

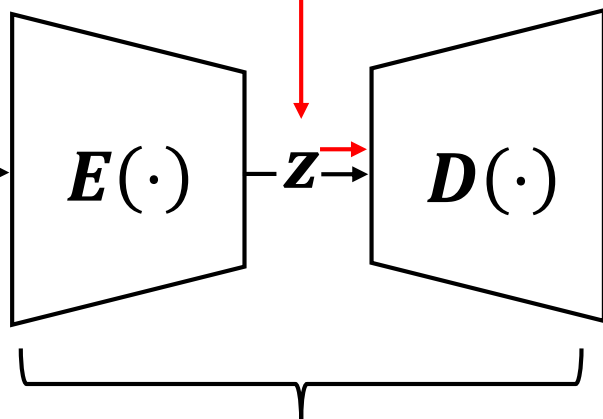
Original image

 $x$ 

Adversarial attack:

$$z \rightarrow z + \delta$$

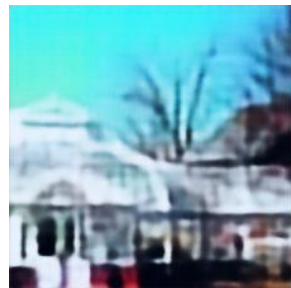
$$\delta = \epsilon \text{sign}(\nabla_{\delta} L_{\text{NCE}})$$

 $G(\cdot)$ 

Identity-disentangled



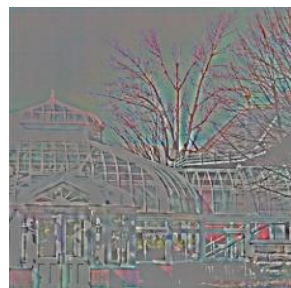
⋮

 $G(x)$ 

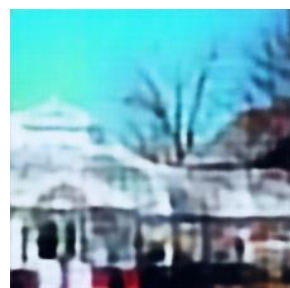
Identity-relevant



⋮

 $x - G(x)$ Adversarially attacked  $G(x)$ 

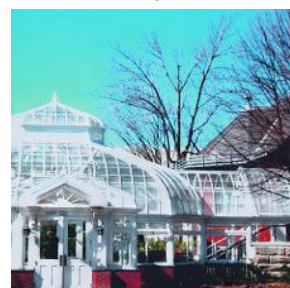
⋮

 $G'(x)$ 

Augmented data



⋮

 $x - G(x) + G'(x)$