

1. Overview

Research Question: How can we generate **plausible** high-fidelity **counterfactuals** of real data, and how do we **evaluate** them?

Motivation:

- Counterfactuals are useful for **explainability**, **data aug.**, **fairness** etc
- High fidelity counterfactuals of **structured variables** are challenging
- Identifiability** guarantees are **absent** in the general case
- Prior work is mostly theoretical, we take a pragmatic approach

Contribution:

- Hierarchical causal mechanisms** for high-dim **structured variables**
- Latent mediator** model for **direct**, **indirect** and **total** effect estimation
- Counterfactual training** to mitigate ignored counterfactual conditioning
- Demonstrate the **axiomatic soundness** of our inferred counterfactuals

2. Causal Mechanisms: Hierarchical VAEs

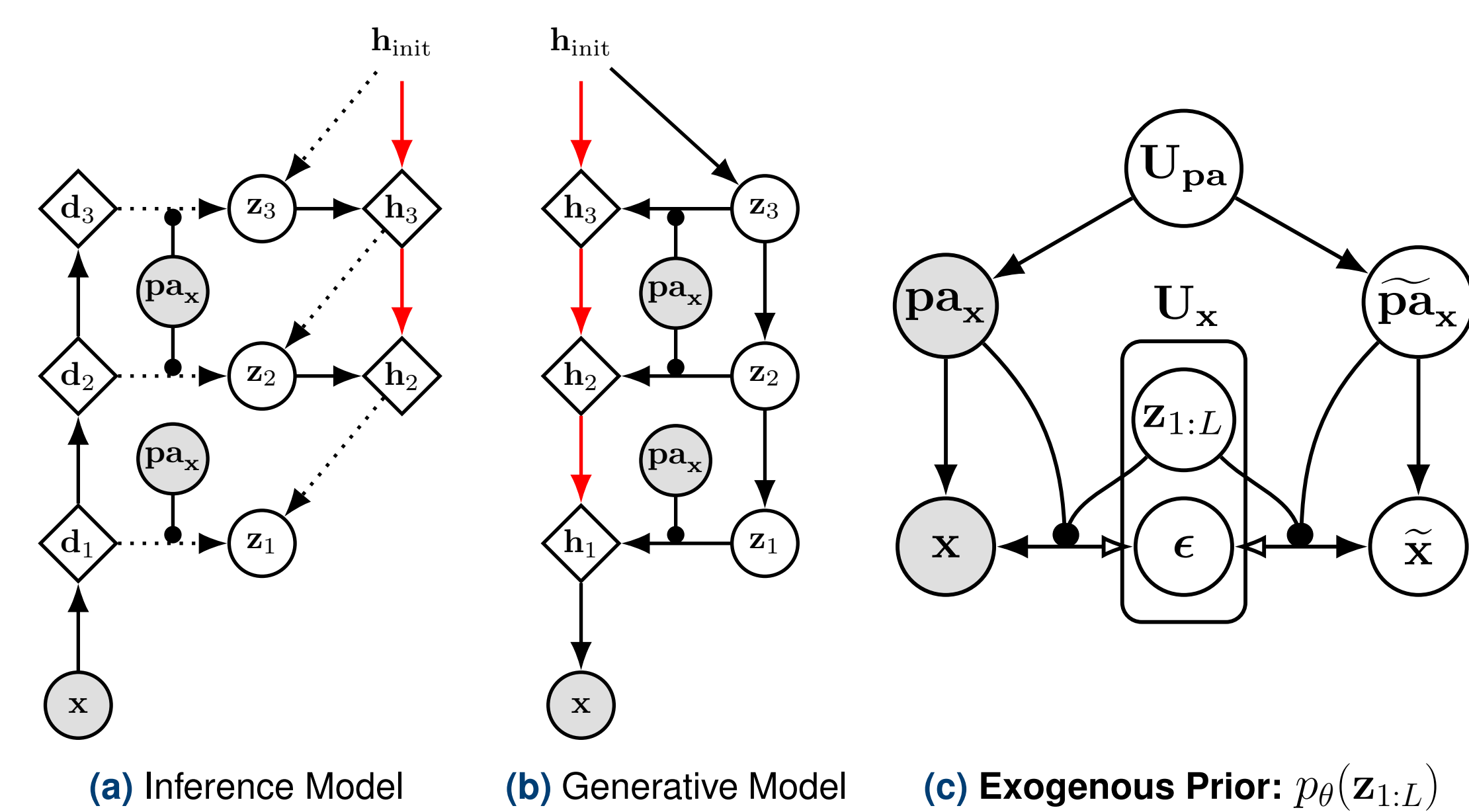


Figure: Twin network of our deep SCM (c) induced by an exogenous prior (b), where $z_{1:L}$ is part of x 's exogenous noise U_x . Directly compatible with Pawlowski et al. (2020)'s framework.

Abduction:

$$z_{1:L} \sim q_\phi(z_{1:L} | x, pa_x)$$

$$\epsilon = h^{-1}(x; g_\theta(z_{1:L}, pa_x)) = \frac{x - \mu(z_{1:L}, pa_x)}{\sigma(z_{1:L}, pa_x)}$$

Action[†]: $do(pa_x := \tilde{pa}_x)$

Prediction:

$$\tilde{x} \sim p_\theta(\tilde{x} | z_{1:L}, \tilde{pa}_x) = h(\epsilon; g_\theta(z_{1:L}, \tilde{pa}_x)) = \mu(z_{1:L}, \tilde{pa}_x) + \sigma(z_{1:L}, \tilde{pa}_x) \odot \epsilon$$

[†]Counterfactual parents \tilde{pa}_x are a result of upstream interventions on the associated causal graph.

3. 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.

Estimating **Direct** (DE), **Indirect** (IE) and **Total** (TE) causal effects[†]:

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

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

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

[†]Identifiability assumptions: sequential ignorability (Imai et al., 2010); iVAE setup (Khemakhem et al., 2020).

4. Counterfactual Training

- A technique for improving **axiomatic effectiveness** of counterfactuals.

$$I(\tilde{pa}_k; \tilde{x}) \geq \mathbb{E}_{p(\tilde{pa}_k, \tilde{x})} [\log q_\psi(\tilde{pa}_k | \tilde{x})] - H(\tilde{pa}_k)$$

Our **constrained optimization** objective:

$$\arg \min_{\theta, \phi} \mathbb{E}_{p_{\text{data}}(x, pa_x)} [\mathcal{L}_{CT}(\mathcal{M}; x, pa_x)] \quad \text{s.t. } \mathcal{F}_{FE}(\theta, \phi; x, pa_x) \leq c,$$

$$\mathcal{L}_{CT}(\mathcal{M}; x, pa_x) = - \sum_k \mathbb{E}_{\tilde{pa}_k \sim p(pa_k), \tilde{x} \sim P_{\mathcal{M}}(\tilde{x} | do(\tilde{pa}_k), x)} [\log q_\psi(\tilde{pa}_k | \tilde{x})]$$

References:

{Pawlowski[†], Castro^{*}} et al. Deep Structural Causal Models for Tractable Counterfactual Inference. NeurIPS 2020.
{Monteiro^{*}, De Sousa Ribeiro^{*}} et al. Measuring Axiomatic Soundness of Counterfactual Image Models. ICLR 2023.
Khemakhem et al. Variational Autoencoders and Nonlinear ICA: A Unifying Framework. AISTATS 2020.
Imai et al. Identification, Inference and Sensitivity Analysis for Causal Mediation Effects. Statistical Science, 2010.

5. Morpho-MNIST

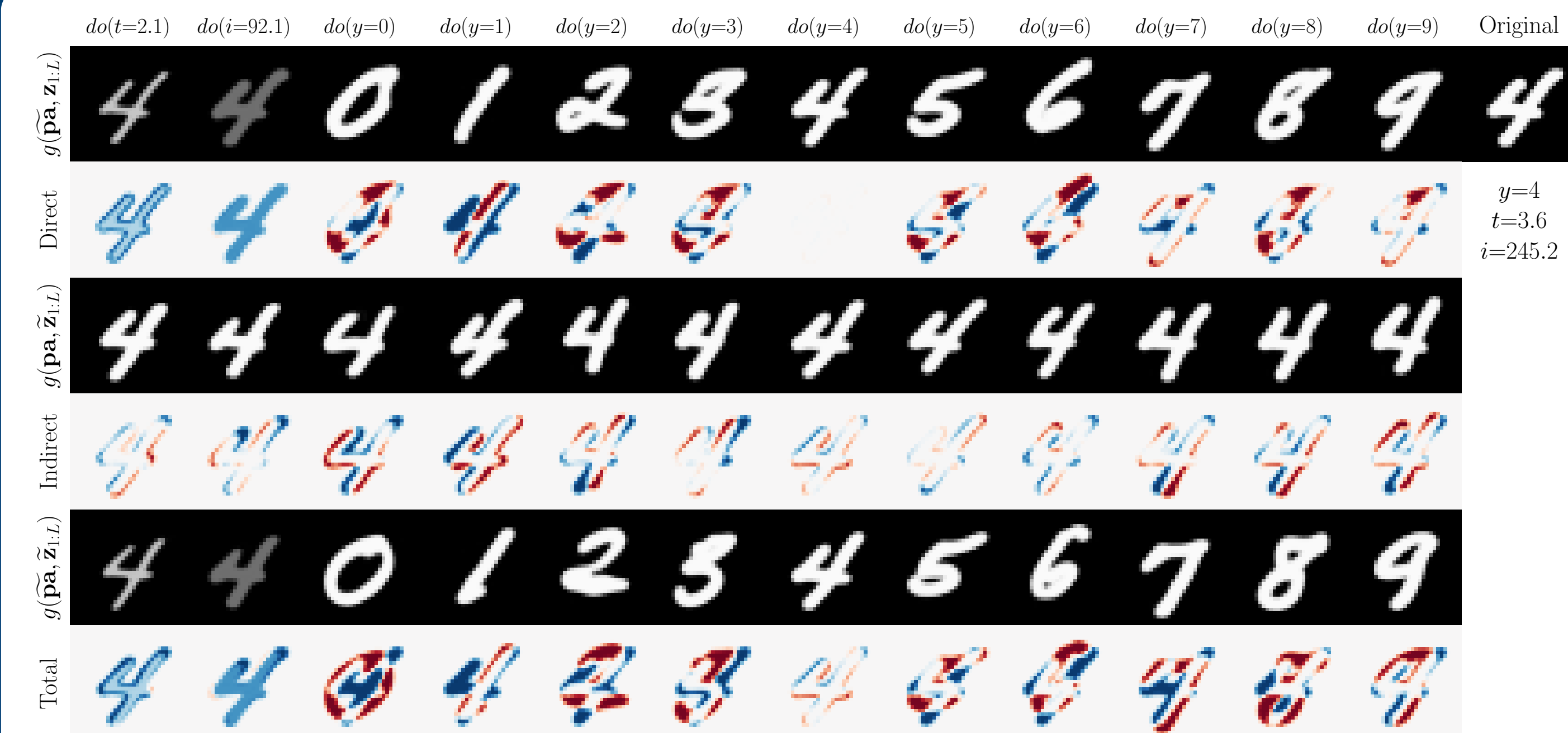
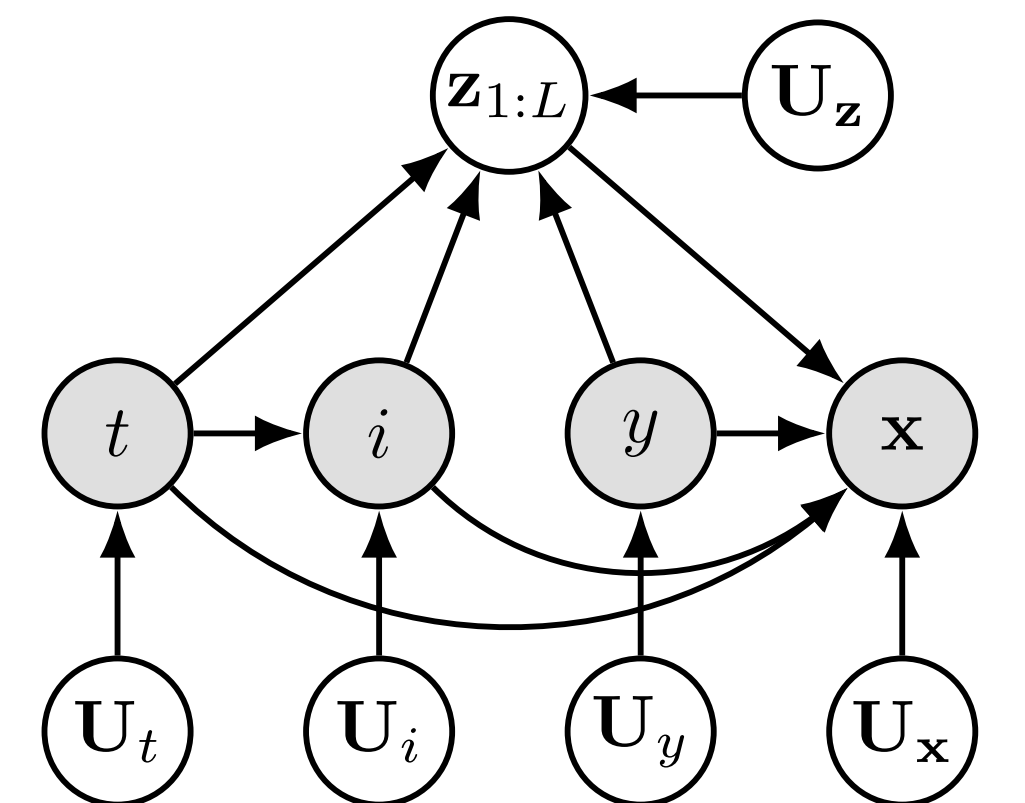


Figure: Morpho-MNIST counterfactuals from our latent mediator SCM. **Direct**, **indirect** and **total** causal effects are shown (red: increase, blue: decrease). **Cross-world counterfactuals** (row 3) are the potential outcome of x given pa_x and the counterfactual mediator $\tilde{z}_{1:L}$ we could've observed had $pa_x := \tilde{pa}_x$.



6. Brain Imaging (UK Biobank)

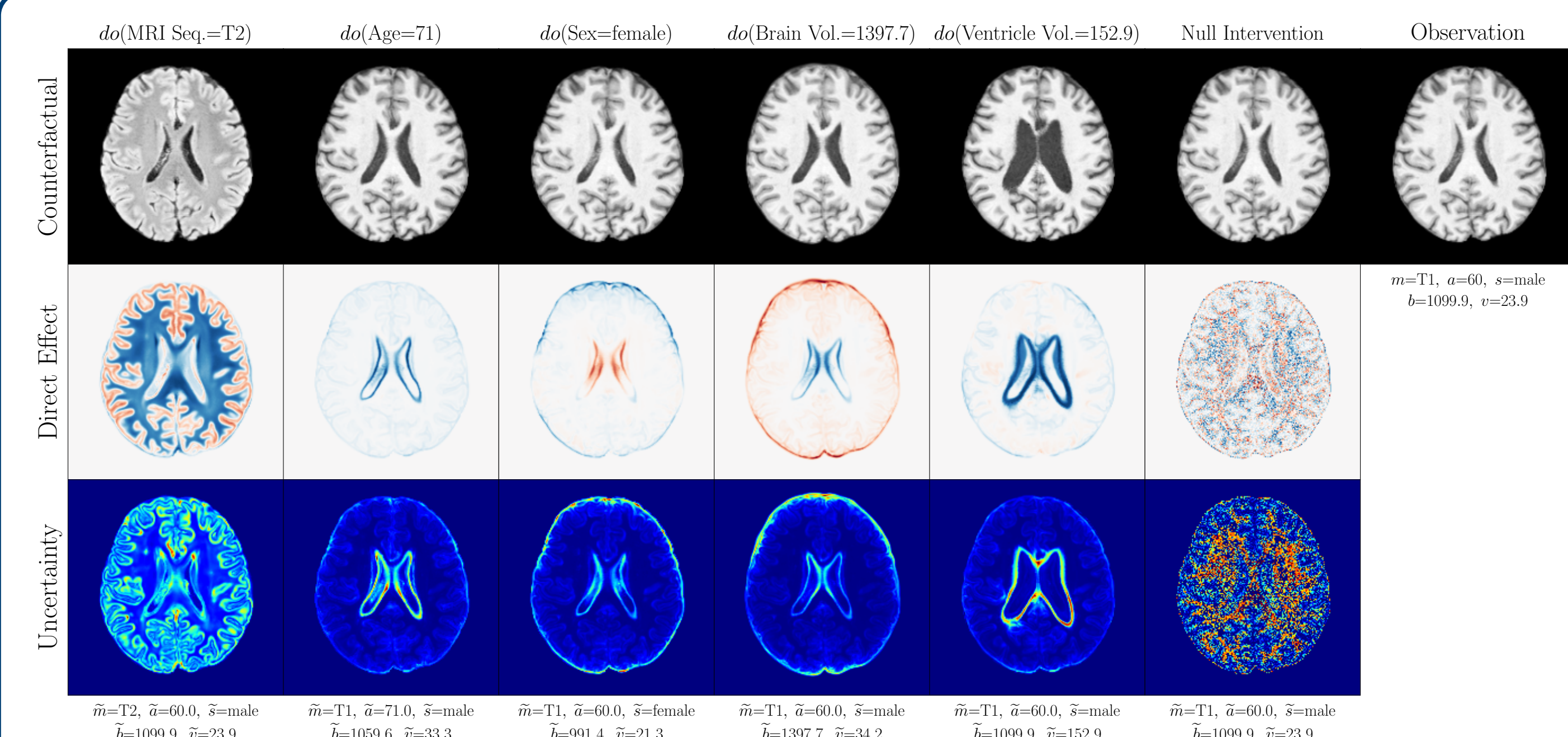
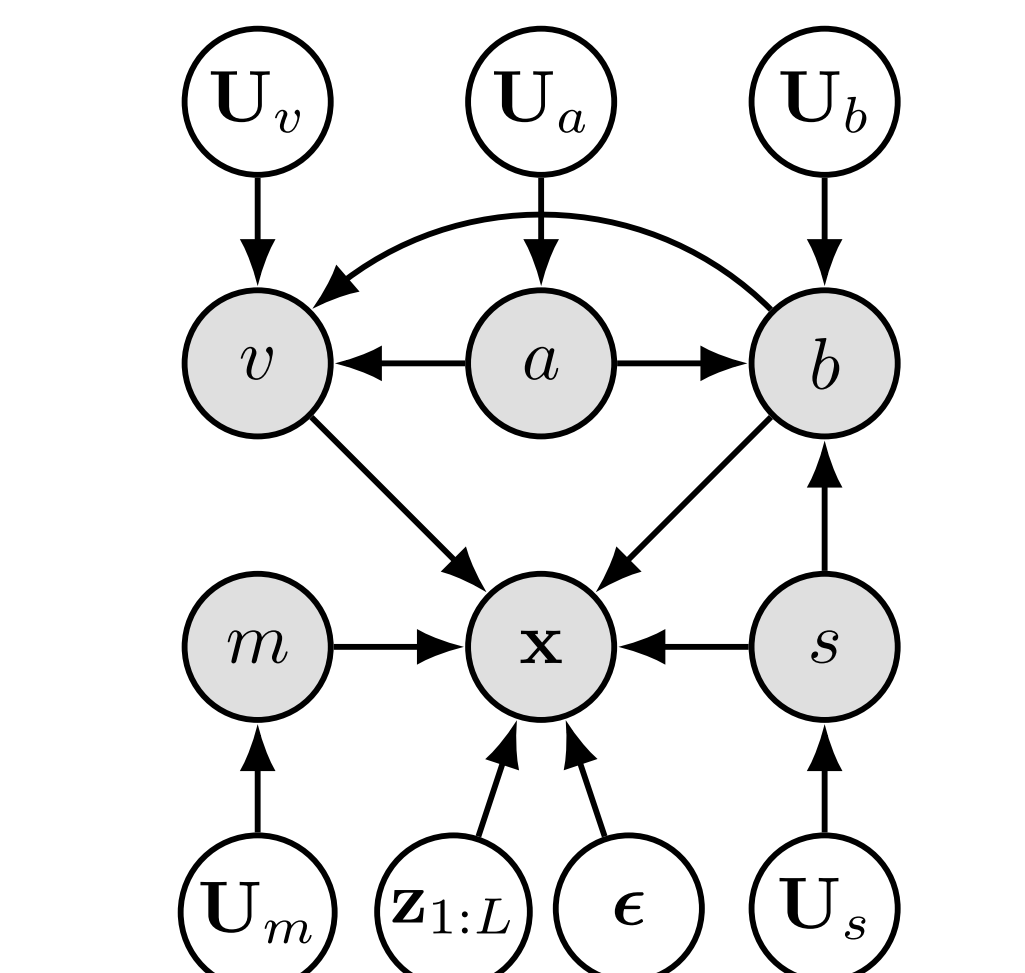


Figure: Brain MRI counterfactuals from our SCM (exogenous prior). Subject identity is preserved. **Counterfactual uncertainty** from stochastic abduction.



7. Chest X-ray (MIMIC)

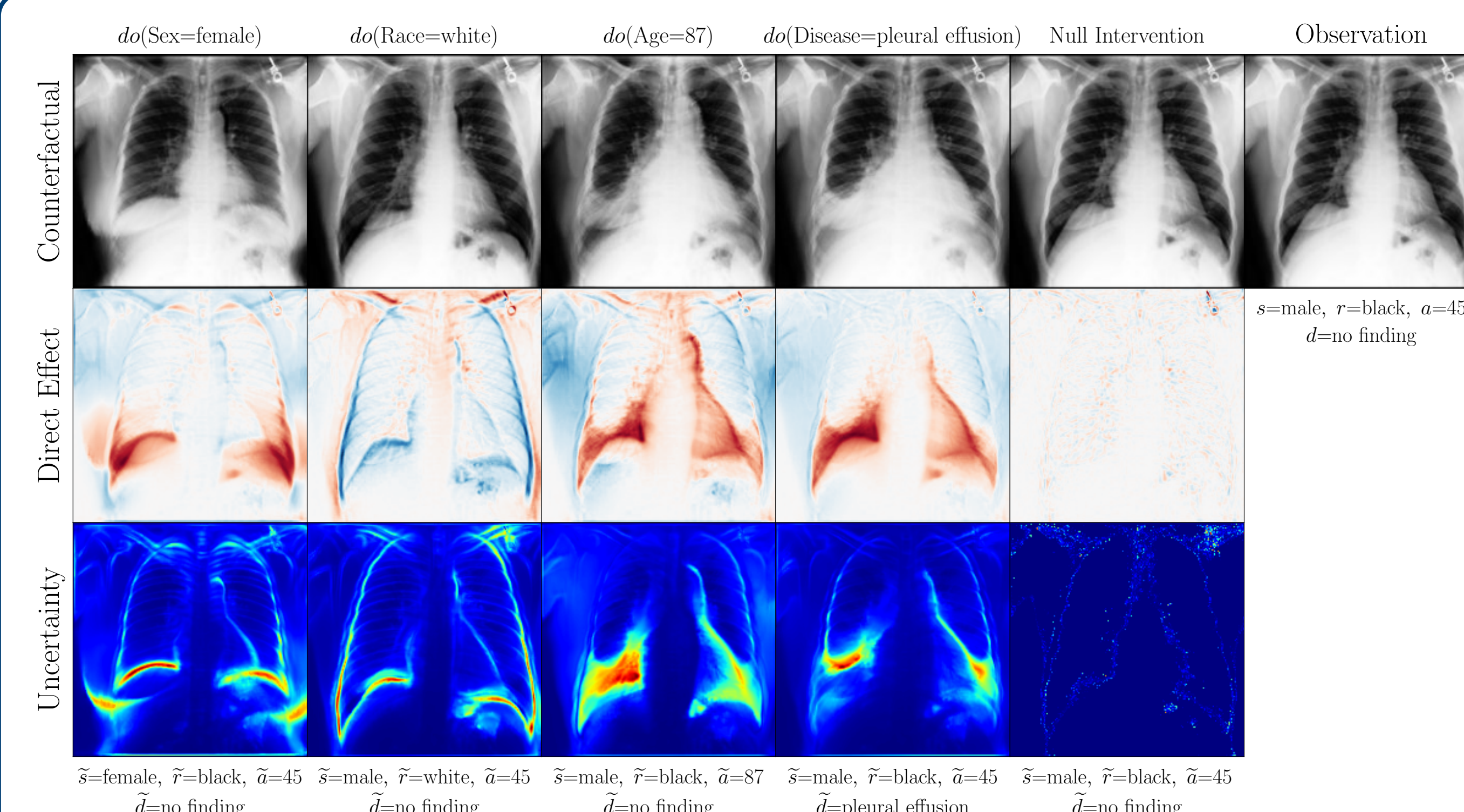


Figure: Chest X-ray counterfactuals from our SCM. **Localized interventional changes** which respect the causal graph (a) and preserve subject identity.

