





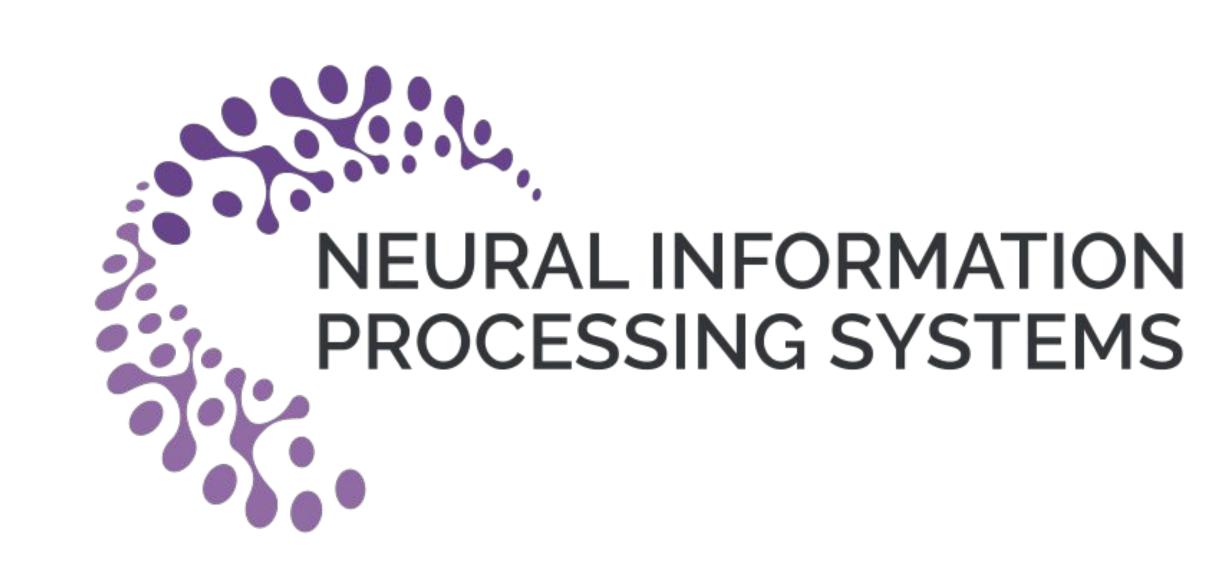


Reliable Visual Counterfactual Explanation Factual Using Conditional Flow Matching

Zhuo Cao¹, Xuan Zhao^{1*}, Lena Krieger^{1,2*}, Hanno Scharr¹, Ira Assent^{1,3}

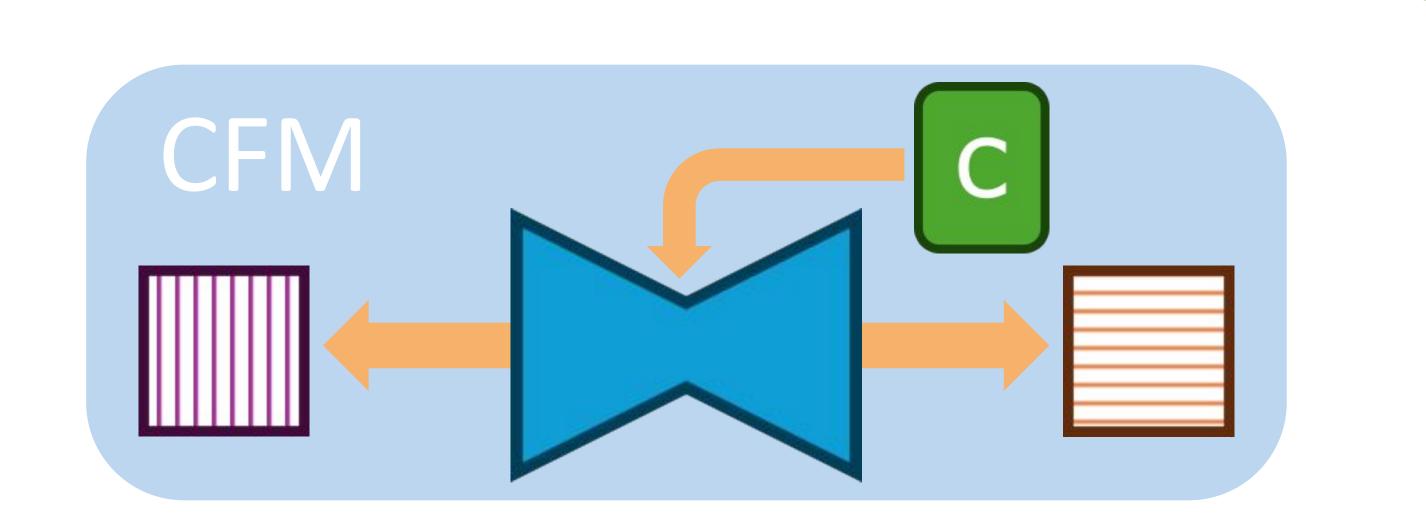
²Munich Center for Machine Learning (MCML), LMU Munich, Germany ¹IAS-8, Forschungszentrum Jülich, Germany ³Aarhus University, Denmark



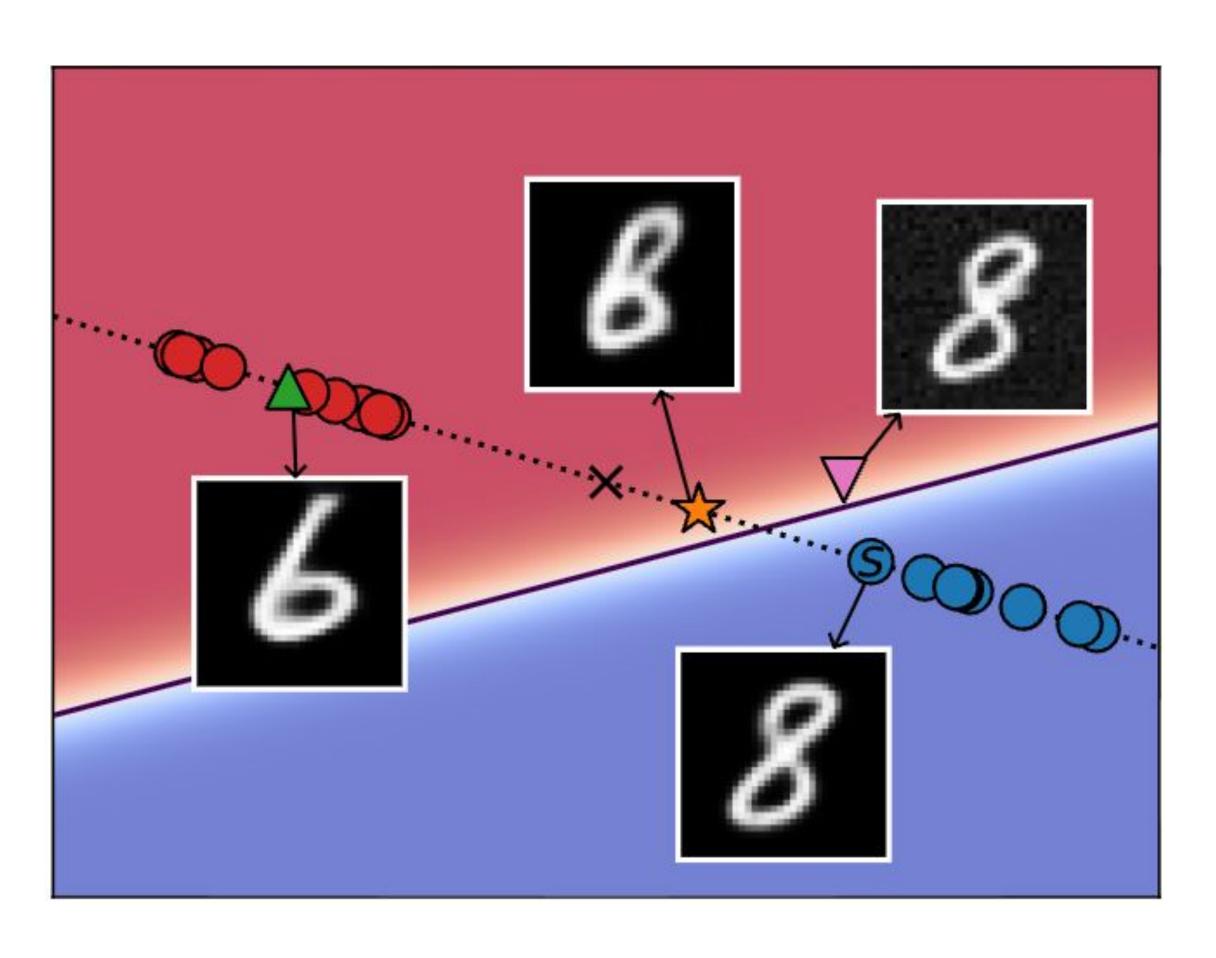


TL; DR

We propose LeapFactual. It offers a flexible, model-agnostic mechanism for generating counterfactual explanations, capable of producing not only high-quality but also reliable results, even in the presence of discrepancies between true and learned decision boundaries.

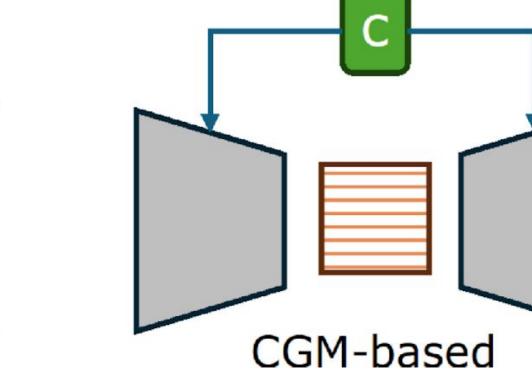


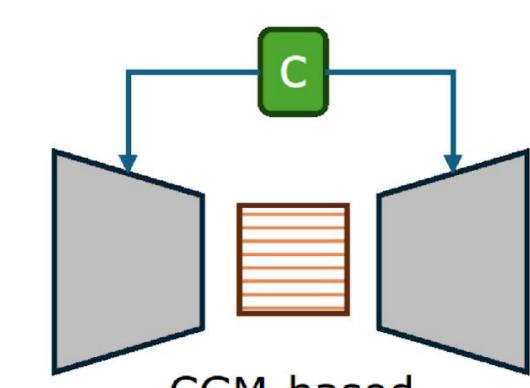
Counterfactual Explanations (CE) -

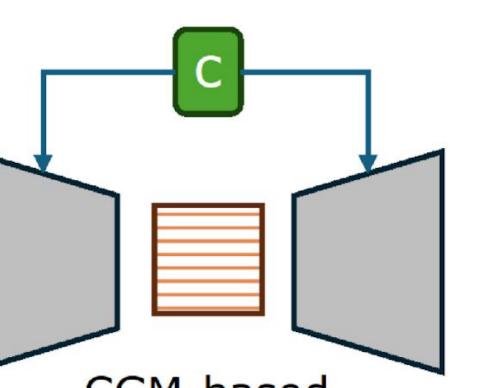


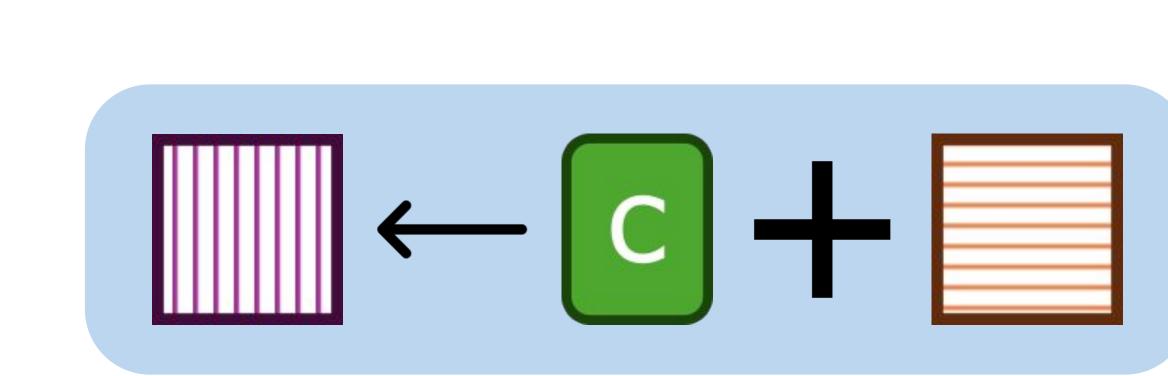
Counterfactual Explanation (CE) seeks to make only semantically meaningful modifications to an input image in order to obtain a similar image with a target label prediction outcome.











Optimization (Opt)-based Methods:

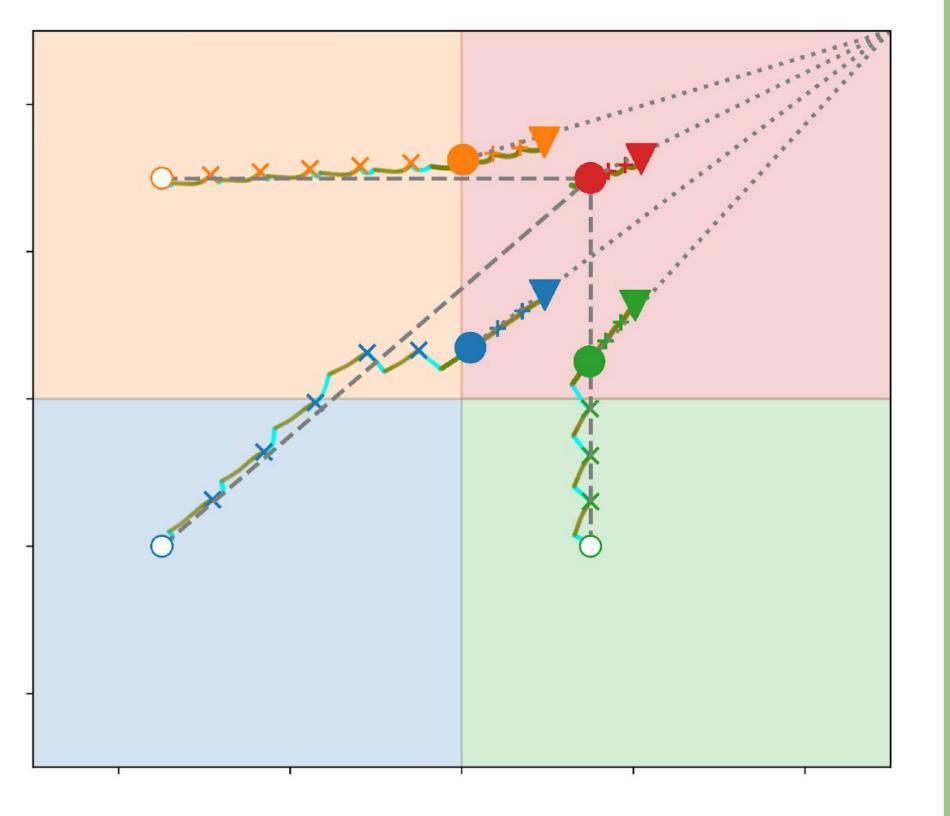
- Stop once the model's prediction matches the target
- Require a differentiable classifier Conditional Generative Model (CGM)-based Methods:
- Without control over boundary distance
- Need fine-tuning whenever the classifier changes

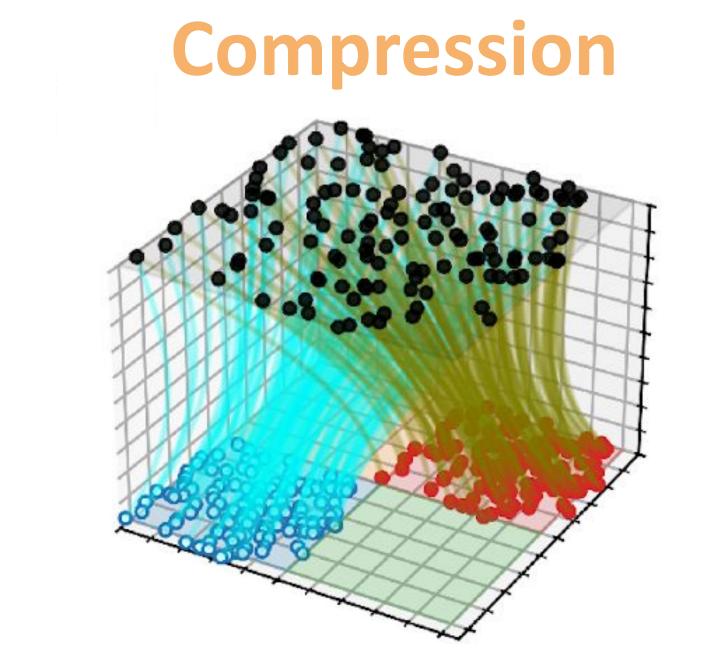
_LeapFactual -

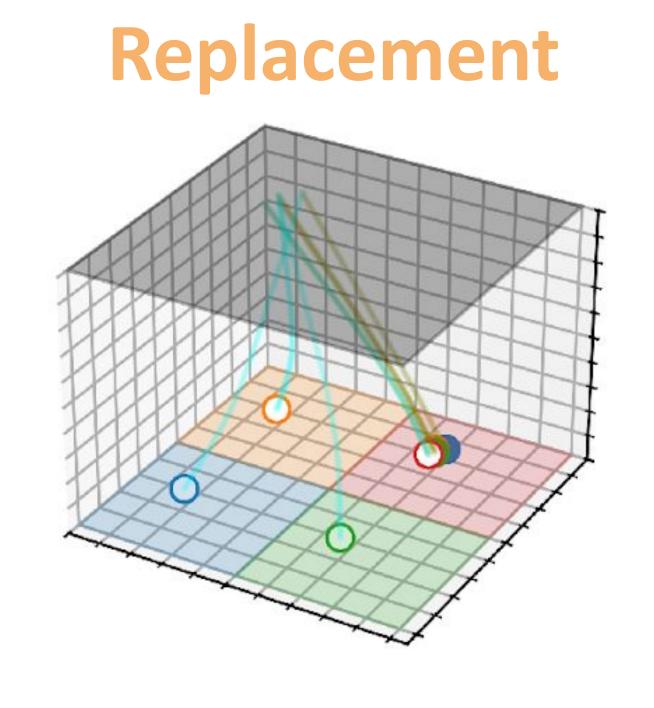
Training Phase: We demonstrate that the classifier's output is disentangled from the residual by training with our proposed Counterfactual Explanation—Conditional Flow Matching (CE-CFM) objective:

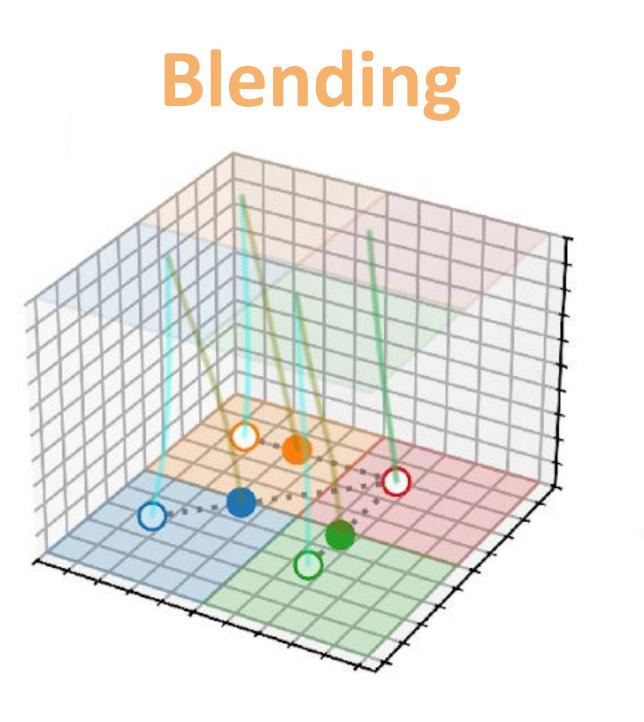
$$\mathcal{L}(\psi) \coloneqq \mathbb{E}_{t,q(h),p_t(Z|h)} \|v_{\psi}(t,z,\mathbf{C}) - u_t(z|h)\|$$

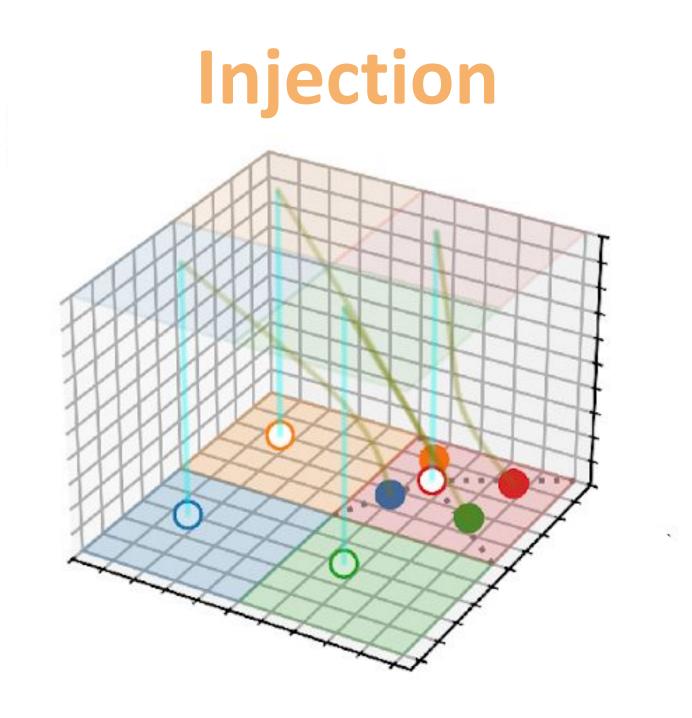
Explaining Phase: We propose two algorithms, LEAP and LeapFactual. Leveraging the flexibility of CFM, LEAP enables information replacement, blending, and injection. LeapFactual further integrates blending and injection leaps to able counterfactual explanations.





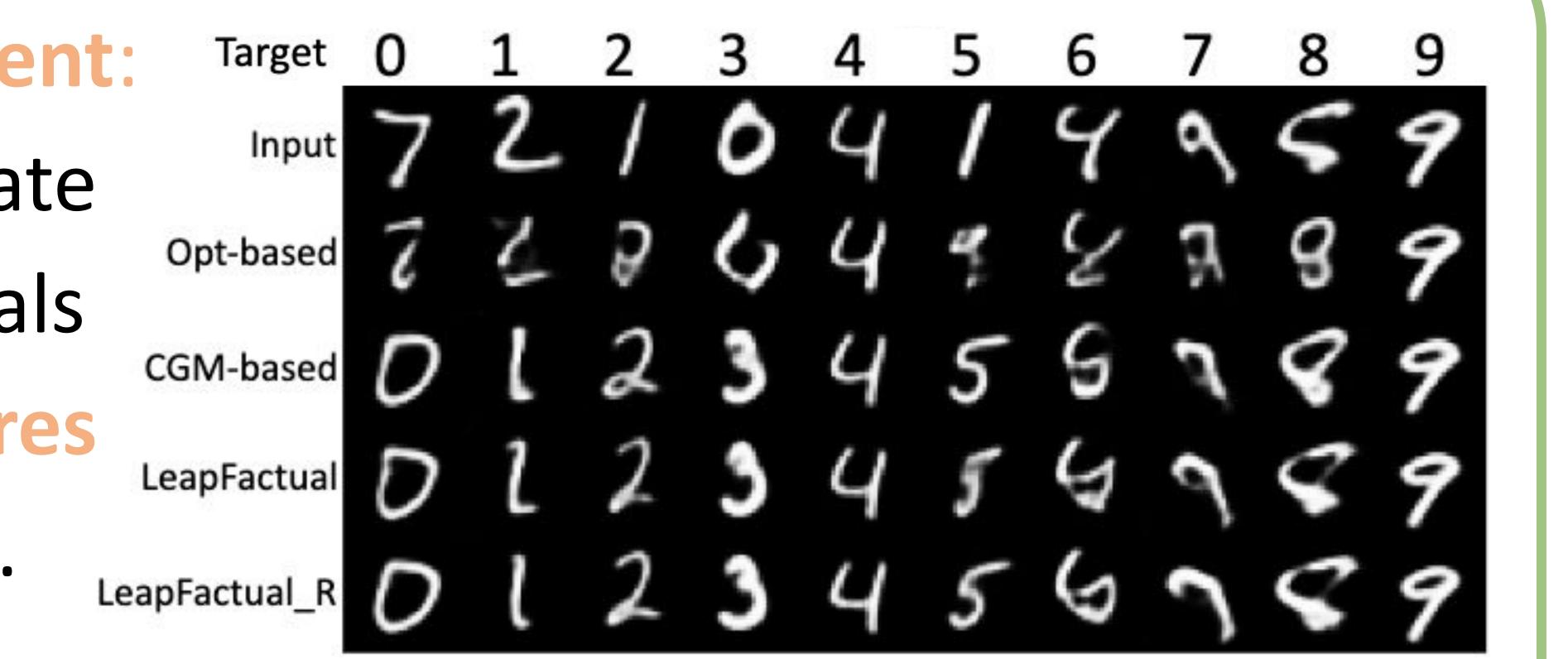






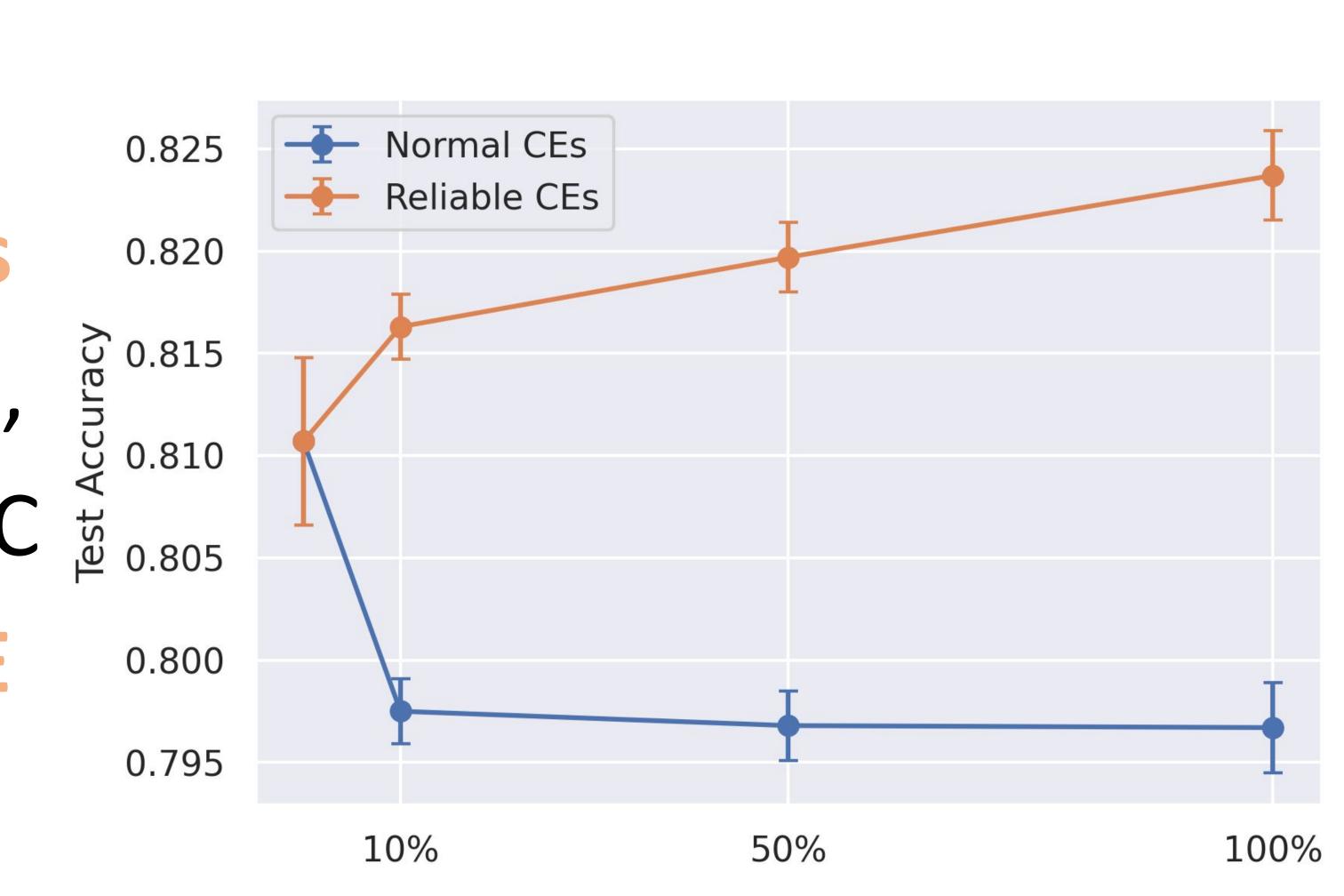
Experiments

Quantitative Assessment: Our methods generate realistic counterfactuals by altering class features while preserving style.



Model Improvement:

Model performance with standard CE blending, \(\frac{5}{2} \) 0.815 but both accuracy and AUC by 0.805 improve with reliable CE 0.800 blending.



Generalization: We show that the proposed method is highly scalable to high-res images and non-differentiable classifiers.













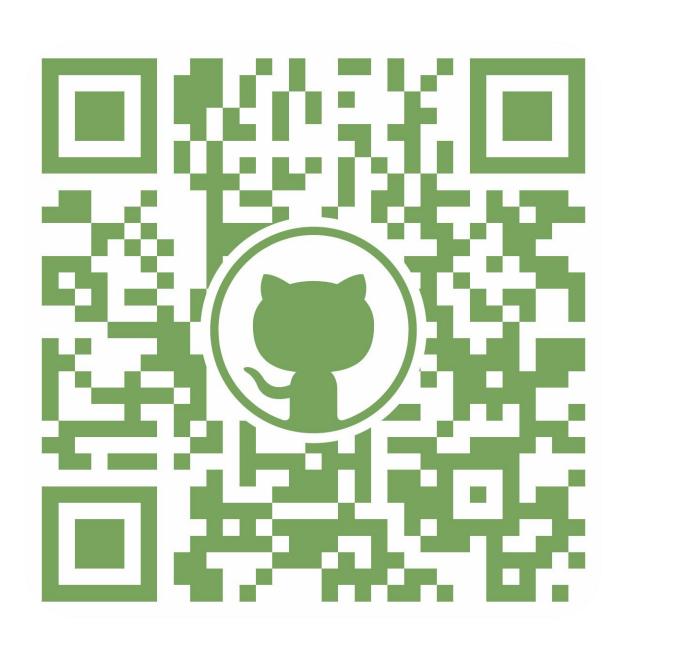


Reliable Visual Counterfactual Explanation Factual Using Conditional Flow Matching

Zhuo Cao¹, Xuan Zhao^{1*}, Lena Krieger^{1,2*}, Hanno Scharr¹, Ira Assent^{1,3}

¹IAS-8, Forschungszentrum Jülich, Germany

²Munich Center for Machine Learning (MCML), LMU Munich, Germany ³Aarhus University, Denmark

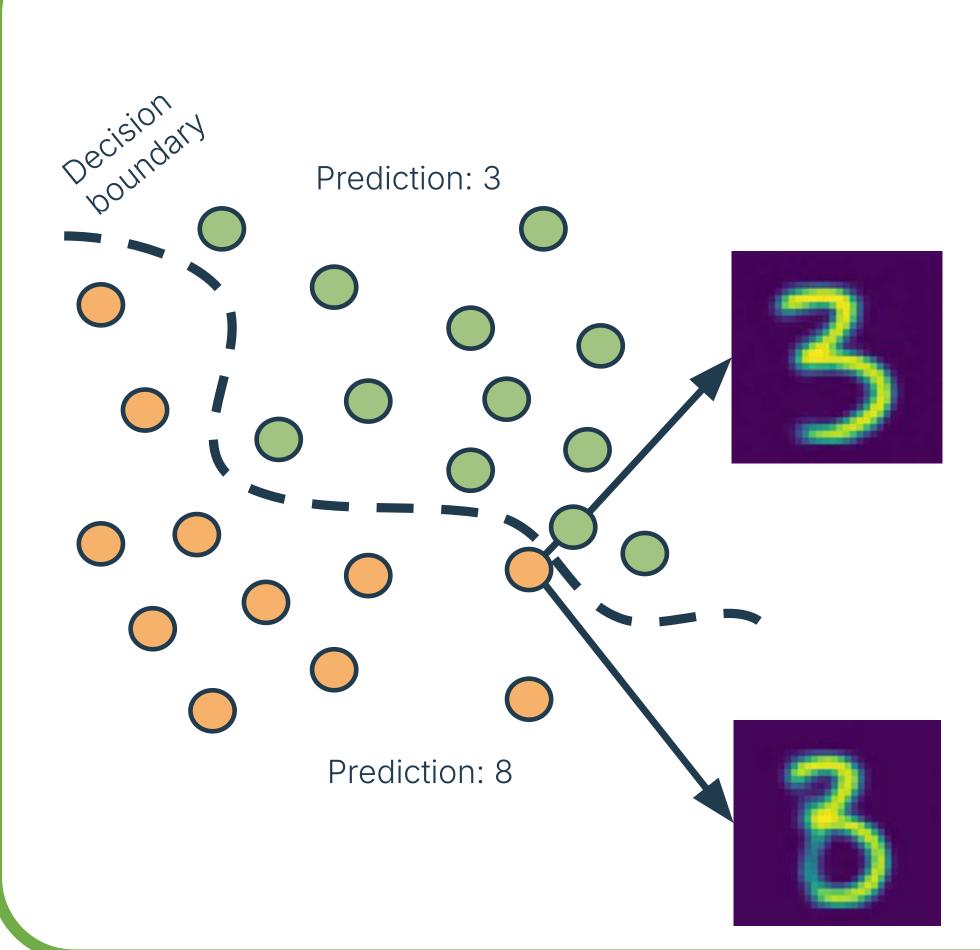




TL; DR-

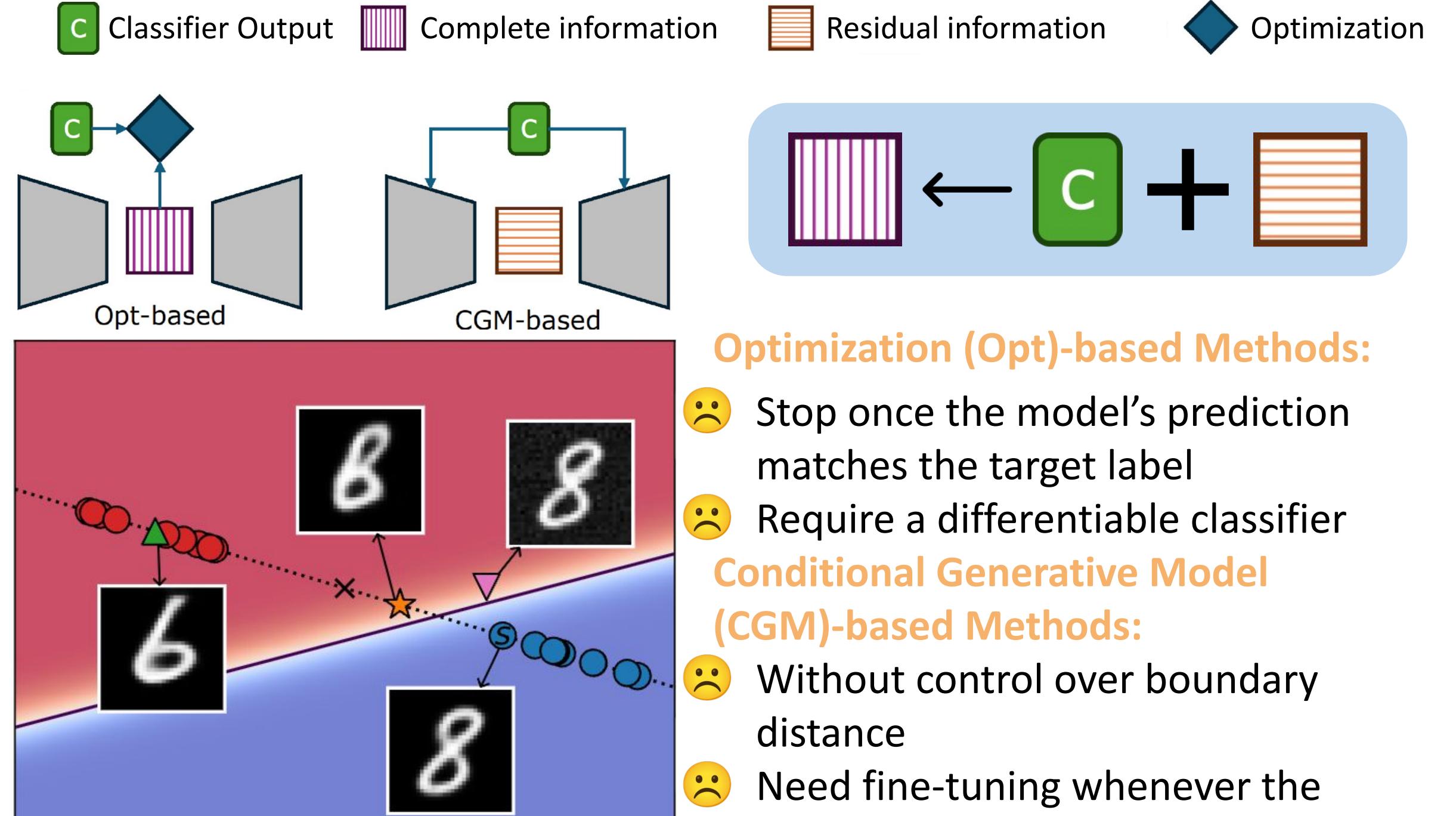
eapFactual. It offers a flexible, mod c mechanism for We propose Lea generating counterfactual explanations, capable of producing not only e results, even in the presence of discrepancies y but also r between true and learned decision boundaries.

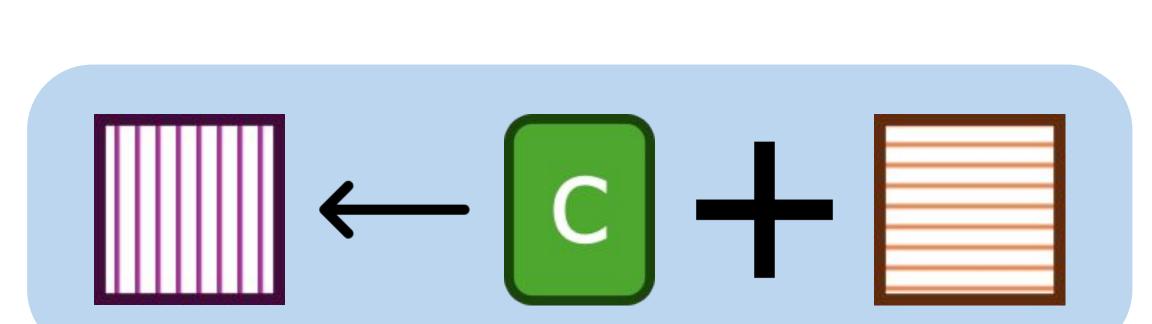
-What is Counterfactual Explanation (CE)



Counterfactual explanation seeks to make only semantically meaningful modifications to an input image in order to obtain a similar image prediction outcome. It exposes decision boundaries more clearly and can also be used to generate new, informative s that improve generalization and fairness.

The Gap: Existing CE methods have limitations-

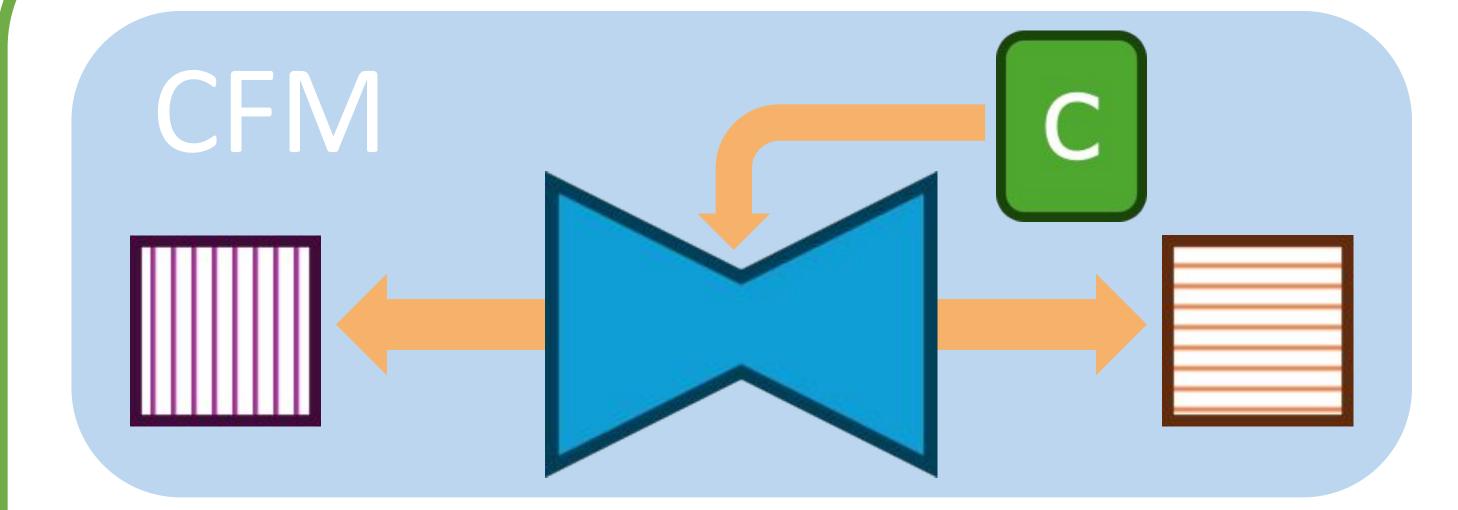




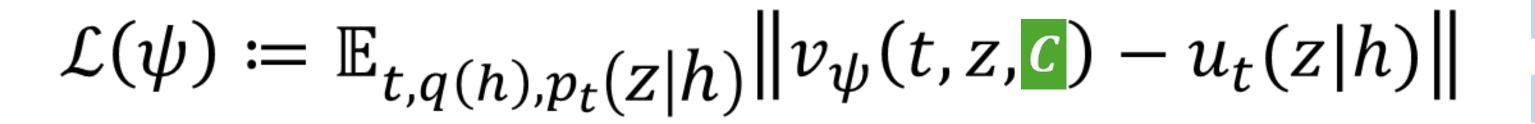
- Stop once the model's prediction matches the target label
- Require a differentiable classifier anditional Generative Model (CGM)-based Methods:
- Without control over boundary distance
- Need fine-tuning whenever the classifier changes

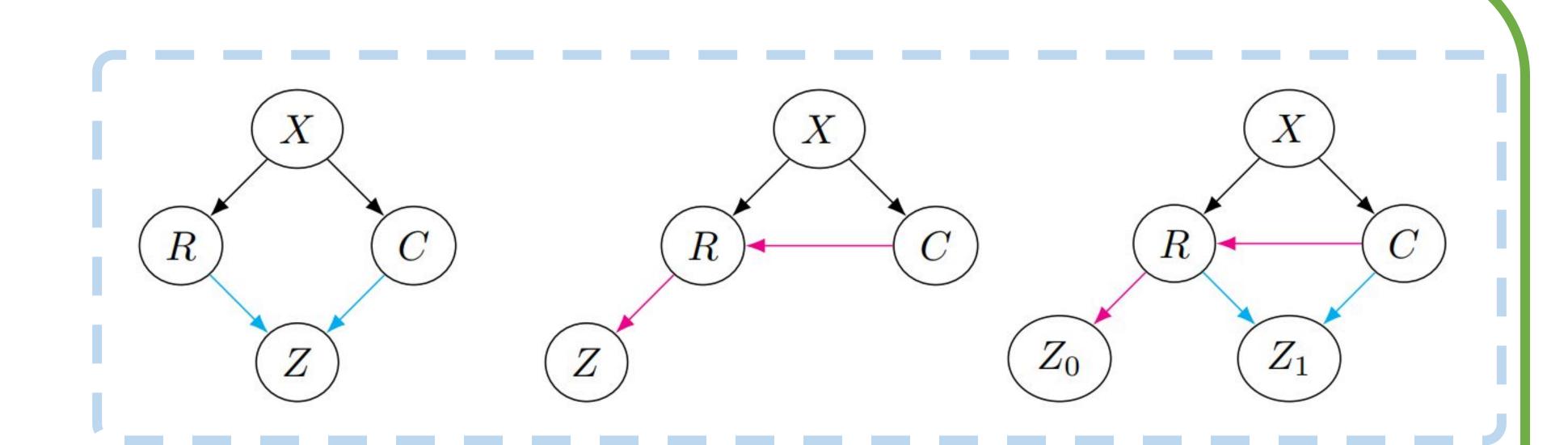
LeapFactual

Compression



Training Phase: We demonstrate that the classifier's output is disentangled from the residual by training with our proposed Counterfactual Explanation—Conditional Flow Matching (CE-CFM) objective:





: We propose two algorithms, LEAP and LEAPFACTUAL. Leveraging the flexibility of CFM, LEAP enables information replacement, blending, and injection. LEAPFACTUAL further integrates blending and injection leaps to generate r e counterfactual explanations.

Algorithm 1 LEAP

 $z \leftarrow \text{LEAP}(v_{\psi}, z, y_c, \hat{y}_c, \gamma_b, \gamma_b)$

 $z \leftarrow \text{LEAP}(v_{\psi}, z, y_{c}, \hat{y}_{c}, \gamma_{i, \text{lift}}, \gamma_{i, \text{land}})$

 $g_c \leftarrow f_{\theta}(g_{\phi}(z))$

Step 3: (Optional) Information Injection

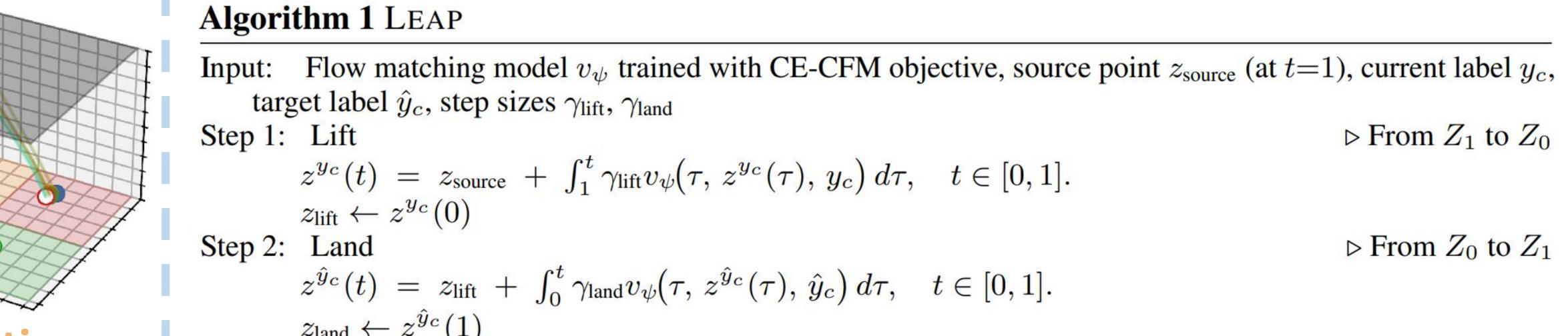
 $g_c \leftarrow f_{\theta}(g_{\phi}(z))$

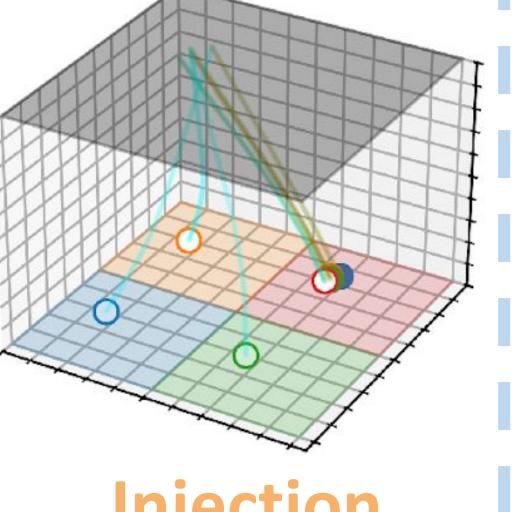
for j = 0 to $N_i - 1$ do

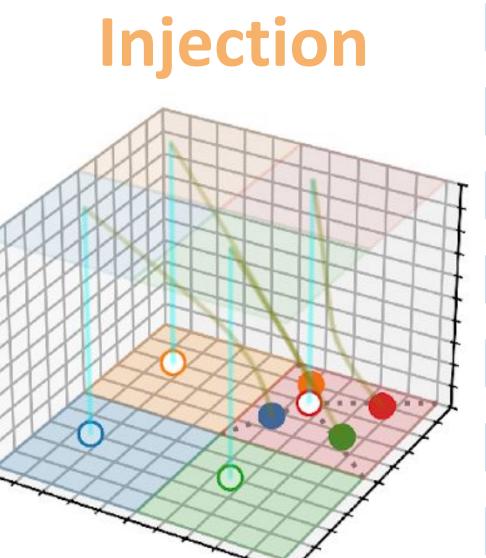
Step 4: Postprocessing

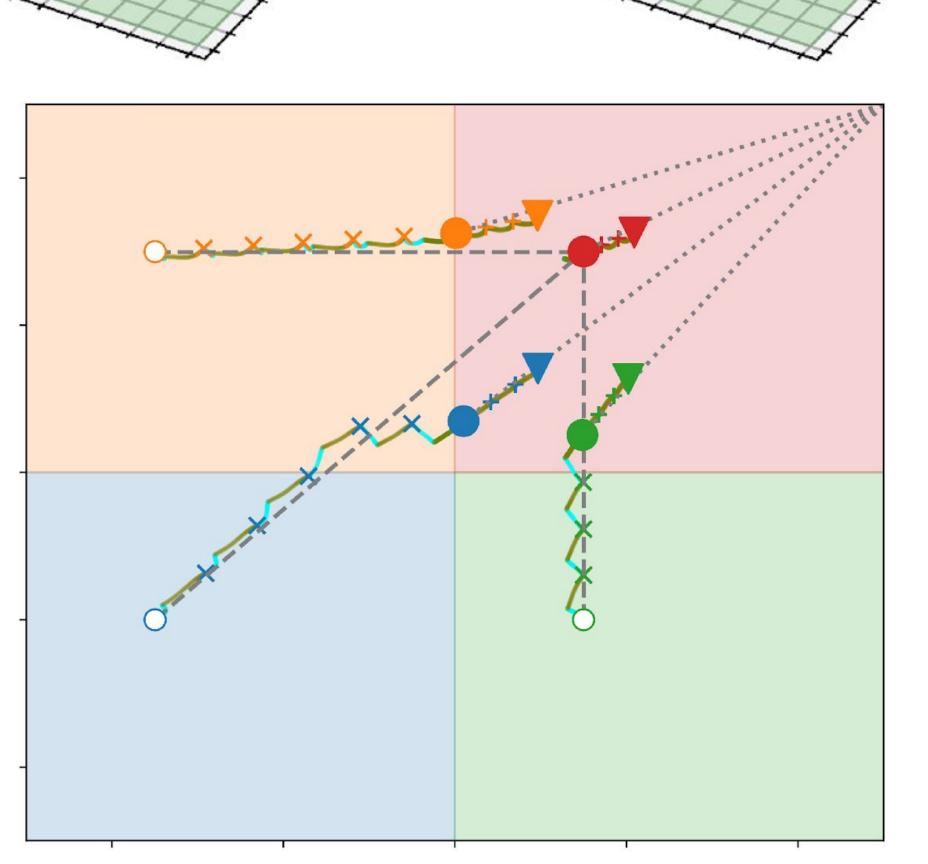
 $x_{\rm CE}=g_{\phi}(z)$

Output: Transported point x_{CE}







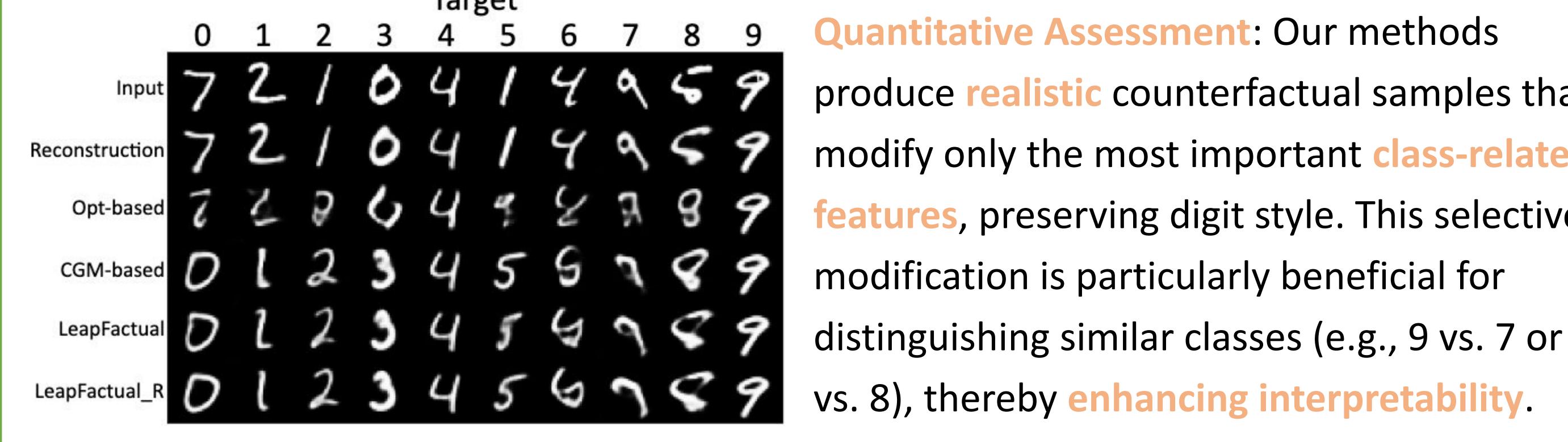


target label \hat{y}_c , step sizes γ_{lift} , γ_{land} \triangleright From Z_1 to Z_0 $z^{y_c}(t) = z_{\text{source}} + \int_1^t \gamma_{\text{lift}} v_{\psi}(\tau, z^{y_c}(\tau), y_c) d\tau, \quad t \in [0, 1].$ \triangleright From Z_0 to Z_1 $z^{\hat{y}_c}(t) = z_{\text{lift}} + \int_0^t \gamma_{\text{land}} v_{\psi}(\tau, z^{\hat{y}_c}(\tau), \hat{y}_c) d\tau, \quad t \in [0, 1].$ Output: Transported point z_{land} Algorithm 2 LEAPFACTUAL Input: Source point x, target label \hat{y}_c , classifier f_{θ} , generative model g_{ϕ} , flow matching model v_{ψ} trained with CE-CFM objective Hyperparameters: Number of blending leaps N_b , blending step size γ_b , number of injection leaps N_i , injection step sizes $\gamma_{i, lift} < \gamma_{i, land}$ Step 1: Preprocessing Determining the current label $y_{\rm c} \leftarrow f_{\theta}(g_{\phi}(z))$ Step 2: Information Blending for j = 0 to $N_b - 1$ do

▶ Blending the source and target classes information

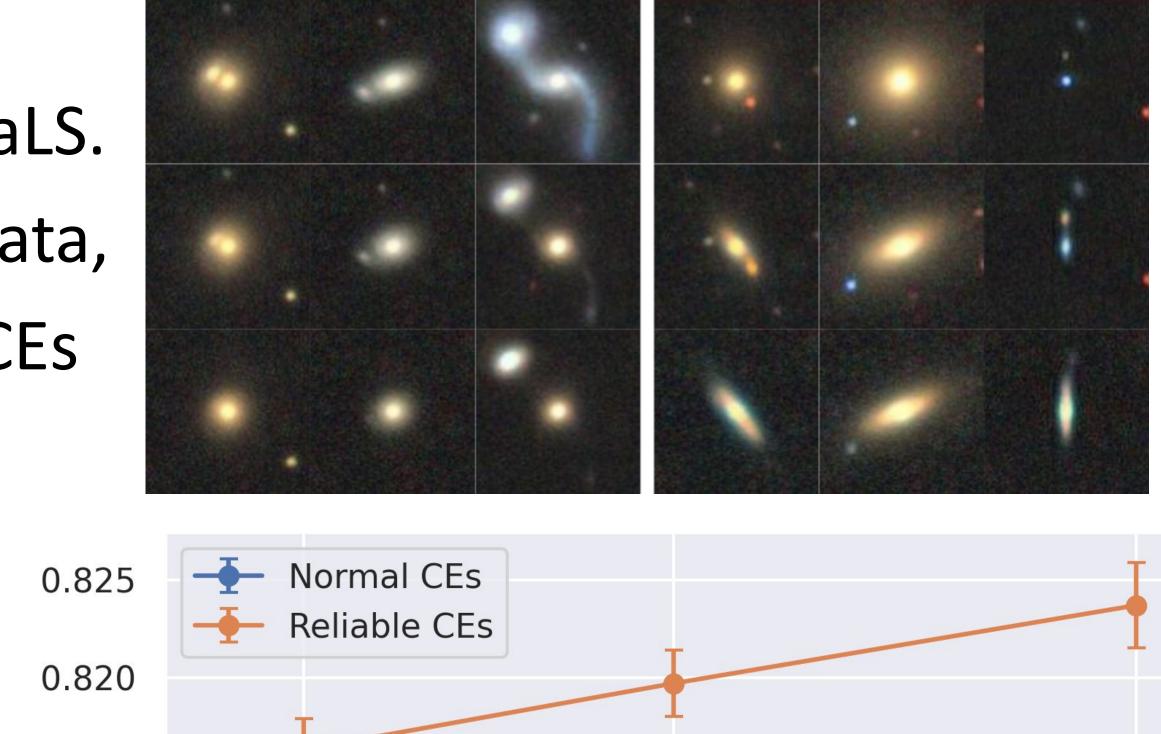
□ Generating Reliable CE

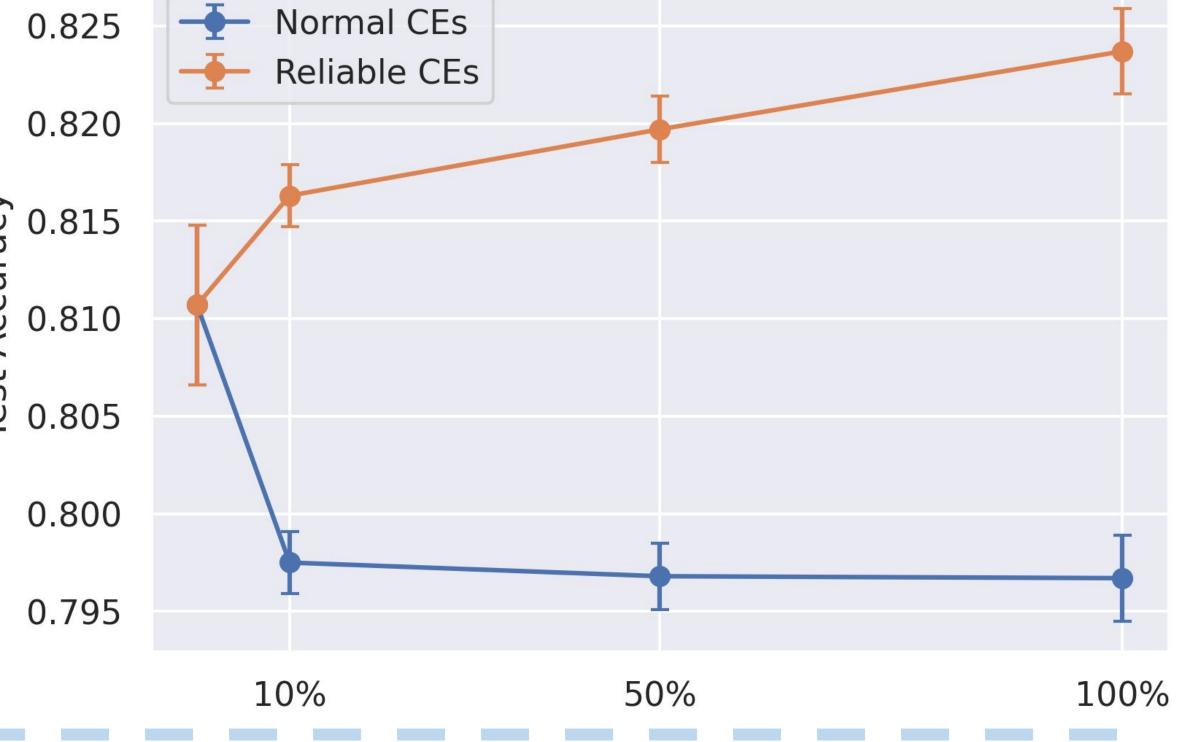
▶ Injecting the target class information

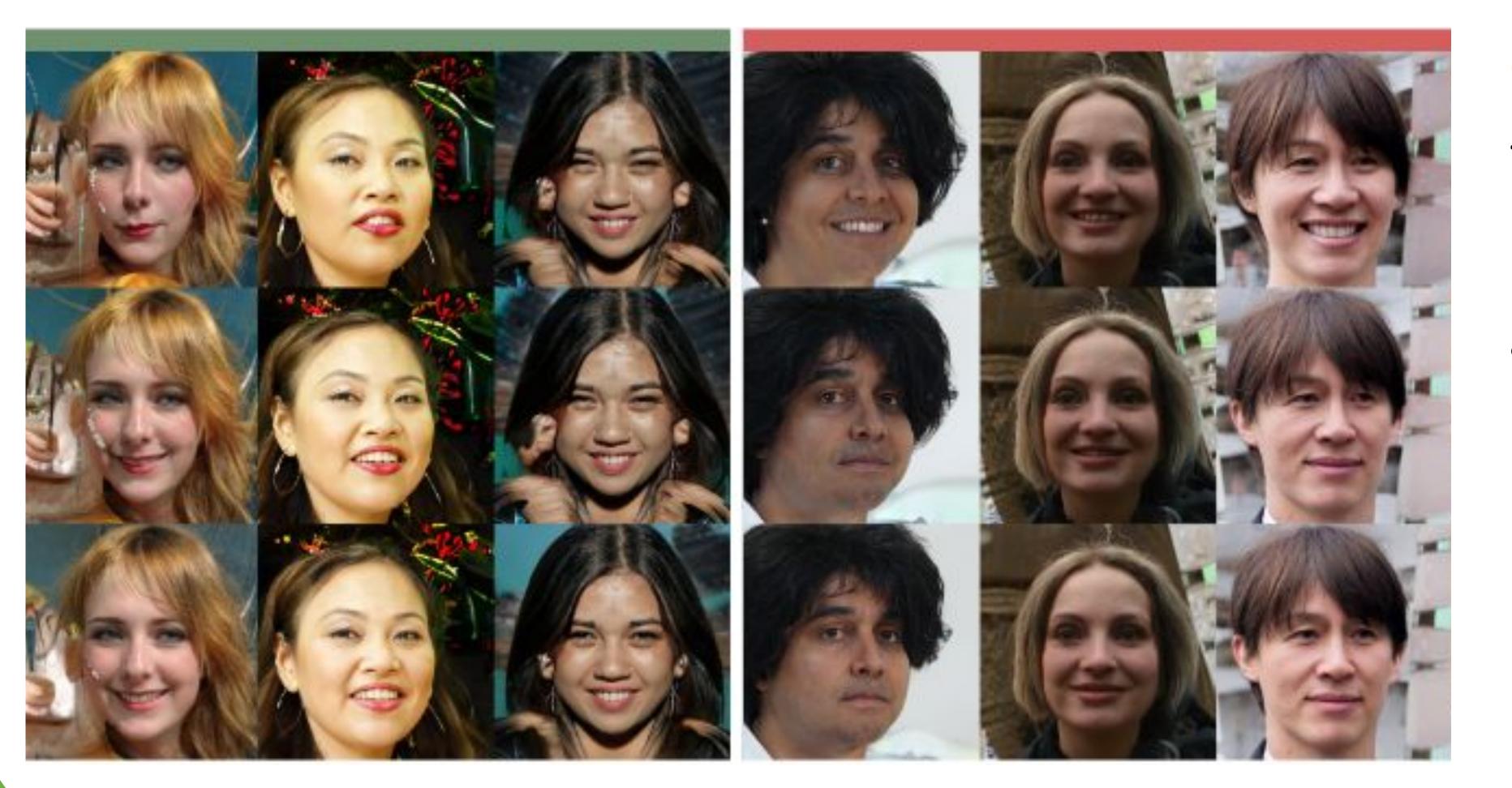


t: Our methods listic counterfactual samples that 72/04/999 modify only the most important class-related Opt-based 7 2 9 4 4 4 5 9 9 features, preserving digit style. This selective Dl 2341 Gasses (e.g., 9 vs. 7 or 5

nt: We show the advantages of reliable CEs in model training on the Galaxy10 DECaLS. When standard CEs are blended into the training data, s as the proportion of CEs model performance increases. In contrast, using re s in both accuracy and AUC with increasing fraction. This is because the standard CEs are on the learned decision boundary of the weak classifier, while reliable CEs align closely with the true decision boundary. They can serve as an augmentation strategy for imbalanced classes, addressing fairness.







: We show that the proposed method is highly scalable (by using an SOTA generative model like Style GAN) and applicable to non-differentiable classifiers, such as human annotators (proxied by a pretrained CLIP