Contrastive Learning for Unpaired Image-to-Image Translation



Taesung Park Alexei A. Efros Richard Zhang Jun-Yan Zhu



UC Berkeley

Adobe Research

ECCV 2020

What is Unpaired Image-to-Image Translation?







Training Set

Test-time behavior





CycleGAN (Zhu et al., ICCV'17) DiscoGAN (Yi et al., ICCV'17) DualGAN (Yi et al., ICCV'17)

Also used in MUNIT (Huang et al., ECCV'18) DRIT (Lee et al., ECCV'18)









invariant



sensitive

What makes for a good output?

Input (horse)





Output (zebra)
?

Retaining input content

Input (horse)



Output (zebra)

Discriminator

Retaining input content



Corresponding patches should have high similarity



- InfoNCE loss (Gutmann et al., AISTATS18, van den Oord et al., 2018) used in MoCo and SimCLR
- To produce positive pairs:
 - Handcrafted data augmentation (MoCo, SimCLR, etc.)
 - Input and synthesized image (ours)

Patchwise contrastive loss



Patchwise contrastive loss



Patchwise contrastive loss



+ No fixed similarity metric (e.g., L1 or perceptual loss)
+ One-sided (no inverse mapping needed)

Internal vs External Patches



Internal Patches

Internal vs External Patches



External patches make things worse

Power of Internal patches

Texture Synthesis by Non-parametric Sampling (Efros & Leung, ICCV'99, Efros & Freeman, SIGGRAPH'01)



IN WORK IN TERRORIM SWORE WE NOT AN ADDRESS t adapters a super Tring cooms," in Heft he first ad it i en det norden avtana sibed it Last at hent bediga Al-1 established in Al Heftarro," as do Lewisdolf I tian A2 Tha," as Lewing question last attractions. He to then All Loss full reserved a Lossy, at "then do Dynam & the erfinitial). Reconveting merrin," as libuas De fair f De and itical communitarihed relast fall. He fall, Heffit in ambroard washis he left a ringing questica Lewis. icara coronna," autour stage of Monica Lewinow are a Than Firing pictors standbrat nowea or left a roowne bourstof Mik Selft a Lent fast agine lisuser tican Mel of it was " Televasid, a stag tad reisenered it a stag give autical data one years of Moung full. He yibof Mouse or your climb Tripp?" That hedue Al Lest Deer you ats Trees? Initial conseilur Alithe free of sing dos Littical room proving of the stoaruss of all 7 get room L oon Lewron Denta show I He fan goent og ing of, at hoos

"Zero-Shot" Super-resolution using Deep Internal Learning (Shocher, Cohen & Irani CVPR'18)





EDSR [

ZSSR (ours)

Internal vs External Patches





DTN (Taigman et al., ICLR'17), CycleGAN (Zhu et al., ICCV'17)



DTN (Taigman et al., ICLR'17), CycleGAN (Zhu et al., ICCV'17)



Lighter Footprint

Training time (sec/iter, lower is better)



Lighter Footprint

Training time (sec/iter, lower is better)



Lighter Footprint

Training time (sec/iter, lower is better)





Dealing with Dataset Bias

Source training set



horse 17.9%

Target training set



zebra 36.8%

Dealing with Dataset Bias



horse 17.9%

detected pixels:

zebra 30.8%

zebra 25.9%

zebra 19.1%

zebra 36.8%



Cat



Yosemite Summer -





GTA

FID evaluating the realism of output images (lower is better)



Segmentation Score evaluating correspondences



Single Image Translation

Claude Monet's painting



Internal contrastive loss is well-suited for single image translation. Also see InGAN (Shocher et al., ICCV'19), SinGAN (Shaham et al., ICCV'19)

Single Image Translation

Reference photo



Claude Monet's painting



Internal contrastive loss is well-suited for single image translation. Also see InGAN (Shocher et al., ICCV'19), SinGAN (Shaham et al., ICCV'19)

Single Image Translation

Reference photo



Internal contrastive loss is well-suited for single image translation. Also see InGAN (Shocher et al., ICCV'19), SinGAN (Shaham et al., ICCV'19)



Painting

Reference Photo



Painting


Gatys et al. CVPR'16





STROTSS (Kolkin et al., CVPR'19)



WCT² (Yoo et al., ICCV'19)











Painting





Gatys et al. CVPR'16



Painting

STROTSS (Kolkin et al., CVPR'19)



Painting

WCT² (Yoo et al., ICCV'19)













Our translation result







Painting



Painting



Painting

Questions or Comments?













Swapping Autoencoder for Deep Image Manipulation Tassung Park - Jun Yan Zhou, Giver Wang, Jingwan Lin, Ek Shekhiman, Alexa Elros ¹, Renard Zhoor

UC Behickey, Adobe Research

inter-image intra-image

Disentanglement?

content



style



MUNIT (Huang, Liu, Belongie, Kautz, ECCV'18)

Structure for each row



Style for each column



Extracting style and structure from an image



Extracting style and structure from an image



Extracting style and structure from an image











Swap







What is Texture?

"An image that can be represented by first and second-order statistics"



Conjecture by Bela Julesz, 1962

Two textures that differ by first-order statistics

What is Texture?

"An image that can be represented by first and second-order statistics"



Conjecture by Bela Julesz, 1962

Two textures that differ by second-order statistics

What is Texture?

"An image that can be represented by first and second-order statistics"



Conjecture by Bela Julesz, 1962

Two textures that differ by third-order statistics
















structure image





Embedding a Real Input Image

Fast

feed-forward pass of the encoder. Magnitudes faster than baselines.

Accurate

Prior distribution not enforced on the latent space.

Spatial resolution is retained.



StyleGAN2 (Karras et al., CVPR'20), Im2StyleGAN (Abdal et al., ICCV'19)

Embedding a Real Input Image

Fast

feed-forward pass of the encoder. Magnitudes faster than baselines.

Accurate

Prior distribution not enforced on the latent space.

Spatial resolution is retained.



Photorealistic and Disentangled Swapping Quality

Structure

Texture



StyleGAN2Im2StyleGANSTROTSSWCT2OursIm2StyleGANIm2StyleGANIm2StyleGANIm2StyleGANIm2StyleGAN



Photorealistic and Disentangled Swapping Quality





StyleGAN2 (Karras et al., CVPR'20), Im2StyleGAN (Abdal et al., ICCV'19) STROTSS (Kolkin et al., CVPR'19), WCT² (Yoo et al., ICCV'19)

Realism of generated images

Method	Runtime	Human Perceptual Study (AMT Fooling Rate) (†)			
	(sec)(1)	Church	FFHQ	Waterfall	Average
Swap Autoencoder (Ours)	0.113	31.3±2.4	29.8 ± 2.5	41.8 ± 2.2	34.4 ± 1.4
Im2StyleGAN [990	8.5±2.1	11.6 ± 2.3	12.8±2.4	11.0±1.3
StyleGAN2 [38]	192	24.3 ± 2.2	22.8 ± 2.5	35.3 ± 2.4	27.5 ± 1.4
STROTSS [41]	166	13.7±2.2	13.6±2.5	23.0±2.1	16.7±1.3
WCT ² [69]	1.35	27.9±2.3	$26.6{\pm}2.4$	$35.8{\pm}2.4$	30.1 ± 1.4
StyleGAN2 Im:	2StyleGAN	STROTSS	WCT ²	Ours	
	Structure		Iro		

Structure





Which do you think is more similar in style?

Which do you think is more similar in structure/content?

 WCT^2



















Texture





Smooth Latent Space



 $Z_{add_snow} = z_{snow} - z_{summer}$

Smooth Latent Space



more snow

input image

less snow

PCA on the Latent Space



discovered edit vectors

Interactive UI



Structural Editing

Extract the structure code at this position



Structural Editing

Overwrite the structure code here



User-Guided Portrait Painting to Photo



input

output

User-Guided Animal Face Transformation



input

same pose, different styles

PCA on the Structure Code



direction

5 o'clock shadow

Editing Landscape Images



Editing Landscape Images



UI with input image



brush stroke visualization



1. remove road



2. draw mountain



PCA on the structure code, with user-drawn mask







PCA with the style (texture) code



Summary



GAN that can embed images

Summary



structure / style disentanglement

Summary



interactive user editing



https://taesung.me/ContrastiveUnpairedTranslation https://taesung.me/SwappingAutoencoder