# Generative Counterfactual Augmentation for Bias Mitigation

Jason Uwaeze, Pranav Kulkarni, Vladimir Braverman, Michael A. Jacobs, Vishwa Parekh

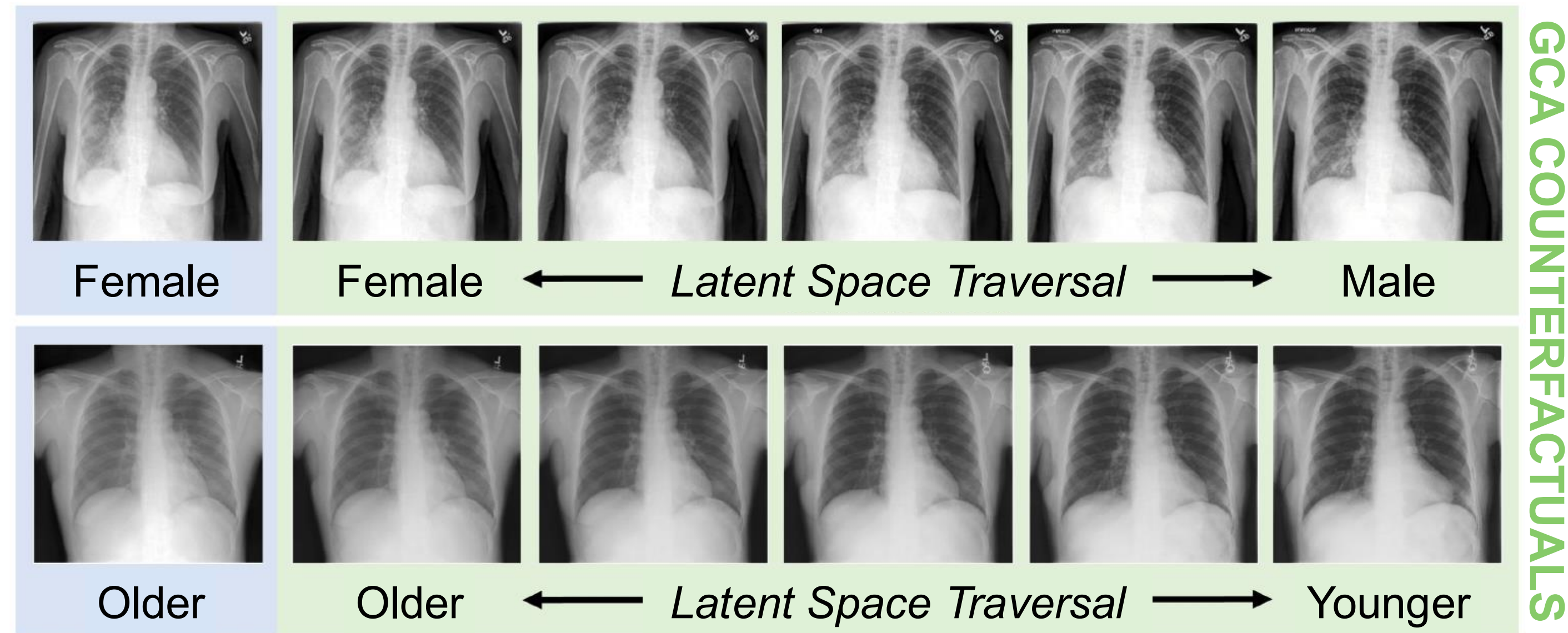ICCV OCT 19-23, 2025 HONOLULU HAWAII

## A new approach to adversarial bias mitigation

**What?** Constructed demographic-complete training data augmentations to mitigate adversarial bias in pneumonia chest x-ray (CXR) **DenseNet121** classifiers
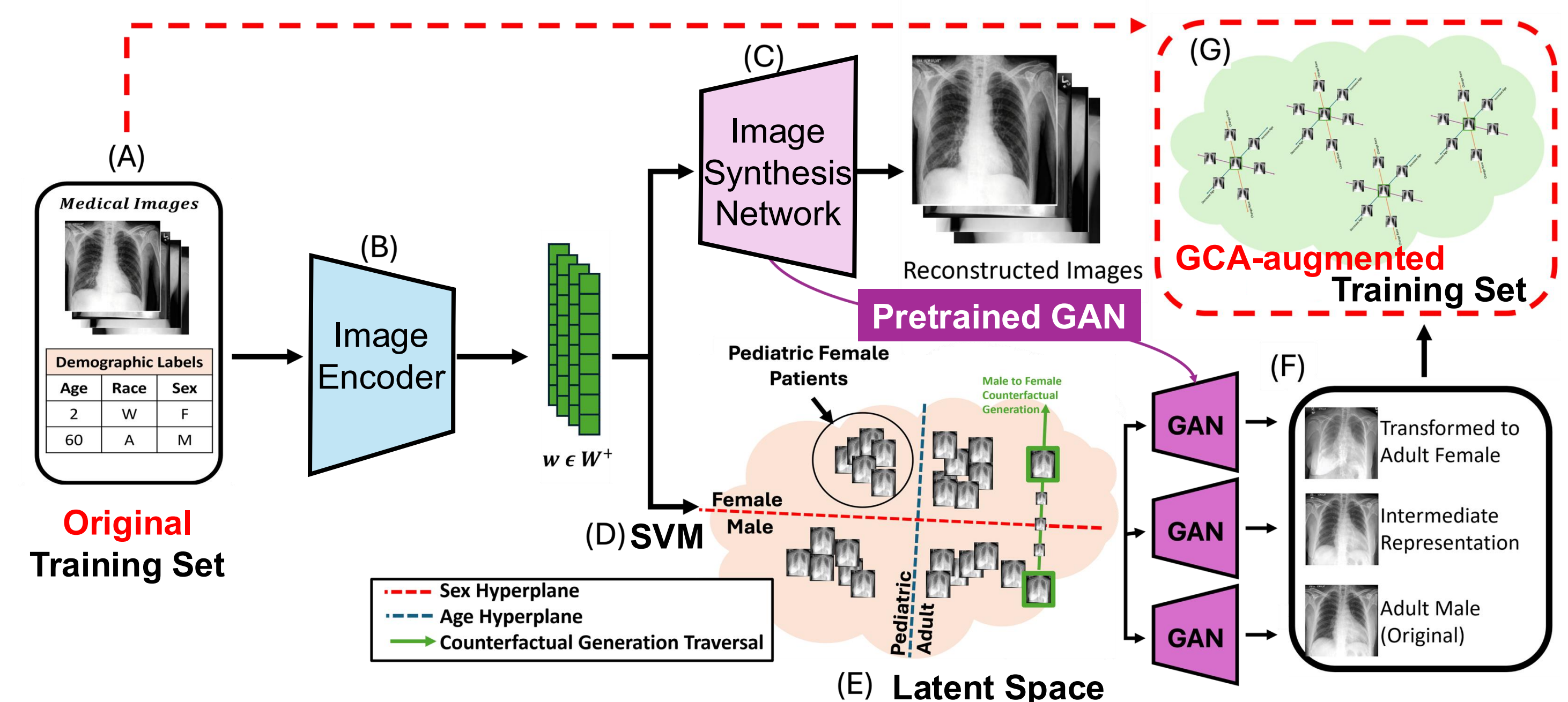
**How?** Generated counterfactuals using StyleGAN3 synthesis network and SVM-guided latent space traversals for **sex** and **age** attributes

ORIGINAL / GCA COUNTERFACTUALS

Female — Female ← *Latent Space Traversal* → Male

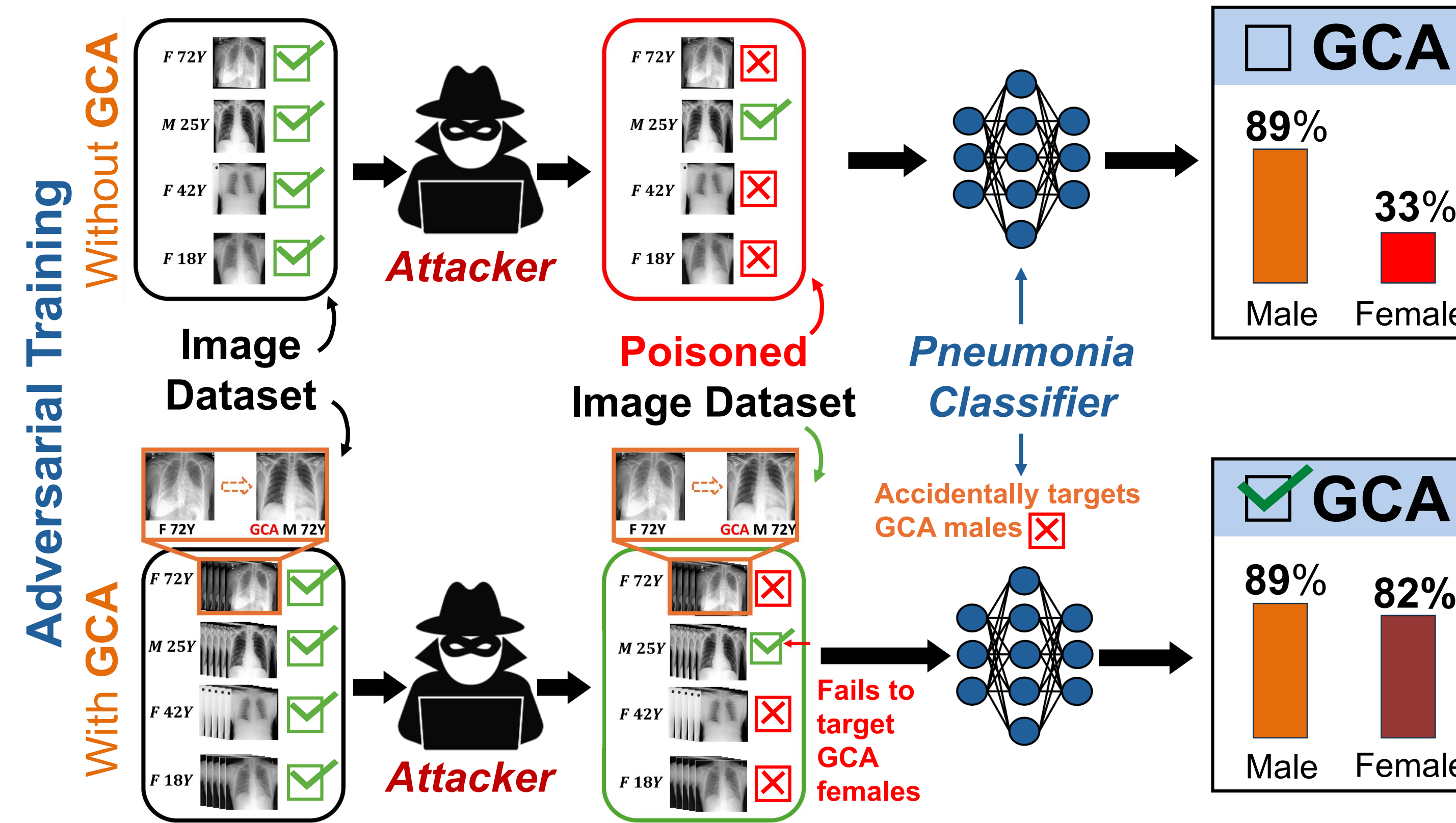Older — Older ← *Latent Space Traversal* → Younger

**Why?**
- Deep learning methods learn protected characteristics (e.g., sex and age) to achieve excellent classification performance
- Potentially amplifying existing systemic disparities in healthcare and worsening patient outcomes
- Thus, we introduce **generative counterfactual augmentations (GCA)**
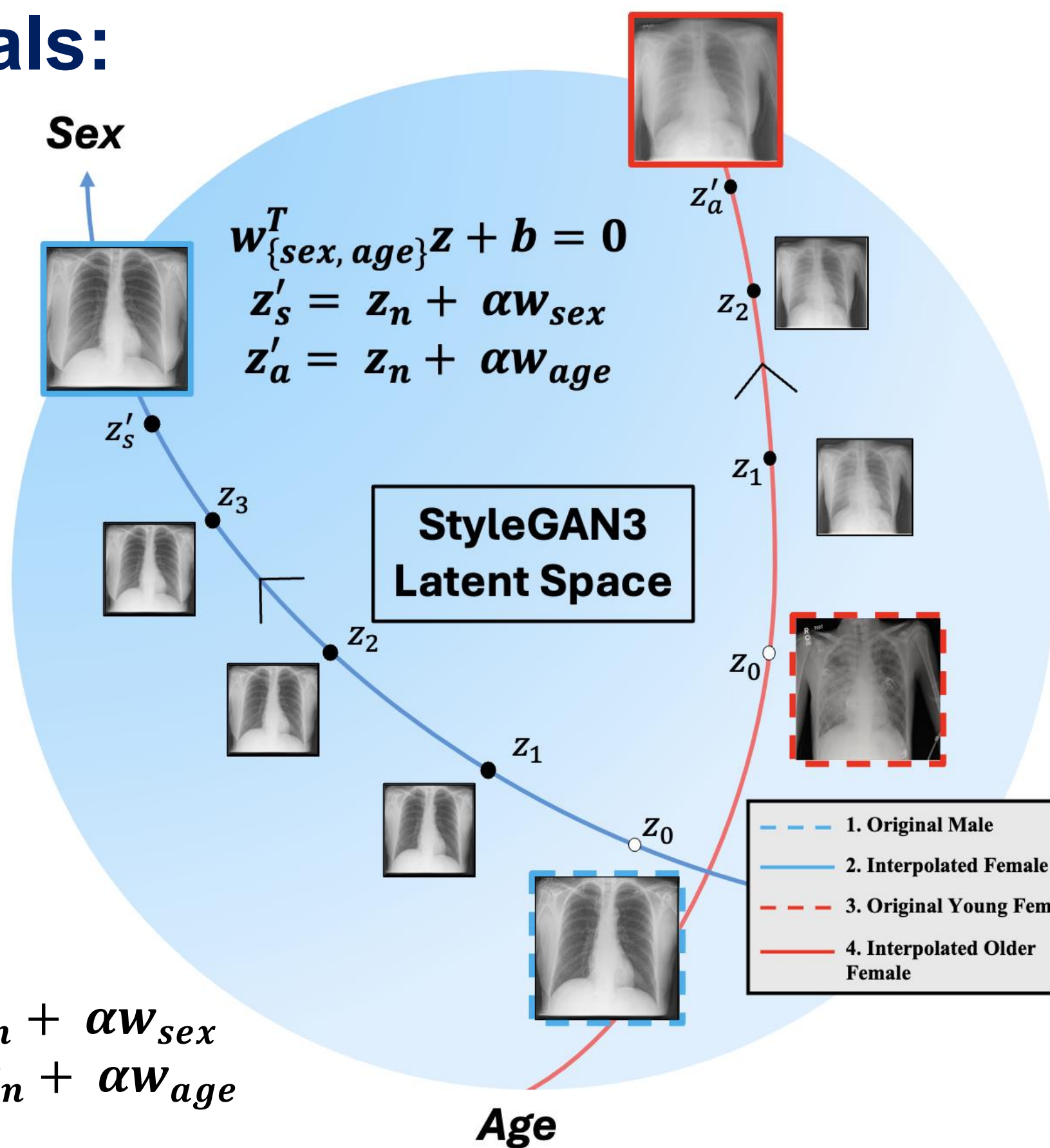
## Generative counterfactual augmentation framework



(A) Medical Images — Demographic Labels (Age, Race, Sex: 2 W F, 60 A M) — **Original Training Set**
(B) Image Encoder — $w \in W^+$
(C) Image Synthesis Network — Reconstructed Images — **Pretrained GAN**
(D) SVM — Sex Hyperplane, Age Hyperplane, Counterfactual Generation Traversal
(E) Latent Space — Pediatric Female Patients, Female/Male, Pediatric Adult
(F) GAN — Transformed to Adult Female / Intermediate Representation / Adult Male (Original)
(G) GCA-augmented **Training Set**

## Medical image classifiers are susceptible to adversarial bias attacks:



**Without GCA** (Adversarial Training): Image Dataset → *Attacker* → Poisoned Image Dataset → Pneumonia Classifier → ☐ GCA: Male 89%, Female 33%

**With GCA** (Adversarial Training): Accidentally targets GCA males ✗ → *Attacker* → Fails to target GCA females → ☑ GCA: Male 89%, Female 82%
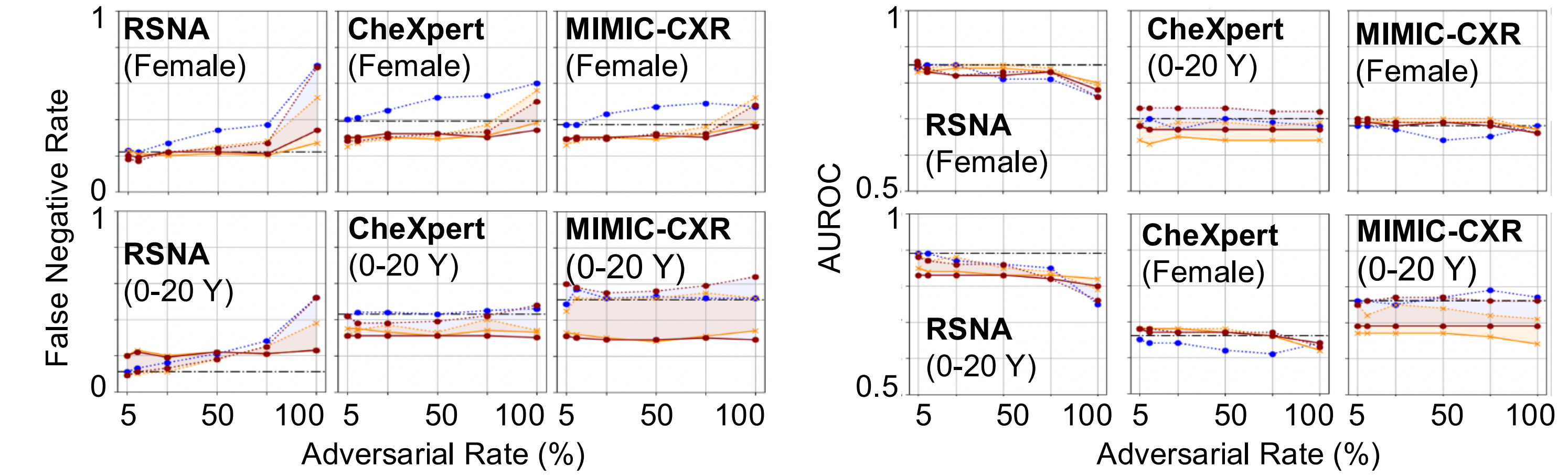
## Latent space interpolations for generating CXR counterfactuals:

- Train SVM classifiers in the StyleGAN3 latent space
- Use learned SVM hyperplanes to guide latent interpolations across **sex** and **age** attributes
- Let $w_{sex}$ and $w_{age}$ denote normal vectors separating demographic groups
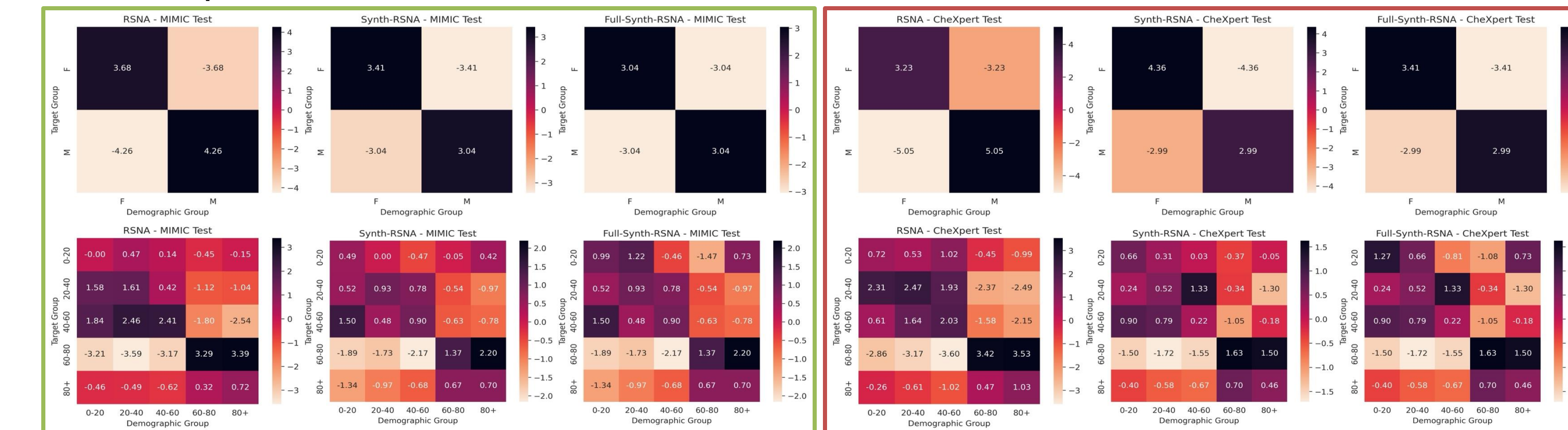- Interpolate by traversing perpendicular to each hyperplane:
  - **Sex interpolation**: $z'_s = z_n + \alpha w_{sex}$
  - **Age interpolation**: $z'_a = z_n + \alpha w_{age}$



$$w_{\{sex, age\}}^T z + b = 0$$
$$z'_s = z_n + \alpha w_{sex}$$
$$z'_a = z_n + \alpha w_{age}$$

StyleGAN3 Latent Space

Legend:
1. Original Male
2. Interpolated Female
3. Original Young Female
4. Interpolated Older Female

## GCA lowers false negative rates (FNR) in CXR datasets while achieving high AUROC:



Legend: Baseline, RSNA, Synth-RSNA, Synth-Synth (Overall), Full-RSNA, Full-Synth-RSNA (Overall)

FNR plots: RSNA (Female), CheXpert (Female), MIMIC-CXR (Female), RSNA (0-20 Y), CheXpert (0-20 Y), MIMIC-CXR (0-20 Y) — Adversarial Rate (%)

AUROC plots: RSNA (Female), CheXpert (Female), MIMIC-CXR (Female), RSNA (0-20 Y), CheXpert (0-20 Y), MIMIC-CXR (0-20 Y) — Adversarial Rate (%)

## GCA reduces vulnerability of the targeted and non-targeted subgroups:

- Models trained on **(1)** original, **(2)** demographic-specific synthetic, and **(3)** demographic-complete synthetic RSNA datasets
- Models tested on **RSNA**, **MIMIC-CXR**, and **CheXpert** datasets
- Impact of GCA remained consistent across external MIMIC-CXR and CheXpert datasets



## Summary/Conclusion:

- We propose GCA, a counterfactual demographic-complete data augmentation method to mitigate adversarial bias in deep learning.
- Using controlled injections of underdiagnosis bias across age and sex groups, we show that GCA reduces FNR disparities, preserves high AUROC, and generalizes across external datasets.
- **Future work**: to adapt GCA for other imaging modalities and tasks, ensuring trustworthiness in real-world clinical settings.

*Paper & Code* — Scan Me