

FineGAN: [Unsupervised fine-grained object category discovery] ^{fine case}
Purpose: model fine grained object categories by hierarchically disentangling the background, shape & appearance.
(No manual annotations)

- Prior work on unsupervised category discovery is not fine-grained & supervised methods are used for fine grained discovery.
- unsupervised learning enables the discovery of latent structure which may not be done by annotators. [Exploit hierarchy]
- Disentangles the background, object shape & appearance.
- * Use information theory to associate each feature to a latent code & then use ~~that~~ condition the rel b/w them in a specific way \Rightarrow induces hierarchy
 - BG + Parent image + Child image where the parent & child capture 1 factor of variation each. [capturing factors of variation - info theory]

- By imposing constraints on the rel b/w the parent & child codes, the parent learns the shape & the child learns \Rightarrow texture
- FineGAN automatically generates masks to help condition the latent codes to only focus on object factors & stitch them together. [Benchmark: CUB, Stanford-dogs, Stanford-cars]
 \hookrightarrow Beats SOTA unsupervised techniques [JULE, DEPICT]

- InfoGAN maximises mutual information whereas FineGAN learns a hierarchical disentanglement

$$\text{obj} \cdot \text{fx}^n = \mathcal{L} = \lambda \mathcal{L}_b + \beta \mathcal{L}_p + \gamma \mathcal{L}_c$$

latent codes $z \sim \mathcal{N}(0, 1)$ $b \sim \text{cat}(k=N_b, p=1/N_b)$ $p \sim \text{cat}(k=N_p, p=1/N_p)$ \uparrow categorical codes
 $c \sim \text{cat}(k=N_c, p=1/N_c)$

- Imposed hierarchy constraints: (i) $N_p < N_c$
- These help 'p' to capture shape & (ii) For each parent code, there are a few child codes tied to it.
 - 'c' to capture appearance.

- (iii) Sets the background code to be the same as the child code

1. Background Stage

- $\hookrightarrow G_b, D_b, D_{aux}$ \hookrightarrow To generate bg, a detector is required that can detect instances of the super-category
 \hookrightarrow The detector is used to ~~detect~~^{locate} non-object patches in the real image x_i . Train G_b & D_b using:

$$L_b = L_{bg-adv} + L_{bg-aux} \rightarrow \text{also computed at the patch level}$$

\hookrightarrow predicts $N \times N$ grid with real/fake scores

$$L_{bg-adv} = \min_{G_b} \max_{D_b} E_x [\log(D_b(x))] + E_{z,b} [\log(1 - D_b(G_b(z, b)))]$$

$$L_{bg-aux} = \min_{G_b} E_{z,b} [\log(1 - D_{aux}(G_b(z, b)))] \Rightarrow \text{assigns bg \& fg probability to each patch}$$

2. Parent Stage: Generate fg entity & stitch it to the bg.

$G_{f,m}$ & $G_{f,m}$ transform F_p into P_f & P_m (foreground, mask)

$$P = P_{f,m} + B_m \rightarrow \text{stitching}$$

$$P_{f,m} = P_m \odot P_f \text{ and } B_m = (1 - P_m) \odot B$$

Masked foreground Inverse masked background

- Maximise mutual info $I(p, P_{f,m})$ with D_p approximating

$$P(p|P_{f,m})$$

$$L_p = L_{p-info} = \max_{D_p, G_{p,f}, G_{p,m}} E_{z,p} [\log(D_p(p|P_{f,m}))]$$

\hookrightarrow This ensures the decision is based on fg only

- $P_{f,m}$ must not capture any fine grained info and is also common to multiple fine grained categories.

3. Child Stage:

Child code 'c' encodes color/texture and is conditioned on 'c' & F_p are fed to G_c . The generated feature rep F_c is used by generators $G_{c,m}$ & $G_{c,f}$ to get C_m & C_f respectively.

$$C = C_{f,m} + P_{c,m} \rightarrow \text{stitching}$$

$$\hookrightarrow C_{f,m} = C_m \odot C_f ; P_{c,m} = (1 - C_m)P$$

$$L_c = L_{adv} + L_{info}$$

$$L_{adv} = E_x \log(D_{adv}(x)) + E_{z,p,c} [\log(1 - D_{adv}(c))]$$

$$L_{info} = \max_{D_c, G_{c,f}, G_{c,m}} E_{z,p,c} [\log(D_c(c|C_{f,m}))]$$

D_{adv} : Diff b/w real x & gen c

D_c : Approximate $p(c|C_{f,m})$