

# Principled Knowledge Extrapolation with GANs

Ruili Feng<sup>1</sup>, Jie Xiao<sup>1</sup>, Kecheng Zheng<sup>1</sup>, Deli Zhao<sup>2</sup>, Jingren Zhou<sup>3</sup>, Qibin Sun<sup>1</sup>, Zheng-Jun Zha<sup>1</sup>

<sup>1</sup>University of Science and Technology of China, Hefei, China

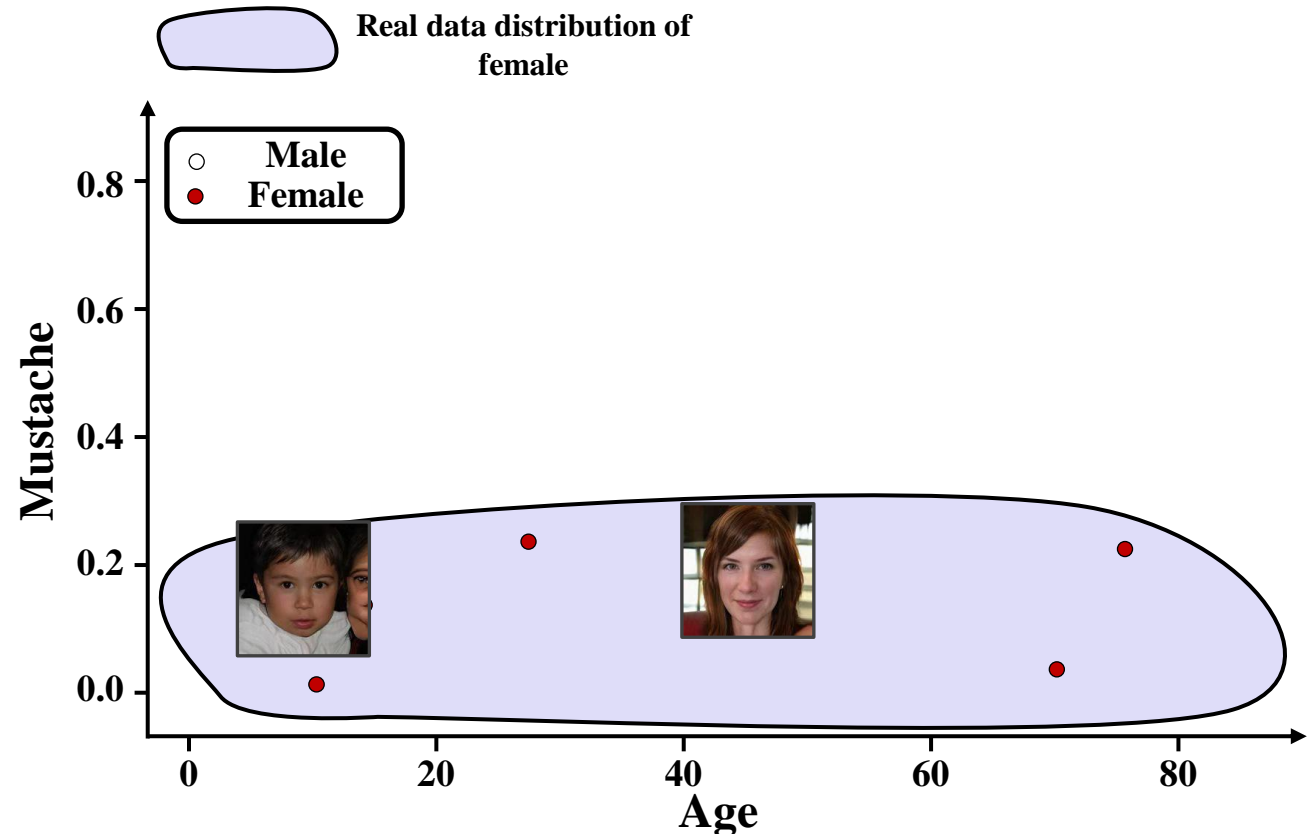
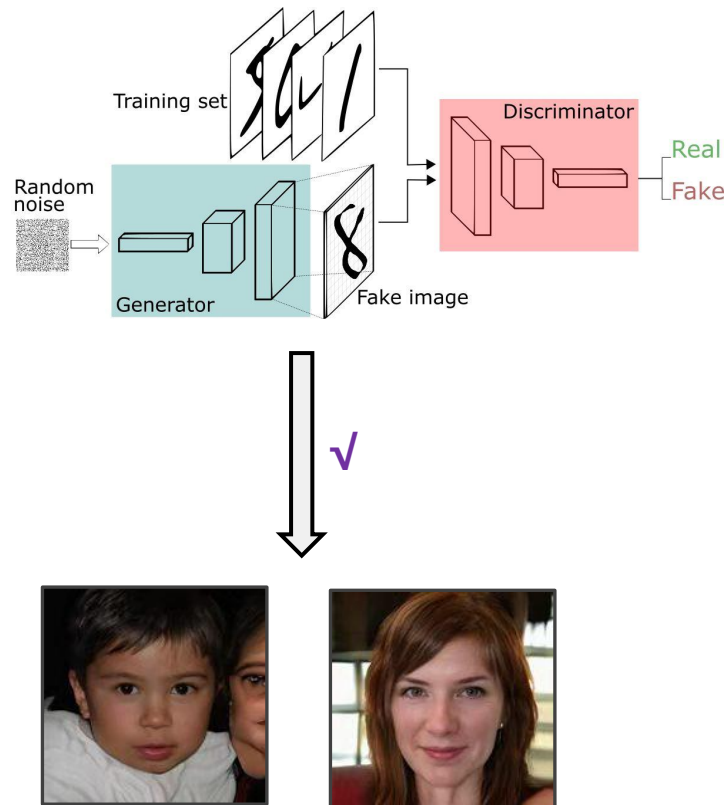
<sup>2</sup>Ant Research, Hangzhou, China

<sup>3</sup>Alibaba Group, Hangzhou, China

ICML 2022, held in Baltimore, Maryland USA

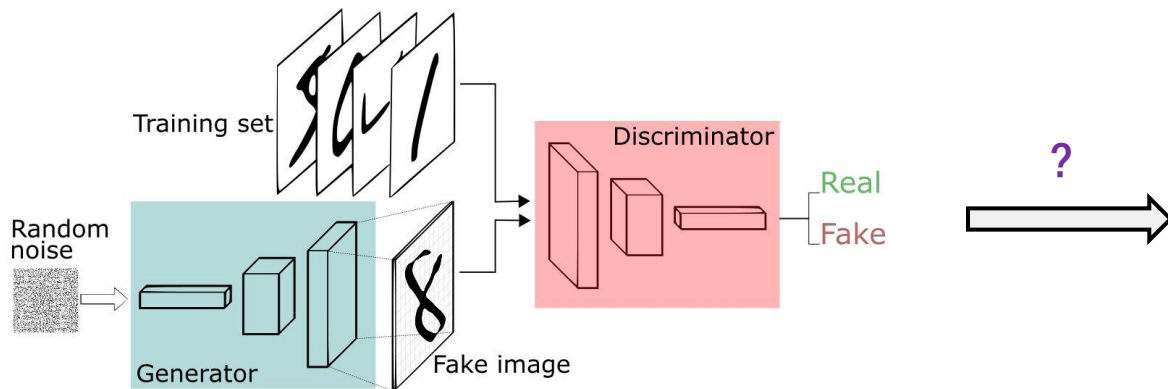
# What is Counterfactual Synthesis

- Generative Networks can well generate high fidelity examples whose distribution is from that of the given training data.



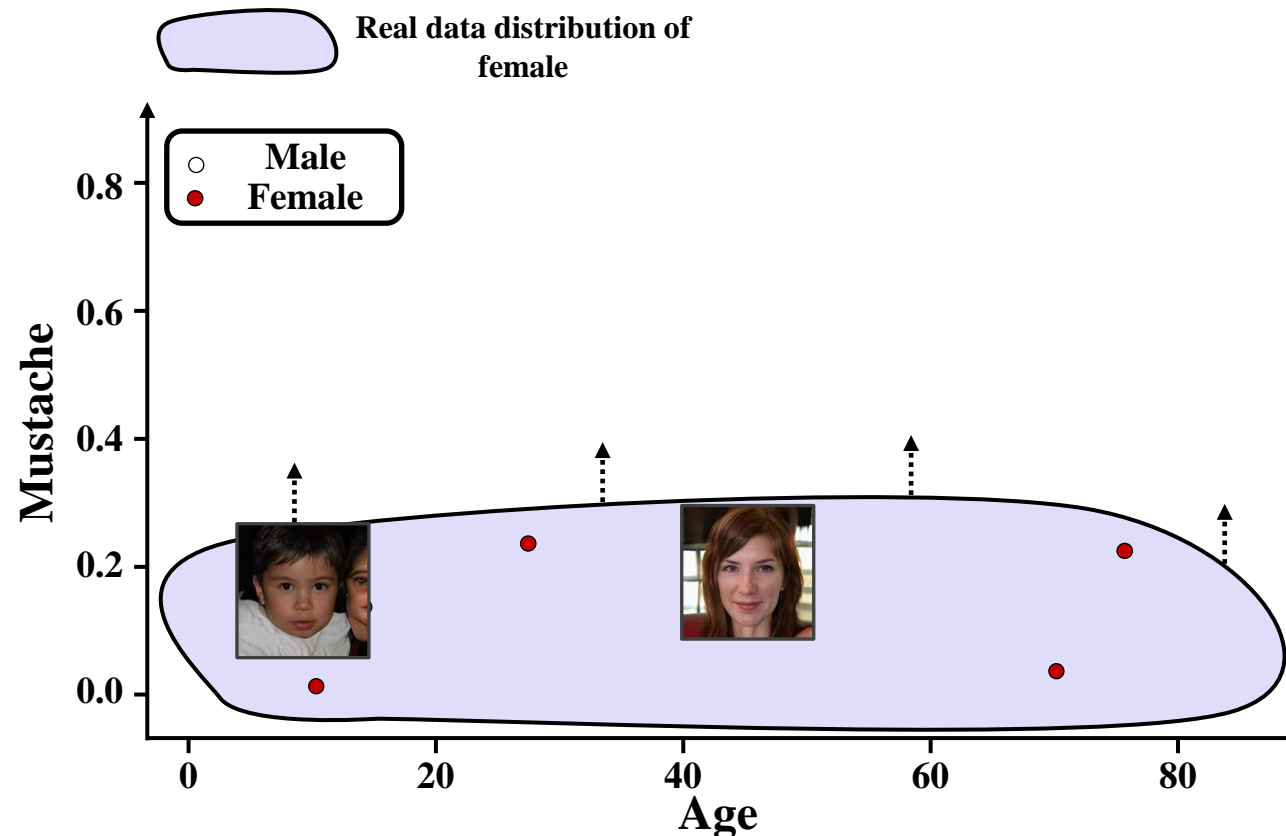
# What is Counterfactual Synthesis

- Current methods fail to generate high fidelity counterfactual results;
- Real data distribution excludes the cases of children or female in mustache.



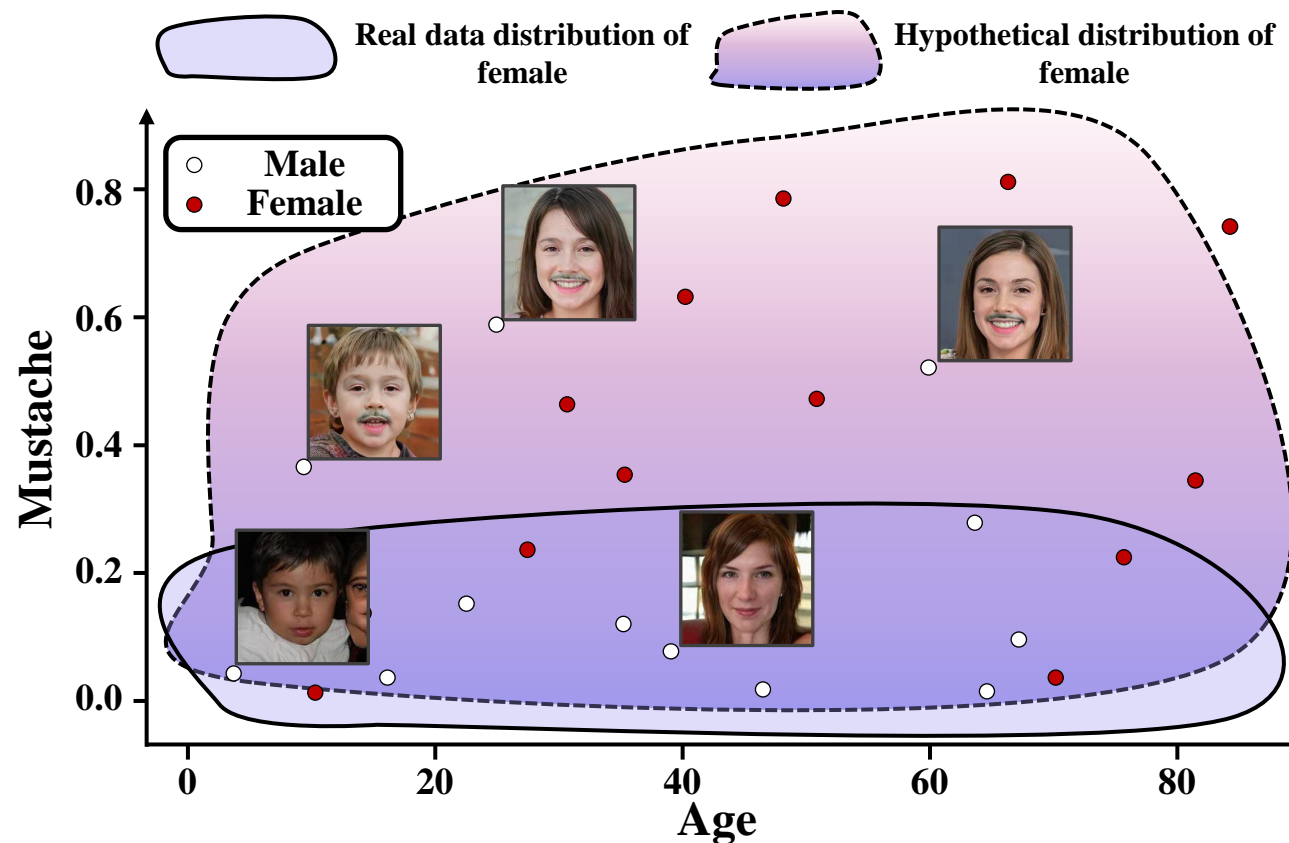
# What is Counterfactual Synthesis

- Counterfactual Synthesis desirably extrapolates to those counterfactual cases, but keeps the other aspects unchanged.



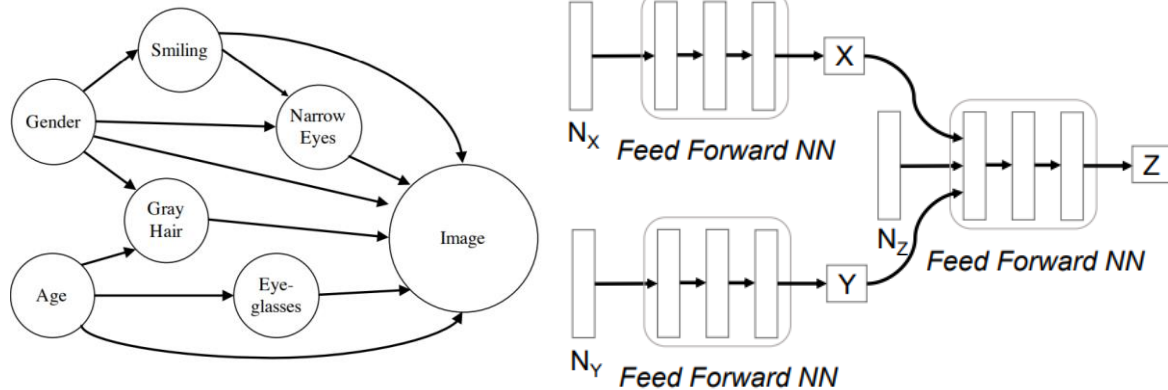
# What is Counterfactual Synthesis

- Counterfactual Synthesis desirably extrapolates to those counterfactual cases, but keeps the other aspects unchanged.

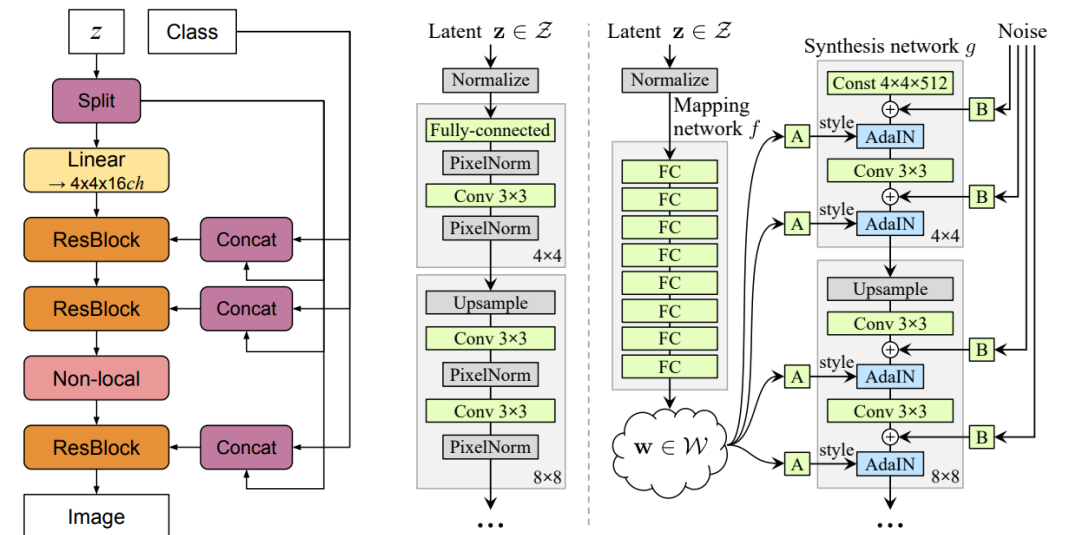


# Structural Causal Model for Counterfactual Synthesis

- Prior SCM to identify causalities
  - many potential factors, obscure entanglements, strong ignorability



- SCM cannot easily leverage state-of-the-art generative models (i.e., BigGAN, StyleGAN)



# Principled GAN Knowledge Extrapolation

- Given:
  - data domain:  $\mathcal{X}$
  - data distribution:  $\mathbb{P}_{\mathcal{X}}$
- Assume:
  - pretrained generator:  $G_{\theta}x: \mathcal{Z} \rightarrow \mathcal{X}$
  - posterior distribution for a knowledge  $l$ :  $\mathbb{P}_l(x) = \mathbb{P}(l|x)$
- Task: infer  $\mathbb{P}_H$ , with only difference from  $\mathbb{P}_{\mathcal{X}}$  in knowledge  $l$

# Principled GAN Knowledge Extrapolation

- Solution: by adversarial training:

$$\min_{G_\theta} \max_{D_\phi} \mathbb{E}_{x \sim \mathbb{P}_H} [\log D_\phi(x)] + \mathbb{E}_{x \sim \mathbb{P}_{G_\theta}} [\log(1 - D_\phi(x))] \quad (1)$$

- Goodfellow et al. have shown that Eq. (1) converges to  $\mathbb{P}_{G_\theta} = \mathbb{P}_H$ !
- Problem: we do not have samples from  $\mathbb{P}_H$ , how to train discriminator?
- If we know the optimal discriminator, then we only need to solve

$$\min_{G_\theta} \mathbb{E}_{x \sim \mathbb{P}_{G_\theta}} [\log(1 - D_{\phi^*}(x))] \quad (2)$$

which do not need samples from  $\mathbb{P}_H$ !



# Indiscernibility Space

- **Indistinguishable assumption**

$\mathbb{P}_\mathcal{X}$  is indistinguishable from  $\mathbb{P}_H$  except for the altered knowledge  $l$

- **Definition**

*the indiscernibility space  $\mathcal{I}^l$  for knowledge  $l$ :*

$\mathcal{I}^l = \{\theta: \mathbb{P}_{G_\theta} \text{ is indistinguishable from } \mathbb{P}_H \text{ except for knowledge } l\}$

- By Indiscernibility Space, solving (1) is equivalent to solving

$$\min_{G_\theta \in \mathcal{I}^l} \max_{D_\phi} \mathbb{E}_{x \sim \mathbb{P}_H} [\log D_\phi(x)] + \mathbb{E}_{x \sim \mathbb{P}_{G_\theta}} [\log (1 - D_\phi(x))]$$

# Optimal Discriminator

**Theorem.** *If  $G_\theta \in \mathcal{I}^l$ , then the optimal discriminator of problem  $\max_{D_\phi} V(D_\phi, G_\theta)$  is  $D_\phi^*(x) = \mathbb{P}_l(x)$  for some probability distribution  $\mathbb{P}_l$  of knowledge  $l$ .*

Problem (2) is equivalent to

$$\min_{\theta \in \mathcal{I}^l} \mathbb{E}_{x \sim \mathbb{P}_{G_\theta}} [\log \mathbb{P}_{\bar{l}}(x)] = -H(\mathbb{P}_{G_\theta}, \mathbb{P}_{\bar{l}}), \quad (3)$$

where  $\mathbb{P}_{\bar{l}} = 1 - \mathbb{P}_l$ ,  $H$  is the cross entropy function.

# Solution to Generator

**Algorithm:** repeatedly update  $\theta$  as

$$\theta^{k+1} = \epsilon \text{Tr} \left( \nabla_{\theta} H \left( \mathbb{P}_{G_{\theta^k}}, \mathbb{P}_{\bar{l}} \right), \lambda \right),$$
$$\text{where } \text{Tr}(v, \lambda) = \begin{cases} v_i = 0, & \text{if } |\nabla_{\theta} H_i| \leq \lambda, \\ v_i = 1, & \text{if } \nabla_{\theta} H_i < 0, 0 \leq \lambda < |\nabla_{\theta} H_i|, \\ v_i = -1, & \text{if } \nabla_{\theta} H_i > 0, 0 \leq \lambda < |\nabla_{\theta} H_i|, \end{cases}$$

and  $\theta^0 = \theta$ , then  $G_{\theta^K}$  generates a distribution that approximates the counterfactual hypothesis  $\mathbb{P}_H$ !

# Principal Knowledge Descent

- **Theorem.** *Let  $\Delta$  be the descent value of the objective (3) by implementing Algorithm 1, and  $\delta$  be the change of the other knowledge, i.e.,*

$$\Delta = H(\mathbb{P}_{G_{\theta^K}}, \mathbb{P}_{\bar{l}}) - H(\mathbb{P}_{G_{\theta}}, \mathbb{P}_{\bar{l}}),$$
$$\delta = \mathbb{E}_{x \sim \mathbb{P}_X} \left[ \left| f_{G_{\theta^K}}^r(x) - f_X^r(x) \right| \right],$$

*where  $K$  is the iteration turns. Assume that  $\mathbb{P}_{G_{\theta^X}} = \mathbb{P}_X$ ,  $\varepsilon$  is small enough,*

*and  $L = \sup_{\|\theta - \theta^X\|_{\infty} < K\varepsilon} \left\| \nabla_{\theta} f_{G_{\theta}}^r \right\|_{\infty}$ . There is  $\lambda_{max} < 0$  such that  $\forall \lambda \in$*

*$(0, \lambda_{max})$ , we have  $\Delta > 0$  and*

$$\frac{\Delta}{\delta} \geq \frac{\lambda}{L} + o(1).$$

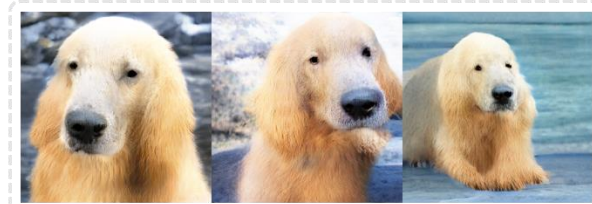
# Findings and Results

- Counterfactual Synthesis

Origin Domain:  
Irish Setter



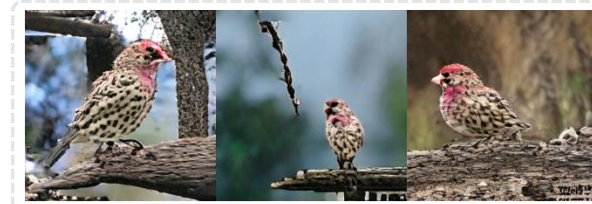
**CF: Polar Bear  
Nose & Color**



Origin Domain :  
Goldfinch



**CF: Cheetah  
Spot**



Origin Domain:  
Husky



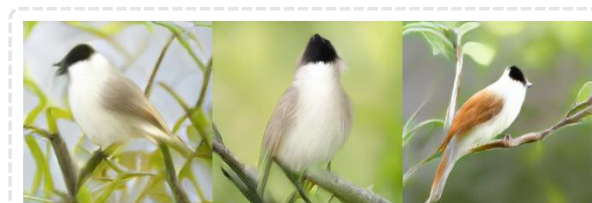
**CF: Ursus  
Arctos Fur**



Origin Domain:  
Junco



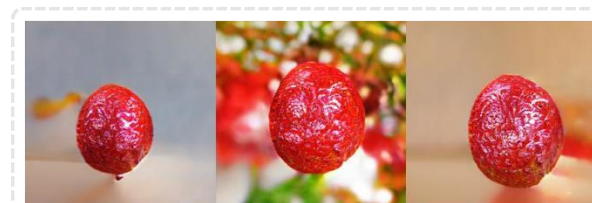
**CF: Guinea Pig  
Fur & Color**



Origin Domain:  
Orange



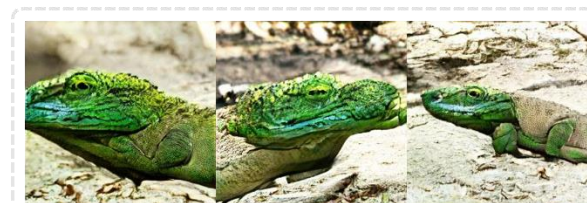
**CF: Strawberry  
Surface**



Origin Domain:  
Komodo Dragon



**CF: Green  
Lizard Color**



- Data domain: ImageNet
- Pretrained generator: BigGAN256
- Posterior distribution: ResNet50



# Findings and Results

- Counterfactual Synthesis



Domain: Female    **Counterfact: Mustache**



Domain: Child    **Counterfact: Gray Hair**



Domain: Child    **Counterfact: Mustache**



Domain: Male    **Counterfact: Lipstick**

- Data domain: FFHQ
- Pretrained generator: StyleGAN2
- Posterior distribution: ResNet50 trained on CelebA-HQ