Spatial Content Alignment For Pose Transfer



Wing-Yin Yu



Lai-Man Po



Yuzhi Zhao



Jingjing Xiong



Kin-Wai Lau

Friday Seminar 2021-10-15



Department of Electrical Engineering

香港城市大學 City University of Hong Kong

Content

- Background knowledge
- The task I have done
- Method
- Experiment result



Deep Learning (DL)

- Deep Learning is a subfield of machine learning
- Automatically extract features by a neural network





香港城市大學 City University of Hong Kong

Computer Vision (CV)

General computer vision tasks



Image Classification

Classify an image based on the dominant object inside it.

datasets: MNIST, CIFAR, ImageNet



Object Localization

Predict the image region that contains the dominant object. Then image classification can be used to recognize object in the region datasets: ImageNet



Object Recognition

Localize and classify all objects appearing in the image. This task typically includes: proposing regions then classify the object inside them. datasets: PASCAL, COCO



Semantic Segmentation

Label each pixel of an image by the object class that it belongs to, such as human, sheep, and grass in the example. datasets: PASCAL, COCO



Instance Segmentation

Label each pixel of an image by the object class and object instance that it belongs to.

datasets: PASCAL, COCO



Keypoint Detection

Detect locations of a set of predefined keypoints of an object, such as keypoints in a human body, or a human face. datasets: COCO

Credit: https://medium.com/@nikasa1889/the-modern-history-of-object-recognition-infographic-aea18517c318



What if...

Image \rightarrow Label



Label \rightarrow Image?





- A framework of Deep Learning model
 - Capture **distribution** from training data
 - Generate new data from same distribution
- Major components:
 - **Generator**: generate real-looking examples
 - Discriminator: classify examples as real or fake



Generator

Discriminator

$$\min_{G} \max_{D} V(D, G)$$

$$V(D,G) = \mathbb{E}_{x \sim p_{data}(x)}[\log D(x)] + \mathbb{E}_{z \sim p_z(z)}[\log(1 - D(G(z)))]$$

Goodfellow, Ian, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. "Generative adversarial networks." *Communications of the ACM* 63, no. 11 (2020): 139-144.



Problem overview

Task 1: Detect the pose keypoints from the image



Task 2: Generate image from the keypoints



• Our task: Generate image with other poses





Image translation





Pose transfer

Input

- Source person image
- Source pose
- Target pose

Output

- A new person image
 - > Target pose
 - Style and Texture of source person image





Motivation

Spatial misalignment problem

- Sparse correspondence between pose landmarks
 - > Warp-based method: blurry effect
 - > Attention-based method: lack of excitation
- Insufficient content guidance
 - Person identity
 - Garment reconstruction
 - > Incomplete texture

Content-driven guidance

- Leverage edge map as an extra constraint with pose heatmap
- Guide the network to produce more photo-realistic person images through texture enhancement



Source



Target





Methodology

Two-phase approach

- Phase 1: Prior Content Transfer Network (PCT-Net)
- Phase 2: Image Synthesis Network (IS-Net)





Phase 1: Prior Content Transfer Network (PCT-Net)

- Objective: generate the content in advance
 - Leverage pixel-level edge map to highlight high-frequency signals
 - Dominate a wide range of spectrum for content information

$$E_g = \mathcal{F}(E_s, P_s, P_t)$$



Methodology

Phase 2: Image Synthesis Network (IS-Net)

- Objective: render a realistic-looking person image
 - Encoder-decoder architecture
 - Appearance style of I_s , prior content E_g , spatial content P_t



$$I_g = \mathcal{G}(I_s, E_g, P_t)$$

Department of

香港城市大學 City University of Hong Kong

Electrical Engineering

Content-Style Spatially-Adaptive Normalization (CS-SPADE)

- Inspired by SPADE
- Accept aligned content and non-aligned style features as Source
- Inject content and style information with learnable parameters

$$\mathcal{W}(\mathbf{S}, f) = \delta(\gamma_{c,h,w}^{i}(\mathbf{S}) \cdot \frac{f_{n,c,h,w}^{i} - u_{nc}^{i}}{\sigma_{nc}^{i}} + \beta_{c,h,w}^{i}(S))$$

 $\delta = {\rm LeakyReLU}$

Mean of
$$f_{n,c,h,w}^{i} = u_{nc}^{i} = \frac{1}{H^{i}W^{i}} \sum_{h,w} f_{n,c,h,w}^{i}$$

SD of $f_{n,c,h,w}^{i} = (\sigma_{nc}^{i})^{2} = \frac{1}{H^{i}W^{i}} \sum_{h,w} f_{n,c,h,w}^{i}^{2} - u_{nc}^{i}^{2}$

Content-Style Decoder Block (Content-Style DeBlk)

$$I_g^l = \mathcal{U}\left(E_g, \mathcal{W}(I_g^{l-1}, I_s^L)\right) + \mathcal{V}\left(P_t, \mathcal{W}(I_g^{l-1}, I_s^L)\right)$$



$$\mathcal{L}_{full} = \lambda_{adv} \mathcal{L}_{adv} + \lambda_1 \mathcal{L}_1 + \lambda_{per} \mathcal{L}_{per} + \lambda_{cx} \mathcal{L}_{cx}$$

Adversarial loss:

$$\mathcal{L}_{adv} = \mathbb{E}_{I_s, I_t, I_g} \left[\log \left(D_s(I_s, I_t) \right) + \log \left(1 - D_s(I_s, I_g) \right) \right] + \mathbb{E}_{I_t, P_t, I_g} \left[\log \left(D_c(P_t, I_t) \right) + \log \left(1 - D_c(P_t, I_g) \right) \right]$$

Appearance loss:

$$\mathcal{L}_1 = \frac{1}{CHW} \sum_{c,h,w} \left\| I_g \right\|_{c,h,w} - I_t |_{c,h,w} \right\|_1$$

Perceptual loss:

$$\mathcal{L}_{per} = \frac{1}{C_{\ell}H_{\ell}W_{\ell}} \sum_{c,h,w} \left\| \theta_{\ell} \left(I_{g} \right) \right\|_{c,h,w} - \theta_{\ell}(I_{t}) \left\|_{c,h,w} \right\|_{1}$$

Contextual loss:

$$\mathcal{L}_{cx} = -\frac{1}{C_{\ell}H_{\ell}W_{\ell}}\sum_{c,h,w}\log\left[CX\left(\theta_{\ell}\left(I_{g}\right)|_{c,h,w},\theta_{\ell}(I_{t})|_{c,h,w}\right)\right]$$



Adversarial loss:

Minimize the KL divergence between two different domains

香港城市大學 City University of Hong Kong

Appearance loss:

Minimize the pixel-wise least absolute distance

$$\mathcal{L}_{1} = \frac{1}{CHW} \sum_{c,h,w} \left\| I_{g} \right\|_{c,h,w} - I_{t} |_{c,h,w} \right\|_{1}$$







Perceptual loss:

minimize the distance in feature space,

$$\mathcal{L}_{per} = \frac{1}{C_{\ell}H_{\ell}W_{\ell}} \sum_{c,h,w} \left\| \theta_{\ell} \left(I_{g} \right) \right\|_{c,h,w} - \theta_{\ell}(I_{t}) \left\|_{c,h,w} \right\|_{1}$$

Minimize the L1 distance for the features

Johnson, Justin, Alexandre Alahi, and Li Fei-Fei. "Perceptual losses for real-time style transfer and super-resolution." In *European conference on computer vision*, pp. 694-711. Springer, Cham, 2016.



Department of Electrical Engineering 香港城市大學 City University of Hong Kong



Contextual loss:

minimize the distance in feature space,

$$\mathcal{L}_{cx} = -\frac{1}{C_{\ell}H_{\ell}W_{\ell}} \sum_{c,h,w} \log \left[CX \left(\theta_{\ell} \left(I_{g} \right) |_{c,h,w}, \theta_{\ell}(I_{t}) |_{c,h,w} \right) \right]$$

Generated Image I_a Target Image I_t

Maximize the similarity of the features

Roey Mechrez, Itamar Talmi, and Lihi Zelnik-Manor, "The contextual loss for image transformation with non-aligned data," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 768–783.



Department of Electrical Engineering 香港城市大學 City University of Hong Kong

Comparison with SOTA

Table 1: Quantitative analysis of comparison against the cur-rent state-of-art methods. The best scores are highlighted withbold format.

Methods	IS↑	SSIM↑	FID↓	LPIPS↓
PG ² [8]	3.202	0.773	-	-
Def-GAN [2]	3.362	0.760	18.457	0.233
RATE-Net [3]	3.125	0.774	14.611	0.218
PATN [4]	3.209	0.773	19.816	0.253
APS [5]	3.295	0.775	15.017	0.178
ADGAN [6]	3.364	0.772	13.224	0.176
SCA-GAN(Ours)	3.497	0.775	11.676	0.167

- Our method outperforms the current SOTA methods
- 4% enhancement of IS score compared to ADGAN
- Generate high-quality images with natural texture synthesis
- Beneficial to the edge content guidance for supporting the geometric transformation of fine-gained appearance details



Fig. 2: Qualitative results of comparison against current stateof-art methods. Please zoom in for details.



Ablation Study

 Table 2: Quantitative results of ablation study.
The best scores are highlighted with bold format.

Methods	IS↑	SSIM↑	FID↓	LPIPS↓
with semantic content	3.388	0.762	18.634	0.195
w/o prior transfer	3.365	0.774	12.402	0.166
w/o content branch	3.274	0.756	19.503	0.194
SPADE ResBlk	3.343	0.772	13.888	0.170
Batch-Norm (encoder)	3.340	0.775	12.197	0.166
In-Norm (encoder)	3.301	0.773	12.738	0.167
Full model	3.497	0.775	11.676	0.167

- Edge content is better than semantic content
- Provide sufficient guidance to perform detailed spatial alignment
- Preserve better boundary reconstruction between some ambiguous objects
- Content-Style DeBlk outperforms the original SPADE method
- Better image synthesis quality



Fig. 3: Qualitative results of ablation study. Please zoom in for details.



Slide 20

香港城市大學 City University of Hong Kong

Open source

Yu, Wing-Yin, Lai-Man Po, Yuzhi Zhao, Jingjing Xiong, and Kin-Wai Lau. "Spatial Content Alignment for Pose Transfer." In *2021 IEEE International Conference on Multimedia and Expo (ICME)*, pp. 1-6. IEEE, 2021.



https://github.com/rocketappslab/SCA-GAN



Thank You Q & A

