

Video-to-Video synthesis

Paper summary by Akilesh B

August 21, 2018

1 Introduction

- First work that tries to address the problem of video-to-video translation.
- Directly applying existing image to image translation works (including state-of-the-art pix2pixHD) on videos lead to output videos that are temporally inconsistent and of low perceptual quality.
- Their key contribution is well-designed generator and discriminator architectures along with spatio-temporal adversarial objective.
- Achieve high-resolution, photorealistic, temporally coherent video results on a diverse set of input formats including segmentation masks, sketches, and poses.

2 Network architecture

2.1 Notation

- $s_1^T \equiv \{s_1, s_2, \dots, s_T\}$ denote a sequence of source images (for example image outlines)
- $x_1^T \equiv \{x_1, x_2, \dots, x_T\}$ be the sequence of corresponding ground truth translated images (for example colored images).
- $\tilde{x}_1^T \equiv \{\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_T\}$ be the sequence of generated images.

2.2 Key idea

- To learn a mapping function that can convert s_1^T to \tilde{x}_1^T so that the conditional distribution of \tilde{x}_1^T given s_1^T is identical to the conditional distribution of x_1^T given s_1^T .

$$p(\tilde{x}_1^T | s_1^T) = p(x_1^T | s_1^T) \quad (1)$$

2.3 Sequential generator

Under Markovian assumption, the video frames can be generated sequentially and the generation of frame \tilde{x}_t at time step t depends only on:

- Current source image(s_t)
- Past L source images(s_{t-L}^{t-1})
- Past L generated images($x_{t-L}^{\sim t-1}$)

L is typically chosen to be 2. A feed-forward network F is trained to model the conditional distribution as $x_t^{\sim} = F(x_{t-L}^{\sim t-1}, s_{t-L}^t)$

$$F(x_{t-L}^{\sim t-1}, s_{t-L}^t) = (\mathbf{1} - m_t^{\sim}) \odot w_{t-1}^{\sim}(x_{t-1}^{\sim}) + m_t^{\sim} \odot h_t^{\sim} \quad (2)$$

where \odot is the element-wise product operator and $\mathbf{1}$ is an image of all ones. The first part corresponds to pixels warped from the previous frame, while the second part aims to hallucinate new pixels.

- $w_{t-1}^{\sim} = W(x_{t-L}^{\sim t-1}, s_{t-L}^t)$ is the estimated optical flow from x_{t-1}^{\sim} to x_t^{\sim} , and W is the optical flow prediction function.
- $h_t^{\sim} = H(x_{t-L}^{\sim t-1}, s_{t-L}^t)$ is the hallucinated image, an image generated from scratch.
- $m_t^{\sim} = M(x_{t-L}^{\sim t-1}, s_{t-L}^t)$ is the occlusion mask with continuous values between 0 and 1.

M , W and H are implemented using residual nets.

2.4 Conditional image discriminator D_I

D_I should output 1 for a true pair (x_t, s_t) and 0 for a fake one (x_t^{\sim}, s_t) .

2.5 Conditional video discriminator D_V

The purpose of D_V is to ensure that consecutive output frames resemble the temporal dynamics of a real video given the same optical flow. Suppose, w_{t-K}^{t-2} be $K-1$ optical flow for the K consecutive real images x_{t-K}^{t-1} . D_V should output 1 for a true pair $(x_{t-K}^{t-1}, w_{t-K}^{t-2})$ and 0 for a fake one $(x_{t-K}^{\sim t-1}, w_{t-K}^{t-2})$

2.6 Learning objective function

$$\min_F (\max_{D_I} L_I(F, D_I) + \max_{D_V} L_V(F, D_V)) + \lambda_W L_W(F). \quad (3)$$

where L_I is the GAN loss on images defined by the conditional image discriminator D_I , L_V is the GAN loss on K consecutive frames defined by D_V , and $L_W(F)$ is the flow estimation loss as given below.

$$L_W = \frac{1}{T-1} \sum_{t=1}^{T-1} (\|w_t^{\sim} - w_t\|_1 + \|w_t^{\sim}(x_t) - x_{t+1}\|_1) \quad (4)$$

2.7 Foreground-background prior

The image hallucination function(H) is further decomposed into a foreground model $h_{F,t} = H_F(s_{t-L}^t)$ and a background model $h_{B,t} = H_B(x_{t-L}^{t-1}, s_{t-L}^t)$.

- Optical flow can be estimated accurately in background motion and hence, background image synthesis can be generated accurately via warping.
- Optical flow estimation on a foreground object becomes difficult as it often has a large motion and occupies a small portion of the image.

Hence, the equation in 2 becomes

$$F(x_{t-L}^{t-1}, s_{t-L}^t) = (\mathbf{1} - m_t^\sim) \odot w_{t-1}^\sim(x_{t-1}^\sim) + m_t^\sim \odot ((1 - m_{B,t}) \odot h_{F,t}^\sim + m_{B,t} \odot h_{B,t}^\sim) \quad (5)$$

3 Experiments

- Evaluate proposed approach on several datasets: Cityscapes, Apolloscape, Face video dataset, Dance video dataset etc.
- Compare their models with two baselines trained on same data:
 - (a) *pix2pixHD*: State-of-the-art image to image translation work [1]
 - (b) *COVST*: Based on coherent video style transfer [2] by replacing the stylization network with pix2pixHD.
- They show both subjective and objective metrics for performance evaluation using human preference scores and Frechet Inception Distance (FID).

References

- [1] Wang, T.C., Liu, M.Y., Zhu, J.Y., Tao, A., Kautz, J. and Catanzaro, B., 2018. High-resolution image synthesis and semantic manipulation with conditional gans. In IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (Vol. 1, No. 3, p. 5).
- [2] Chen, D., Liao, J., Yuan, L., Yu, N. and Hua, G., 2017, March. Coherent online video style transfer. In Proc. Intl. Conf. Computer Vision (ICCV).