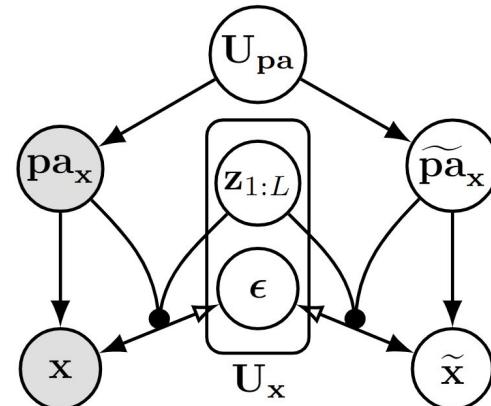
The study of how a treatment effect is mediated by another variable, to help explain why or how an individual may respond to certain stimulus.

- Enables estimation of **Direct (DE)**, **Indirect (IE)** and **Total (TE)** causal effects:

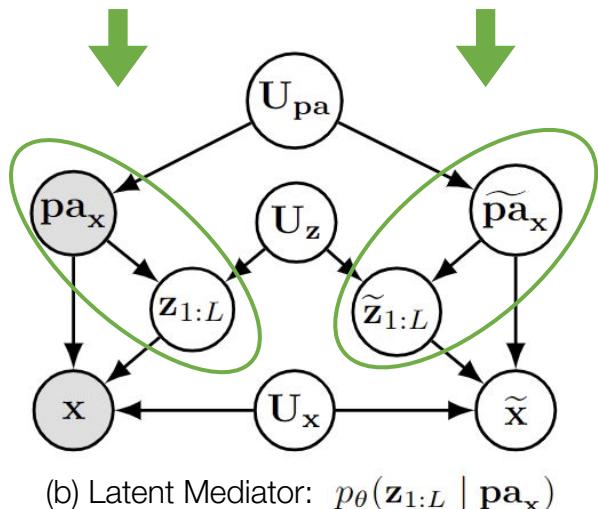
$$DE_x(\tilde{pa}_x) = \mathbb{E} [g_\theta(\tilde{pa}_x, z_{1:L}) - g_\theta(pa_x, z_{1:L})]$$

$$IE_x(\tilde{z}_{1:L}) = \mathbb{E} [g_\theta(pa_x, \tilde{z}_{1:L}) - g_\theta(pa_x, z_{1:L})]$$

$$TE_x(\tilde{pa}_x, \tilde{z}_{1:L}) = \mathbb{E} [g_\theta(\tilde{pa}_x, \tilde{z}_{1:L}) - g_\theta(pa_x, z_{1:L})]$$



(a) Exogenous Prior:  $p_\theta(z_{1:L})$

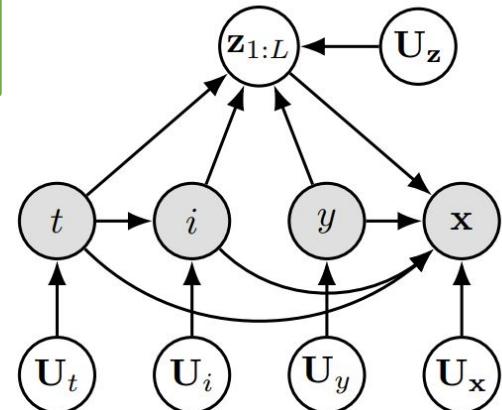
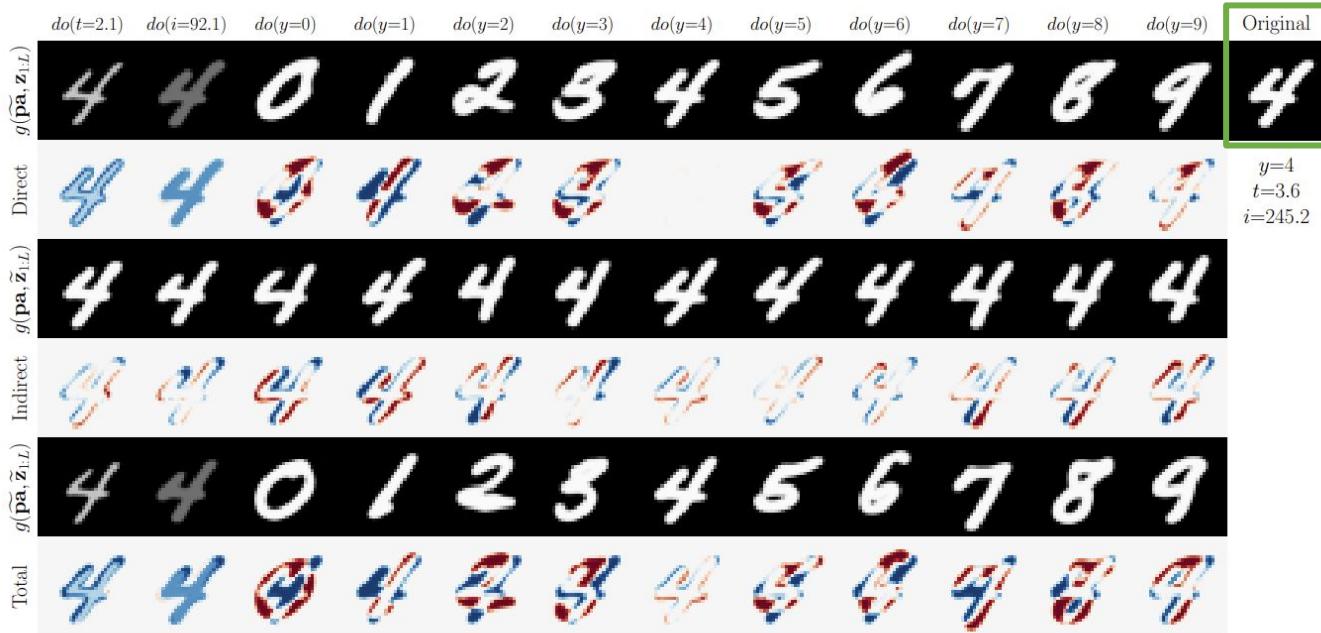


(b) Latent Mediator:  $p_\theta(z_{1:L} | pa_x)$

# Case Study: Morpho-MNIST



**Hugging Face**  
Code & Demo



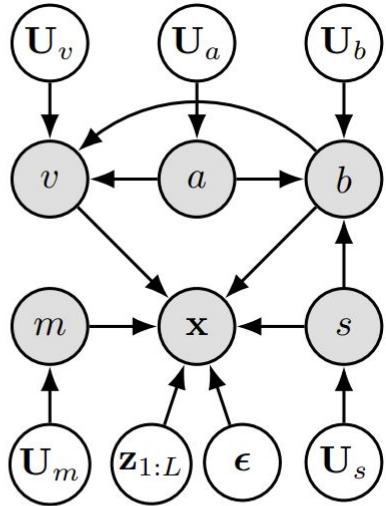
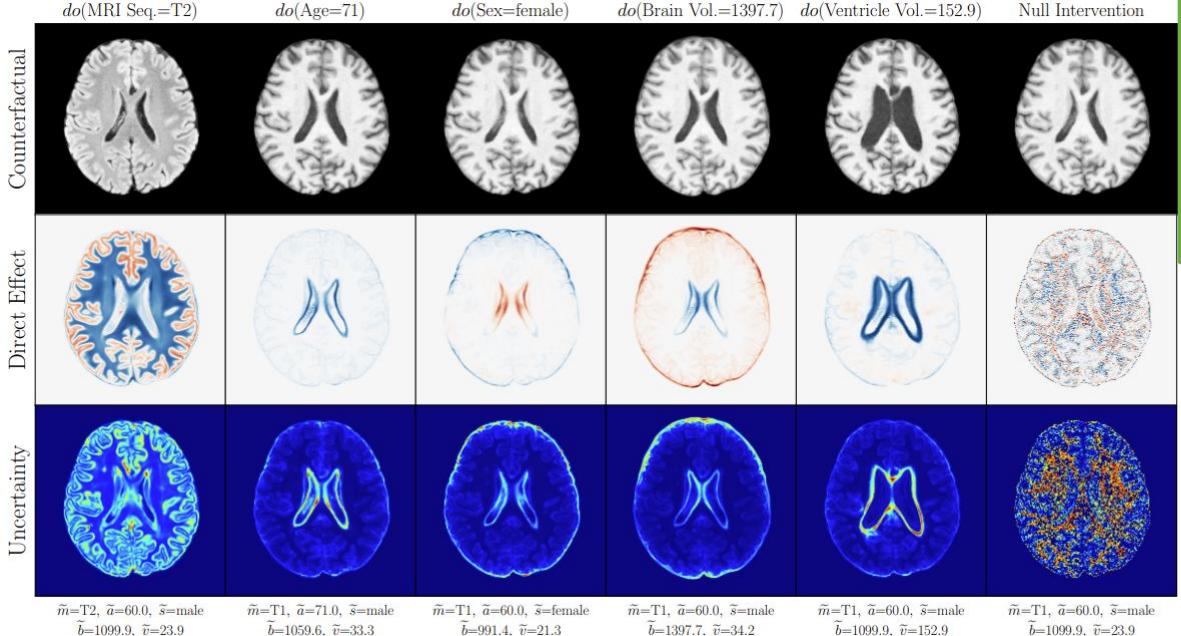
(a) Latent mediator SCM for Morpho-MNIST. Observed variables in the graph: image ( $x$ ), digit ( $y$ ), stroke thickness ( $t$ ) and pixel intensity ( $i$ ).

Figure 3: Morpho-MNIST counterfactuals from our latent mediator SCM. *Direct*, *indirect* and *total* causal effects of interventions are shown (red: increase; blue: decrease). Recall that  $\tilde{x}_{IE} \sim g(\mathbf{p}_x, \tilde{\mathbf{z}}_{1:L})$  are cross-world counterfactuals, i.e. the potential outcome of  $x$  given  $\mathbf{p}_x$  and the (counterfactual) mediator we would have observed  $\tilde{\mathbf{z}}_{1:L}$  had  $\mathbf{p}_x$  been  $\tilde{\mathbf{p}}_x$ .

# Case Study: Brain Imaging



**Hugging Face**  
Code & Demo

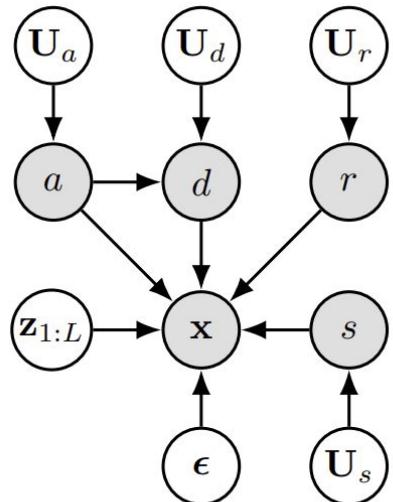
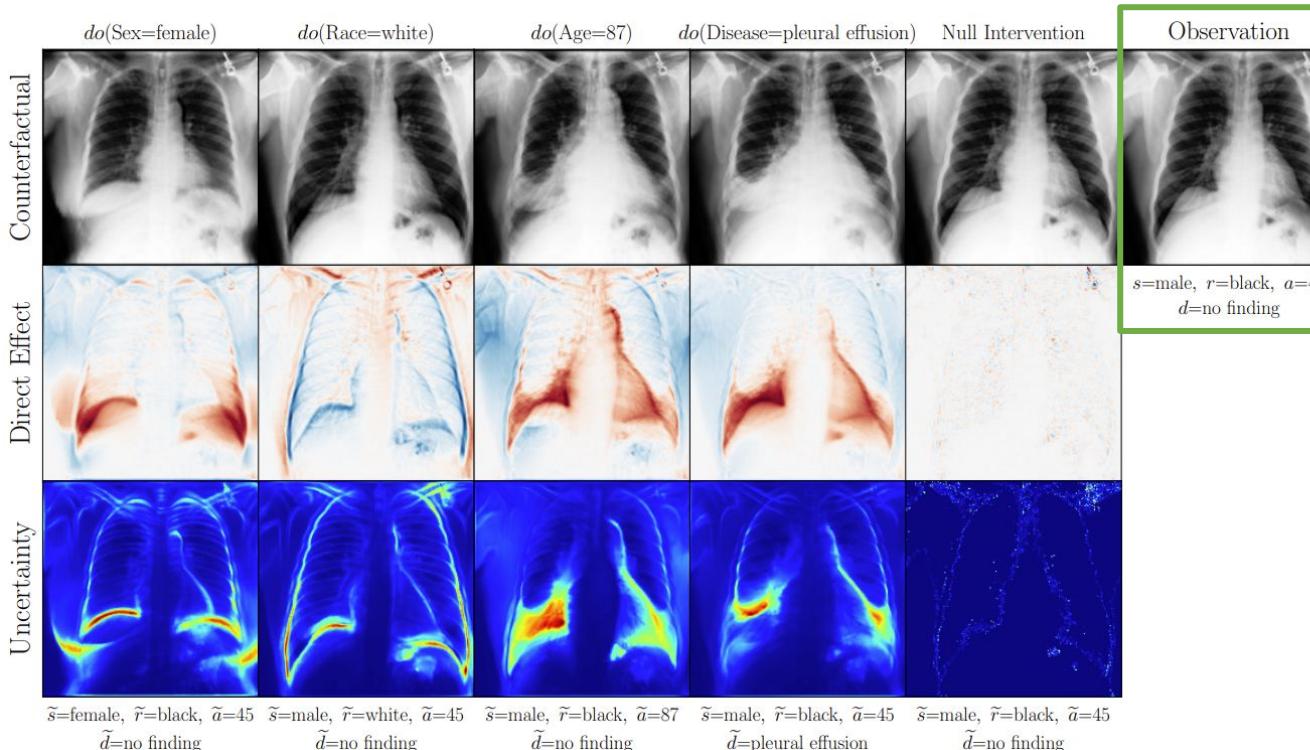


(a) Deep SCM for UK Biobank.  
MRI Seq. ( $m$ ), age ( $a$ ), sex ( $s$ ),  
brain ( $b$ ) & ventricle ( $v$ ) volume.

# Case Study: Chest X-ray Imaging



**Hugging Face**  
Code & Demo



(a) Deep SCM for MIMIC-CXR. The variables in the causal graph are: age ( $a$ ), sex ( $s$ ), race ( $r$ ), disease ( $d$ ) and chest x-ray ( $x$ ). The disease  $d$  is pleural effusion.

# Evaluating Counterfactuals: Axiomatic Properties

- The **soundness theorem** states that the properties of **composition**, **effectiveness**, and **reversibility** are necessary in all causal models (Galles & Pearl, 1998).
- The **completeness theorem** states that these properties are sufficient (Halpern, 1998).
- We can measure counterfactual soundness using these axiomatic properties

Published as a conference paper at ICLR 2023

## MEASURING AXIOMATIC SOUNDNESS OF COUNTERFACTUAL IMAGE MODELS

Miguel Monteiro<sup>1†</sup> Fabio De Sousa Ribeiro<sup>1†</sup> Nick Pawlowski<sup>2</sup> Daniel C. Castro<sup>1,2</sup> Ben Glocker<sup>1</sup>

<sup>1</sup>Imperial College London, <sup>2</sup>Microsoft Research Cambridge. <sup>†</sup>Joint first authors  
(miguel.monteiro, f.de-sousa-ribeiro, b.glocker}@imperial.ac.uk

### ABSTRACT

We present a general framework for evaluating image counterfactuals. The power and flexibility of deep generative models make them valuable tools for learning mechanisms in structural causal models. However, their flexibility makes counterfactual identifiability impossible in the general case. Motivated by these issues, we revisit Pearl's axiomatic definition of counterfactuals to determine the necessary constraints of any counterfactual inference model: *composition*, *reversibility*, and *effectiveness*. We frame counterfactuals as functions of an input variable, its parents, and counterfactual parents and use the axiomatic constraints to restrict the set of functions that could represent the counterfactual, thus deriving distance metrics between the approximate and ideal functions. We demonstrate how these metrics can be used to compare and choose between different approximate counterfactual inference models and to provide insight into a model's shortcomings and trade-offs.

# Example: Composition Axiom

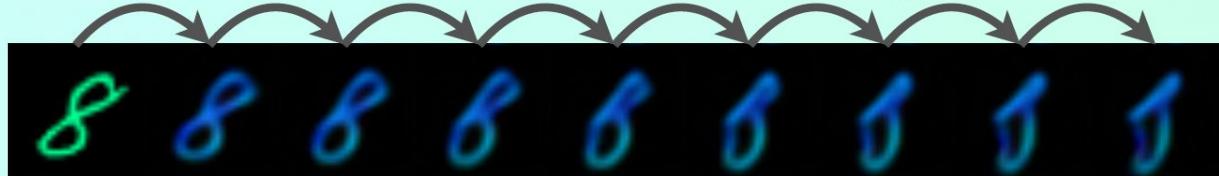
## Composition

High composition (de-biased model)

null-intervention: do\_colour = original\_colour  
do\_digit = original\_digit



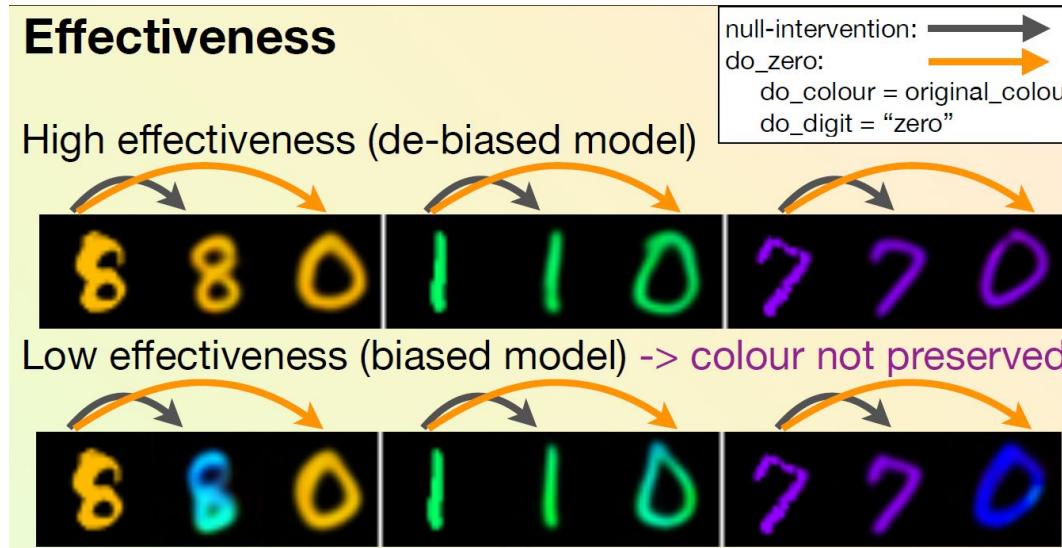
Low composition (biased model) -> identity is lost



**Definition.** Intervening on a variable to have the value it would otherwise have without the intervention will not affect other variables.

→ Implies the existence of a null-intervention

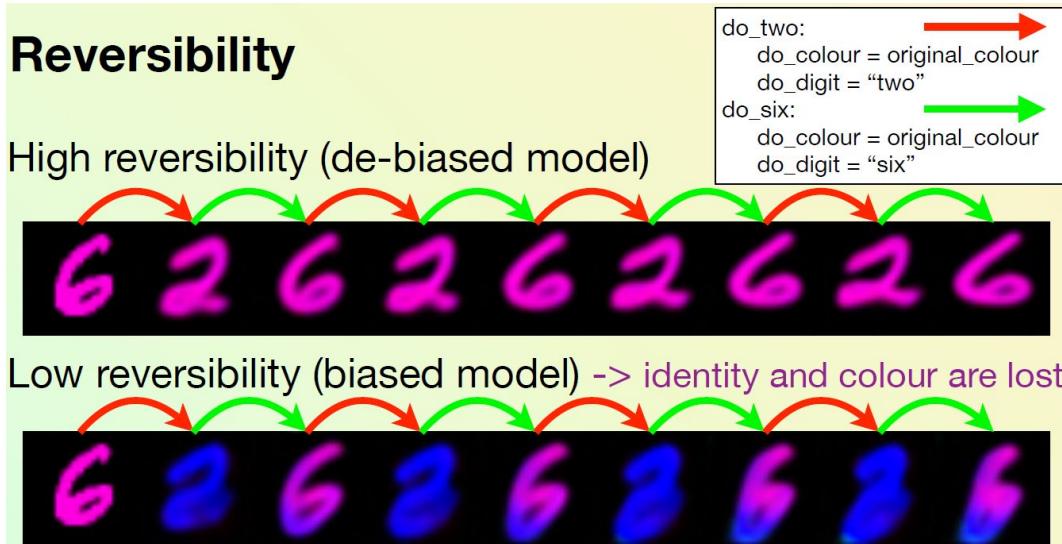
# Example: Effectiveness Axiom



**Definition.** Intervening on a variable to have a specific value will cause the variable to take on that value.

→ Caveat: often relies on a **pseudo-oracle**, e.g. a classifier/regressor

# Example: Reversibility Axiom

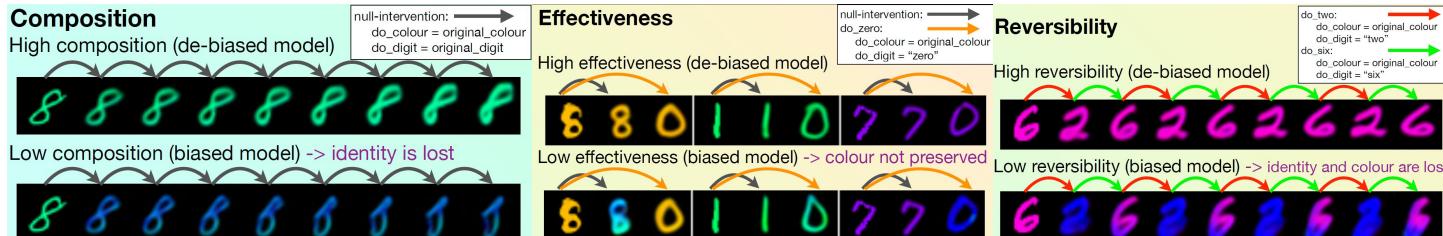
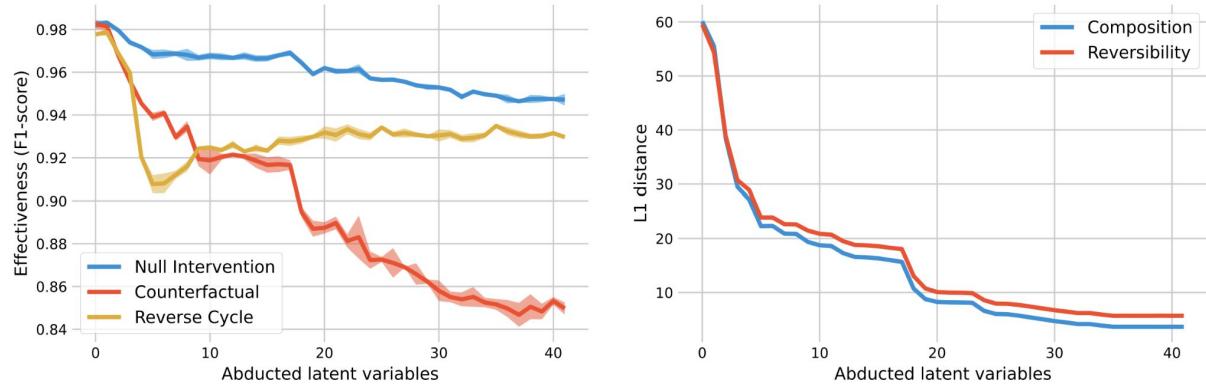


**Definition.** Precludes multiple solutions due to feedback loops. If setting a variable  $X$  to a value  $x$  results in a value  $y$  for a variable  $Y$ , and setting  $Y$  to a value  $y$  results in value  $x$  for  $X$ , then  $X$  and  $Y$  will naturally take on the values  $x$  and  $y$ .

- Follows directly from composition in recursive systems such as DAGs

# Evaluating Counterfactuals: Axiomatic Properties

→ **TLDR:** We identify an inherent trade-off



# Conclusion & Outlook

- Deep SCMs can generate plausible high-fidelity **counterfactuals** of medical images as measured by **axiomatic soundness** of counterfactuals
- Tractable estimation of **direct**, **indirect** and **total** causal effects for high-dimensional structured variables
- **Limitations:**
  - ◆ Only consider Markovian SCMs; although Markovianity is a common assumption in causality literature, it is strong in most cases
  - ◆ Measuring counterfactual effectiveness relies on separately trained classifiers

# Conclusion & Outlook

- Deep SCMs can generate plausible high-fidelity **counterfactuals** of medical images as measured by **axiomatic soundness** of counterfactuals

## Future Work

- I. Targeted data augmentation to improved robustness, sample efficiency & fairness
- II. Providing principled causal explanations, e.g. through mediation analysis
- III. Improving **counterfactual soundness**, via ML techniques and/or medical guidance
- IV. Provide theoretical guarantees of identifiability under plausible causal assumptions
- V. Extensions to **Semi-Markovian** and **Non-Markovian** settings



**Hugging Face**  
Code & Demo



# Thank you for listening

Happy to take questions!

fdesousa@ic.ac.uk



Fabio De Sousa Ribeiro



Miguel Monteiro



Nick Pawlowski



Daniel C. Castro



Tian Xia



Ben Glocker

