# Cross-View Exocentric to Egocentric Video Synthesis

Gaowen Liu
gaoliu@cisco.com
Cisco Systems
San Jose, CA, USA

Hao Tang
hao.tang@unitn.it
DISI, University of Trento
Trento, Italy

Hugo Latapie
hlatapie@cisco.com
Cisco Systems
San Jose, CA, USA

Jason Corso
jcorso@stevens.edu
Stevens Institute of Technology
Hoboken, NJ, USA

Yan Yan
yyan34@iit.edu
Illinois Institute of Technology
Chicago, IL, USA

## ABSTRACT

Cross-view video synthesis task seeks to generate video sequences of one view from another dramatically different view. In this paper, we investigate the exocentric (third-person) view to egocentric (first-person) view video generation task. This is challenging because egocentric view sometimes is remarkably different from the exocentric view. Thus, transforming the appearances across the two different views is a non-trivial task. Particularly, we propose a novel Bi-directional Spatial Temporal Attention Fusion Generative Adversarial Network (STA-GAN) to learn both spatial and temporal information to generate egocentric video sequences from the exocentric view. The proposed STA-GAN consists of three parts: temporal branch, spatial branch, and attention fusion. First, the temporal and spatial branches generate a sequence of fake frames and their corresponding features. The fake frames are generated in both downstream and upstream directions for both temporal and spatial branches. Next, the generated four different fake frames and their corresponding features (spatial and temporal branches in two directions) are fed into a novel multi-generation attention fusion module to produce the final video sequence. Meanwhile, we also propose a novel temporal and spatial dual-discriminator for more robust network optimization. Extensive experiments on the Side2Ego and Top2Ego datasets [11] show that the proposed STA-GAN significantly outperforms the existing methods.

## CCS CONCEPTS

• **Computer systems organization** → **Embedded systems**; *Redundancy*; Robotics; • **Networks** → Network reliability.

## KEYWORDS

Cross-view Video Synthesis, Exocentric, Egocentric

**ACM Reference Format:**
Gaowen Liu, Hao Tang, Hugo Latapie, Jason Corso, and Yan Yan. 2021. Cross-View Exocentric to Egocentric Video Synthesis. In *MM '18: ACM*
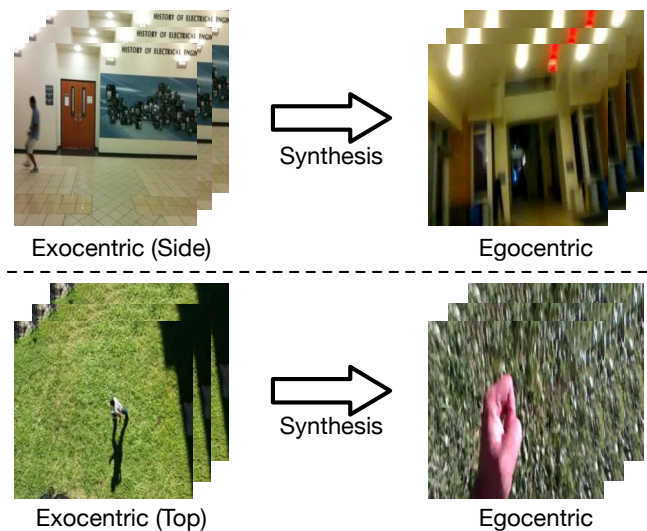
**Figure 1: The goal of exocentric to egocentric cross-view video synthesis is to generate video sequence from exocentric perspective (Side/Top) to egocentric perspective.**

## 1 INTRODUCTION

Wearable cameras, also known as first-person cameras, are widely used in our daily lives since the appearance of low price but high-quality wearable products such as Google Clips, GoPro cameras [1]. Meanwhile, egocentric (first-person) vision has become a critical research topic in the computer vision field [3, 11, 19, 23, 53]. As we know, first-person egocentric views have some unique properties other than third-person exocentric views. Traditional exocentric cameras usually give a wide and global view of the high-level appearance in a video. However, egocentric cameras are able to reveal the focus of attention, behavior, and goal of its wearer. Early egocentric vision studies [19] found that humans are able to seamlessly transfer knowledge between egocentric and exocentric perspectives when performing different activities or interacting with objects. Therefore, understanding the relationship between egocentric and exocentric views is a critical need in computer vision.

However, there is little research to address this important problem in the literature. One likely reason is the difficulty in collecting
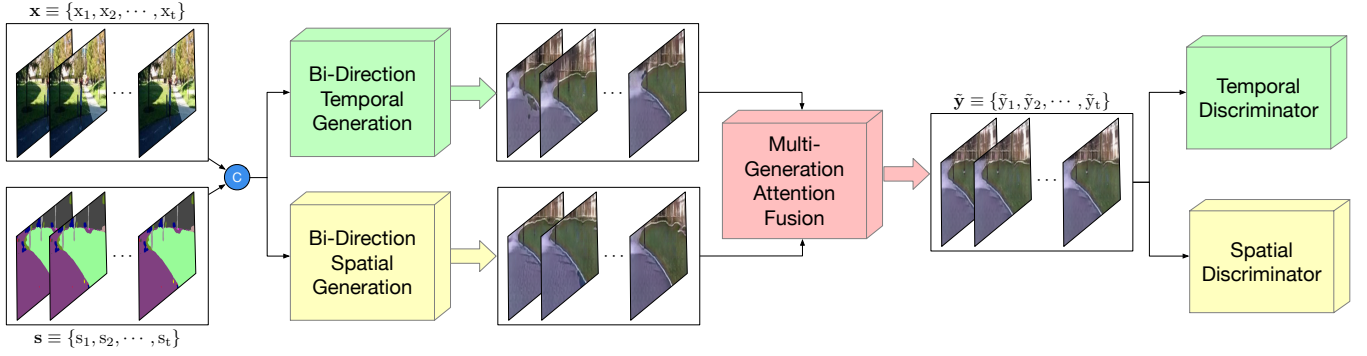
**Figure 2: The framework of the proposed STA-GAN, which consists of four parts, i.e., a temporal generation branch, a spatial generation branch, an attention fusion module and a temporal and spatial dual-discriminator. The bi-directional spatial and temporal generation branches accept exocentric video sequence and conditional semantic maps as inputs and simultaneously synthesizes egocentric video sequence. Then multi-generation attention fusion module fuses the synthesized video sequences that obtained from temporal and spatial generation branches and outputs the final egocentric view video sequence. The proposed dual-discriminator aims to distinguish the generated videos from two spaces, i.e., temporal space and spatial space.**

a large amount of high quality egocentric view data for the wide variety of computer vision problems in different emerging context [31, 53]. To address such a limitation, the technique of generating egocentric videos from exocentric videos, as shown in Figure 1, becomes an alternative solution. Nevertheless, generating egocentric data is extremely challenging because significant differences of visual appearance are expected between egocentric and exocentric videos. Moreover, the sharp changes in the viewpoint makes it an even more difficult task.

Based on these observations, in this paper, we propose a Generative Adversarial Network (GAN) based video generation approach to bridge egocentric and exocentric video analysis. In the past few years, GANs [16] are extremely successful in various image generation problems [5, 10, 17, 18, 20, 32, 38, 46, 55]. Particularly, X-Fork [32], X-Seq [32] and SelectionGAN [45] are proposed to tackle cross-view image generation tasks. However, these approaches do not aim at more challenging cross-view video generation tasks. GANs are also developed for video generation in literature. RecycleGAN [4] works on unsupervised video retargeting for various domains such as face retargeting. Wang et al. [50] propose vid2vid framework which is able to transform a sequence of semantic representations, e.g., semantic label map and sketch map, to a sequence of video frames. StarGAN [8] allows simultaneous training of multiple datasets with different domains within a single network. However, these methods are not able to generate satisfactory results on cross-view video generation tasks due to the dramatically differences between exocentric and egocentric views.

Overcoming these limitations, we propose a novel Bi-directional Spatial Temporal Attention fusion GAN (STA-GAN) to generate egocentric video from an exocentric perspective. The proposed framework consists of three parts: temporal generation branch, spatial generation branch, and attention fusion module, as shown in Figure 2. First, both temporal and spatial branches take a sequence of exocentric view frames as the input and generate a sequence of fake frames of egocentric view and their feature maps for each

module, respectively. The fake frames and the corresponding feature maps are generated simultaneously in both downstream and upstream directions. Therefore, each input frame are corresponding to four pairs of fake frames and feature maps (spatial and temporal modules in two directions). Second, the generated four pairs of fake frames and feature maps are fed into the attention fusion module. Finally, we generate the fused output of fake frames. The proposed framework is able to learn both spatial and temporal information from the forward and backward directions in the time domain simultaneously. We demonstrate that the proposed STA-GAN outperforms other baselines such as X-Fork [32], X-Seq [32], SelectionGAN [45], RecycleGAN [4], vid2vid [50] through extensive experimental evaluations. We establish state-of-the-art results on the Side2Ego and Top2Ego datasets [11]. To the best of our knowledge, we are the first to attempt to incorporate a bi-directional spatial temporal generative network for exocentric to egocentric cross-view video synthesis.

The contributions of this paper can be summarized as follows:

- A novel Bi-directional Spatial Temporal Attention Fusion Generative Adversarial Network (STA-GAN) is proposed. It aims at deploying temporal and spatial information in video and learning both spatial and temporal information between different views simultaneously.
- A group of novel downstream/upstream temporal and spatial loss functions are designed for neural network training. Moreover, a novel attention fusion module is proposed to fuse the generated fake frames to obtain refined final results. Meanwhile, a novel temporal and spatial dual-discriminator are proposed for network training.
- Experimental results on different cross-view datasets show the effectiveness of the proposed model. Our approach outperforms state-of-the-art results by a large margin for the cross-view exocentric to egocentric video synthesis. To the best of our knowledge, we are the first to attempt to tackle the cross-view exocentric to egocentric video generation task.

## 2 RELATED WORK

**Generative Adversarial Networks (GANs).** Over the last few years, GANs [16] have been shown effectively in many image generation and translation tasks [18, 24, 37, 39–44, 57]. For example, Isola et al. [18] propose Pix2Pix adversarial learning framework for paired image generation. Zhu et al. [57] introduce CycleGAN which developed cycle-consistency constraint to deal with unpaired image generation. However, these works aim to generate images which have a large degree of overlapping in the appearance and view with input images. Synthesis is much more challenging when the generation is conditioned on images with drastically different views. Recently, researchers investigate cross-view image generation problems [33]. This is a more challenging task since different views share little overlap information. To tackle this problem, Krishna et al. [32] propose X-Fork and X-Seq GAN-based architecture using an extra semantic map to facilitate generation. Tang et al. [45] propose a semantic-guided multi-channel attention selection module within a GAN framework for cross-view image generation. However, these methods are limited to cross-view image generation task, they are not able to generate satisfactory results for cross-view video generation.

**Egocentric Vision.** Egocentric vision has been recently explored in the computer vision field [2, 13–15, 27, 30, 31, 47, 54]. Aghazadeh et al. [2] propose an approach for discovering anomalous events from videos captured from a small camera attached to a person's chest. Fathi et al. [14] introduce a method for individuating social interactions in first-person videos collected during social events. Some recent works [13, 27, 30, 35, 47] have focused on activity analysis considering different scenarios (e.g., kitchen, office, home). Xu et al. [12] propose a semi-Siamese CNN architecture to address the person-level correspondences across first- and third-person videos. They formulate the problem as learning a joint embedding space for first- and third-person videos that considers both spatial- and motion-domain cues.

**Video-to-Video Synthesis.** There is few recent work investigate video generation problem [6, 9, 15, 36]. TGAN [36] directly generate video clips from noise by using generative adversarial networks. MoCoGAN [49] employ unsupervised adversarial training to decompose motion and content to control the image-to-video generation. Pan et al. [28] work on video-to-video translation to generate a sequence of frames from a sequence of aligned semantic representations. Some recent works such as RecycleGAN [4] and Vid2Vid [50] learn mapping between different videos and transferred motion between faces and from poses to body, respectively. Frameworks [25, 29, 48, 56] propose image generation networks for 3D view synthesis.

However, existing frameworks on video generation require the input and output video scenes sharing the similar architecture, which were insufficient for cross-view video generation. Particularly, exocentric to egocentric cross-view video generation has not yet been studied in literature yet. Our method investigates both cross-view generation and video generation in the exocentric to egocentric perspective setting, which is more challenging than various video generation problems. To the best of our knowledge, this is the first attempt in literature.

## 3 BI-DIRECTIONAL SPATIAL TEMPORAL ATTENTION FUSION GANS

In this section, we present the details of the proposed Bi-directional Spatial Temporal Attention Fusion GAN (STA-GAN). The overall framework of the proposed STA-GAN is illustrated in Figure 2, which contains four different modules, i.e., temporal generation module, spatial generation module, attention fusion, and a dual-discriminator. The bi-directional temporal generation module learns the temporal information of the target video. Meanwhile, the bi-directional spatial generation module models the spatial information of video frames. Moreover, the multi-generation attention fusion module fuses the information from temporal and spatial modules. Lastly, the proposed temporal and spatial dual discriminator aims to distinguish the generated videos from two spaces, i.e., temporal space and spatial space.

### 3.1 Semantic-guided Cross-view Video Generation

Our goal is to generate a video sequence of egocentric view $\mathbf{y} \equiv \{y_1, ..., y_t\}$ from a video sequence of exocentric view $\mathbf{x} \equiv \{x_1, ..., x_t\}$, where $y_t$ and $x_t$ are corresponding real video frames. Our task is to learn a video generator $G$ receives $\mathbf{x}$ and output $\tilde{\mathbf{y}}$ close to the real video $\mathbf{y}$. This process can be formulated as:

$$\tilde{\mathbf{y}} = G(\mathbf{x}). \tag{1}$$

However, the cross-view exocentric to egocentric video synthesis task is challenging due to several reasons. First, exocentric and egocentric views have little overlapping, which leads to ambiguity issues in the generation process. Second, the existing egocentric view datasets are rare and collected by wearable devices which leads huge amount of blurry videos. To alleviate both limitations, in this work, we employ a semantic-guided strategy. We incorporate semantic maps as a conditional guidance. Specifically, we adopt RefineNet [21, 22, 45] to generate semantic maps on both Side2Ego and Top2Ego datasets [11]. The generated semantic maps are used as the conditional input of the generator $G$, as shown in Figure 2. We concatenate the input video $\mathbf{x}$ from the exocentric view and the semantic map $\mathbf{s} \equiv \{s_1, ..., s_t\}$ from a egocentric view, and input them into the video generator $G$ and synthesize the egocentric view video sequence $\tilde{\mathbf{y}}$ as:

$$\tilde{\mathbf{y}} = G(\mathbf{x}, \mathbf{s}). \tag{2}$$

In this way, the semantic maps provide stronger supervision to guide the cross-view video synthesis.

### 3.2 Bi-directional Temporal Generation

Due to the complexity and particularity of video generation task, we take into account of time information in video sequences. Temporal information is crucial in video analysis, a video frame is usually highly correlated to adjacent frames. In this work, we enforce a temporal coherence between adjacent frames by integrating adjacent semantic maps as guidance input in both downstream and upstream directions, which are shown along the arrows in Figure 3(left). The conditional semantic map $s_1$ together with the input frame $x_1$ are input into the generator $G$, and produce the synthesized frame $\tilde{y}_1^1 = G(x_1, s_1)$. Then generated $\tilde{y}_1^1$ and the next semantic map $s_2$ are further fed into the generator $G$ which reconstructs a
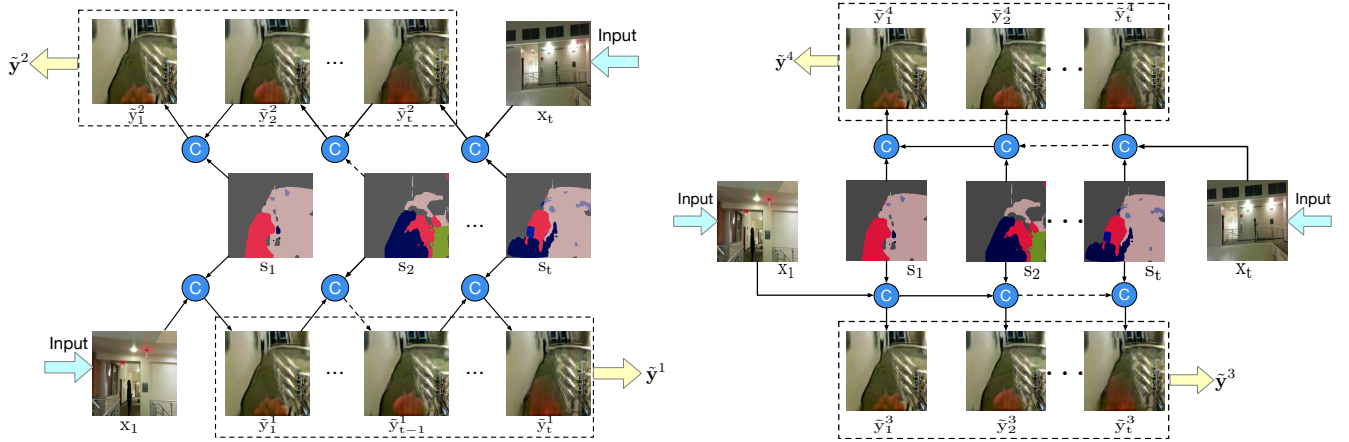
**Figure 3: Illustration of the proposed (left) temporal module, and (right) spatial module. The symbol ⓒ denotes concatenation. This channel-wise concatenation operating is able to avoid mixing of different modalities of information (RGB images and semantic maps). The concatenation will be used in the later attention fusion module.**

new egocentric frame $\tilde{y}_2^1$. We generate downstream direction of egocentric frame at time t based on $\tilde{y}_{t-1}^1$ and semantic map $s_t$. We can formalize the downstream process as:

$$\tilde{y}_t^1 = G(\tilde{y}_{t-1}^1, s_t). \tag{3}$$

In this way, we are able to obtain the generated video sequence $\tilde{\mathbf{y}}^1 \equiv \{\tilde{y}_1^1, ..., \tilde{y}_t^1\}$.

In the opposite direction, the upstream process is time reverse generation process and starts from time t. We input $x_t$ and $s_t$ to the generator $G$ to produce the first frame $\tilde{y}_t^2 = G(x_t, s_t)$. Then other frames can be obtained by using the following formulation:

$$\tilde{y}_{t-1}^2 = G(\tilde{y}_t^2, s_{t-1}). \tag{4}$$

By doing so, we are able to obtain the generated video sequence $\tilde{\mathbf{y}}^2 \equiv \{\tilde{y}_1^2, ..., \tilde{y}_t^2\}$.

To better model the temporal information, we further propose a novel temporal loss. The temporal loss composes of both downstream process and upstream process. The key idea is to add temporal constraints to the generation process. For each frame, the downstream temporal loss is:

$$\mathcal{T}_{dn}(\tilde{y}_t^1) = \mathbb{E}_{x_t, \tilde{y}_t^1} \left[ \|y_t - \tilde{y}_t^1\|_1 \right] + \mathbb{E}_{x_t, y_t} \left[ \log D_S(x_t, y_t) \right] \\ + \mathbb{E}_{x_t, \tilde{y}_t^1} \left[ \log(1 - D_S(x_t, \tilde{y}_t^1)) \right]. \tag{5}$$

Overall, the downstream temporal loss can be represented as:

$$\mathcal{T}_{dn}(\mathbf{x}, \tilde{\mathbf{y}}^1) = \sum_{t=1}^{t} \lambda_u \mathcal{T}_{dn}(\tilde{y}_t^1). \tag{6}$$

The upstream temporal loss is formulated similar as downstream temporal loss, however they computed in an opposite direction for each frame as:

$$\mathcal{T}_{up}(\tilde{y}_t^2) = \mathbb{E}_{x_t, \tilde{y}_t^2} \left[ \|y_t - \tilde{y}_t^2\|_1 \right] + \mathbb{E}_{x_t, y_t} \left[ \log D_S(x_t, y_t) \right] \\ + \mathbb{E}_{x_t, \tilde{y}_t^2} \left[ \log(1 - D_S(x_t, \tilde{y}_t^2)) \right]. \tag{7}$$

The overall upstream temporal loss is:

$$\mathcal{T}_{up}(\mathbf{x}, \tilde{\mathbf{y}}^2) = \sum_{t=t}^{1} \lambda_d \mathcal{T}_{up}(\tilde{y}_t^2). \tag{8}$$

Finally, the temporal loss is the sum of Equation (6) and Equation (8),

$$\mathcal{L}_T(\mathbf{x}, \tilde{\mathbf{y}}) = \mathcal{T}_{dn}(\mathbf{x}, \tilde{\mathbf{y}}^1) + \mathcal{T}_{up}(\mathbf{x}, \tilde{\mathbf{y}}^2). \tag{9}$$

## 3.3 Bi-directional Spatial Generation

Spatial information plays a crucial role in various video related tasks, such as activity recognition [51], object recognition [7], etc. In this work, we incorporate the effects of spatial information by generating non-adjacent frames using the corresponding semantic maps. As illustrated in the spatial module of Figure 3(right), the process of downstream direction is generated along the arrows:

$$\tilde{y}_1^3 = G(x_1, s_1), \cdots, \tilde{y}_t^3 = G(x_1, s_t). \tag{10}$$

In this way, we are able to obtain the generated video sequence $\tilde{\mathbf{y}}^3 \equiv \{\tilde{y}_1^3, ..., \tilde{y}_t^3\}$. The opposite upstream direction in which the sequences is formulated as:

$$\tilde{y}_t^4 = G(x_t, s_t), \cdots, \tilde{y}_1^4 = G(x_t, s_1), \tag{11}$$

where we are able to obtain the generated video sequence $\tilde{\mathbf{y}}^4 \equiv \{\tilde{y}_1^4, ..., \tilde{y}_t^4\}$.

To learn the spatial information better, we propose a new spatial loss. The spatial loss composes of two parts, which are downstream spatial loss and upstream spatial loss. The intuition is that activities and events in videos are spatially related between adjacent frames and are reversible. For each frame, the downstream spatial loss is formulated as follows:

$$\mathcal{S}_{dn}(\tilde{y}_t^3) = \mathbb{E}_{x_{t-i}, \tilde{y}_t^3} \left[ \|y_t - \tilde{y}_t^3\|_1 \right] + \mathbb{E}_{x_{t-i}, y_t} \left[ \log D_S(x_{t-i}, y_t) \right] \\ + \mathbb{E}_{x_{t-i}, \tilde{y}_t^3} \left[ \log(1 - D_S(x_{t-i}, \tilde{y}_t^3)) \right]. \tag{12}$$

The overall downstream spatial loss is:

$$\mathcal{S}_{dn}(\mathbf{x}, \tilde{\mathbf{y}}^3) = \sum_{t=1}^{t} \lambda_n \mathcal{S}_{dn}(\tilde{y}_t^3). \tag{13}$$

The upstream spatial loss is generated reversely, which is formulated similar as:

$$\mathcal{S}_{up}(\tilde{y}_t^4) = \mathbb{E}_{x_{t+i}, \tilde{y}_t^4} \left[ \|y_t - \tilde{y}_t^4\|_1 \right] + \mathbb{E}_{x_{t+i}, y_t} \left[ \log D_S(x_{t+i}, y_t) \right] \\ + \mathbb{E}_{x_{t+i}, \tilde{y}_t^4} \left[ \log(1 - D_S(x_{t+i}, \tilde{y}_t^4)) \right]. \tag{14}$$

**Figure 4: Illustration of the proposed multi-generation attention fusion module. The symbol ⊕, ⊗, ⓒ and ⓢ denote element-wise addition, element-wise multiplication, channel-wise concatenation and channel-wise softmax operation.**

The overall upstream spatial loss:

$$\mathcal{S}_{up}(\mathbf{x}, \tilde{\mathbf{y}}^4) = \sum_{t=t}^{1} \lambda_p \mathcal{T}_{up}(\tilde{y}_t^4).\tag{15}$$

The overall spatial loss is the sum of Equation (13) and Equation (15), where $i$ is the time truncate during training:

$$\mathcal{L}_S(\mathbf{x}, \tilde{\mathbf{y}}) = \mathcal{S}_{dn}(\mathbf{x}, \tilde{\mathbf{y}}^3) + \mathcal{S}_{up}(\mathbf{x}, \tilde{\mathbf{y}}^4).\tag{16}$$

### 3.4 Multi-generation Attention Fusion

After the Bi-directional temporal and spatial modules, we obtain four generated video sequences $[\tilde{\mathbf{y}}^1, \tilde{\mathbf{y}}^2, \tilde{\mathbf{y}}^3, \tilde{\mathbf{y}}^4]$. To combine relevant information from the four generated video sequences, we propose a novel multi-generation attention fusion module to obtain a refined video sequence. The details of the proposed multi-generation attention fusion module is shown in Figure 4. First, four different features are extracted by four convolutional layers simultaneously from the four generated video sequences. Then the obtained four feature maps are concatenated to a new feature as:

$$\mathcal{F} = \text{Concat}(\text{Conv}(\tilde{\mathbf{y}}^1), \text{Conv}(\tilde{\mathbf{y}}^2), \text{Conv}(\tilde{\mathbf{y}}^3), \text{Conv}(\tilde{\mathbf{y}}^4)),\quad(17)$$

where $\text{Concat}(\cdot)$ and $\text{Conv}(\cdot)$ denote channel-wise concatenation operation and convolutional operation. Next, the concatenated feature $\mathcal{F}$ is fed into a de-convolutional layer to obtain the new size feature $\mathcal{F}'$ for attention fusion purpose. Then the attention maps

$[\mathbf{a}^1, \mathbf{a}^2, \mathbf{a}^3, \mathbf{a}^4]$ are learned from a convolutional layer and softmax activation layer, where the softmax activation layer guarantees normalization of attention maps in channel-wise. Finally, the final video sequence $\tilde{\mathbf{y}}$ can be obtained as follows:

$$\tilde{\mathbf{y}} = (\tilde{\mathbf{y}}^1 \otimes \mathbf{a}^1) \oplus (\tilde{\mathbf{y}}^2 \otimes \mathbf{a}^2) \oplus (\tilde{\mathbf{y}}^3 \otimes \mathbf{a}^3) \oplus (\tilde{\mathbf{y}}^4 \otimes \mathbf{a}^4),\tag{18}$$

where $\tilde{\mathbf{y}}$ represents the final synthesized frame sequence, and the symbol $\otimes$ denotes the element-wise multiplication, and $\oplus$ is the element-wise addition.

Instead of multiplying feature maps and real images, we multiply feature maps $(\mathbf{a}^1, \mathbf{a}^2, \mathbf{a}^3, \mathbf{a}^4)$ and the intermediate results $(\tilde{\mathbf{y}}^1, \tilde{\mathbf{y}}^2, \tilde{\mathbf{y}}^3, \tilde{\mathbf{y}}^4)$ to obtain the final result $\tilde{\mathbf{y}}$. We observe that a single-sequence generation space may not be suitable enough for learning a good mapping for this complex synthesis problem. Thus, we explore a larger sequence generation space to have a richer synthesis via constructing multiple intermediate sequence results. We first produce several intermediate sequences which can be regarded as the candidate sequence pool, then learn a set of attention maps. These attention maps are used to spatially select from the intermediate generations and are combined to synthesize the final output.

### 3.5 Temporal and Spatial Dual-discriminator

Traditional image-to-image generation methods use vanilla discriminator [18]. In this paper, we propose a novel Temporal and Spatial Dual-discriminator, which contains two discriminators, i.e., temporal discriminator ($D_T$) and spatial discriminator ($D_S$). $D_T$ takes the real input $x_1/x_t$ and the generated image sequence $\tilde{\mathbf{y}}$ as input. $D_T$ learns to tell whether a sequence of fake output $\tilde{\mathbf{y}}$ and input $x_1/x_t$ are related or not. Meanwhile, $D_S$ takes the real input image x and fake $\tilde{\mathbf{y}}$ as input. $D_S$ learns to tell whether the two frame sequences from different domains are associated with each other or not. $D_T$ and $D_S$ take temporal and spatial information into consideration, respectively.

Assuming we target to learn a mapping $G: \mathbf{x} \to \tilde{\mathbf{y}}$ from input exocentric view $\mathbf{x} \equiv \{x_1, ..., x_t\}$ to output egocentric view $\tilde{\mathbf{y}} \equiv \{\tilde{y}_1, ..., \tilde{y}_t\}$. The generator $G$ is trained to produce fake outputs $\tilde{\mathbf{y}}$ to fool the discriminator $D_T$ and $D_S$. The adversarial loss can be expressed as:

$$\begin{aligned}\mathcal{L}_{cGAN_S}(\mathbf{x}, \tilde{\mathbf{y}}) = \sum_{t=1}^{t} (&\mathbb{E}_{x_t, y_t}\left[\log D_S(x_t, y_t)\right] \\ &+ \mathbb{E}_{x_t, \tilde{y}_t}\left[\log(1 - D_S(x_t, \tilde{y}_t))\right]).\end{aligned}\tag{19}$$

$$\begin{aligned}\mathcal{L}_{cGAN_T}(\mathbf{x}, \tilde{\mathbf{y}}) = &\mathbb{E}_{x, y}\left[\log D_T(x_t, \mathbf{y})\right] \\ &+ \mathbb{E}_{x, \tilde{y}}\left[\log(1 - D_T(x_t, \tilde{\mathbf{y}}))\right] \\ &+ \mathbb{E}_{x, y}\left[\log D_T(x_1, \mathbf{y})\right] \\ &+ \mathbb{E}_{x, \tilde{y}}\left[\log(1 - D_T(x_1, \tilde{\mathbf{y}}))\right].\end{aligned}\tag{20}$$

$x_1/x_t$ in Equation (20) is the starting and ending frame in the temporal synthesis. The total adversarial loss is formulated as follows:

$$\mathcal{L}_{cGAN}(\mathbf{x}, \tilde{\mathbf{y}}) = \mathcal{L}_{cGAN_S} + \lambda_g \mathcal{L}_{cGAN_T}.\tag{21}$$

### 3.6 Optimization Objective

The training objective can be decomposed into four main components which are adversarial loss, temporal loss, spatial loss and reconstruction loss.

**Figure 5: Video frames generated from exocentric view to egocentric view on Side2Ego dataset.**
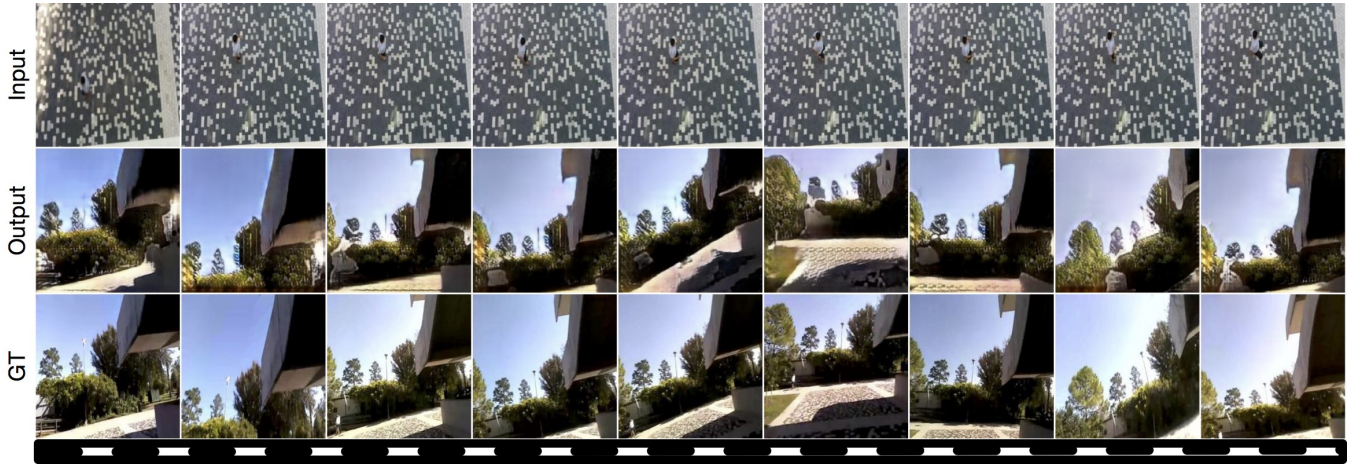


**Figure 6: Video frames generated from exocentric view to egocentric view on Top2Ego dataset.**

**Reconstruction Loss.** The task of the generator is to reconstruct an video sequences $\tilde{\mathbf{y}} \equiv \{\tilde{y}_1, ..., \tilde{y}_t\}$ as close as the target sequences $\mathbf{y} \equiv \{y_1, ..., y_t\}$. We use $\mathcal{L}1$ distance denoted as $\| \cdot \|_1$ in the reconstruction loss:

$$\mathcal{L}_{re}(G) = \sum_{t=1}^{t} \lambda_r \mathbb{E}_{\mathbf{x},\tilde{\mathbf{y}}} \left[ \|y_t - \tilde{y}_t)\|_1 \right]. \tag{22}$$

**Overall Loss.** The final optimization loss is a weighted sum of the above losses. Generator $G$ and discriminator $D$ are trained in an end-to-end fashion to optimize the following objective function:

$$\mathcal{L} = \mathcal{L}_{cGAN} + \mathcal{L}_{re} + \mathcal{L}_T + \mathcal{L}_S, \tag{23}$$

where $\lambda_i$'s in Equations (6), (8), (13), (15), (21), (22) are the regularization parameters.

**Network Architecture.** We employ U-Net [34] as the architecture of our generator $G$. We impose the skip connection strategy from down-sampling path to up-sampling path to avoid the vanishing gradient problem. We adopt PatchGAN [18] for the discriminators $D_S$ and $D_T$. The feature maps for attention fusion are extracted by the up-sampling layers of the U-Net during the training of generator

$G$. We adopt RefineNet [22] to generate semantic maps on the Side2Ego and Top2Ego datasets as in [32, 45].

## 4 EXPERIMENTS

**Datasets.** To explore the effectiveness of the proposed STA-GAN, we conduct extensive experiments on the Side2Ego and Top2Ego datasets [11]. These datasets simultaneously recorded egocentric and exocentric videos. Each video pair contains one egocentric and one exocentric video (side-view or top-view). The pair of videos are temporally aligned. These datasets are challenging due to two reasons. First, it contains dramatically different indoor and outdoor scenes. Second, the datasets are collected simultaneously by an exocentric camera (side or top view) and an egocentric wearable camera. The datasets include a huge amount of blurred image frames for egocentric view. For Side2Ego dataset, there are 124 videos containing 26,764 pairs of frames for training and 13,788 pairs for testing. For Top2Ego dataset, there are 135 videos containing 28,408 pairs of frames for training and 14,064 pairs for testing. All image frames are in high-resolution with 1280×720 pixels.

**Figure 7: Video frames generated from exocentric view to egocentric view on Side2Ego dataset using different methods.**

**Table 1: Quantitative evaluation of different image and video generation methods on the Side2Ego dataset. For these metrics except KL score and FID, higher is better.**

| Method | SSIM ↑ | PSNR ↑ | SD ↑ | KL ↓ | FID ↓ | Top-1 ↑ Accuracy (%) | | Top-5 ↑ Accuracy (%) | |
|---|---|---|---|---|---|---|---|---|---|
| X-Fork [32] | 0.4499 | 17.0743 | 20.4443 | 51.20 ± 1.94 | 216.5575 | 4.49 | 9.76 | 11.63 | 19.44 |
| X-Seq [32] | 0.4763 | 17.1462 | 20.7468 | 45.10 ± 1.95 | 184.2808 | 6.51 | 12.70 | 11.97 | 19.36 |
| SelectionGAN [45] | 0.5128 | 18.3021 | 20.9426 | **7.26 ± 1.27** | **139.1429** | 20.84 | 37.49 | 42.51 | 65.22 |
| RecycleGAN [4] | 0.3446 | 15.9242 | 18.9429 | 42.40 ±1.61 | 186.5897 | 2.32 | 2.40 | 9.13 | 10.98 |
| Vid2Vid [50] | 0.3955 | 15.9012 | 19.7169 | 59.41 ± 1.93 | 196.9749 | 7.52 | 14.70 | 13.97 | 24.36 |
| STA-GAN (Ours) | **0.5607** | **20.7027** | **20.9491** | 9.44 ± 1.48 | 169.3514 | **26.83** | **39.83** | **42.72** | **69.30** |

**Table 2: Quantitative evaluation of different image and video generation methods on the Top2Ego dataset. For these metrics except KL score and FID, higher is better.**

| Method | SSIM ↑ | PSNR ↑ | SD ↑ | KL ↓ | FID ↓ | Top-1 ↑ Accuracy (%) | | Top-5 ↑ Accuracy (%) | |
|---|---|---|---|---|---|---|---|---|---|
| X-Fork [32] | 0.2952 | 15.8849 | 18.7349 | 63.96±1.74 | 208.6632 | 0.8 | 1.22 | 3.16 | 4.08 |
| X-Seq [32] | 0.3522 | 16.9439 | 19.2733 | 54.91 ± 1.81 | 198.3948 | 1.07 | 1.77 | 4.29 | 6.94 |
| SelectionGAN [45] | 0.5047 | 22.0244 | 19.1976 | 10.07 ± 1.29 | 199.1441 | 8.85 | 16.55 | 24.32 | 33.90 |
| RecycleGAN [4] | 0.3264 | 16.4767 | 19.4543 | 31.55 ±1.32 | 235.8220 | 0.43 | 2.04 | 2.47 | 3.70 |
| Vid2Vid [50] | 0.3895 | 17.1233 | 19.6043 | 31.48±1.77 | 208.9400 | 5.65 | 9.77 | 12.38 | 18.97 |
| STA-GAN (Ours) | **0.5383** | **22.5816** | 19.2895 | **10.02 ±1.30** | **175.7446** | **8.93** | **24.80** | **26.37** | **46.74** |

**Parameter Settings.** All images are scaled to 256×256. We enable image flipping and random crops for data augmentation. We train 200 epochs with the batch size of 8. In our experiments, we set $\lambda_u$=1, $\lambda_d$=0.1, $\lambda_n$=1, $\lambda_p$=0.1, $\lambda_g$=10, $\lambda_r$=10 in Equations (6), (8), (13), (15), (21), (22) respectively. The number of time truncat $i$ in Eq. (12) and (14) is set to 3. The proposed STA-GAN is implemented by PyTorch. We perform our experiments on NVIDIA Geforce GTX 1080 Ti GPU with 11 GB memory to accelerate training process.

**Evaluation Metrics.** We follow [32, 45, 50] and apply metrics such as top-k prediction accuracy, KL score and Fréchet Inception Distance (FID) for evaluating the proposed method. These metrics evaluate the generated images in a high-level feature space. We

also employ pixel-level similarity metrics in the experiments, i.e., Structural-Similarity (SSIM) [52], Peak Signal-to-Noise Ratio (PSNR) and Sharpness Difference (SD) [26].

### 4.1 State-of-the-Art Comparisons

**Quantitative Comparisons.** We compare our STA-GAN with both cross-view image generation methods, i.e., X-Fork [32], X-Seq [32], SelectionGAN [45], and video generation methods, i.e., RecycleGAN [4] and Vid2Vid [50]. The quantitative results compared with state-of-the-art methods on the Side2Ego and Top2Ego datasets are presented in Table 1 and Table 2. We observe that the proposed STA-GAN achieves better results than state-of-the-art

| Top | Semantic Map | X-Seq | X-Fork | SelectionGAN | RecycleGAN | Vid2Vid | **Ours** | GT |



**Figure 8: Video frames generated from exocentric view to egocentric view on Top2Ego dataset using different methods.**

methods in most cases. Accuracies are computed in two ways, the first column considers all images, whereas the second column computes accuracy of real images whose top-1 prediction is greater than 0.5. Particularly, we increase the metric PSNR by 2.4, Top-1 accuracy by 6%/2%, Top-5 accuracy by 4%/5%, compared to the second best baseline on Side2Ego dataset. Meanwhile, we increase the metric PSNR by 0.5, Top-1 accuracy by 0.1%/8%, Top-5 accuracy by 2%/13%, compared with the second best baseline on Top2Ego dataset. These comparisons show the effectiveness of STA-GAN for cross-view exocentric to egocentric video synthesis.

**Qualitative Comparisons.** Video sequence synthesis results are shown in Figure 5 and Figure 6. We can observe that the synthesized egocentric video sequence change smoothly from left to right, and are visually close to the ground truth. The qualitative results compared with state-of-the-art methods are shown in Figure 7 and Figure 8. We observe that our method generates more clear and reasonable video frames than other methods. It is obvious that objects are in the correct positions for generated egocentric video frames using our method. Moreover, the structure and layout of the generated video frames are closer to the ground truth. Results show that egocentric video frames generated by STA-GAN are visually better compared with other baselines, which further confirm that the proposed STA-GAN network has the ability to transfer the video sequences from exocentric to egocentric perspective.

### 4.2 Ablation Study

To evaluate the performance of proposed STA-GAN, we conduct experiments with different settings on the Top2Ego dataset. As shown in Table 3, the proposed STA-GAN considers six different settings. Baseline A utilizes only spatial information during generation process while baseline B utilizes only temporal information. We observe that baseline B performs better than baseline A which demonstrates that temporal information is more important for video generation task. Baseline C by considering both temporal and spatial information improves the SSIM, PSNR and SD metrics to 0.3098/17.0236/18.6043 respectively, meaning that both spatial

**Table 3: Ablation study of STA-GAN on the Top2Ego dataset.**

| No. | Setting of STA-GAN | SSIM ↑ | PSNR ↑ | SD ↑ |
|-----|-------------------|--------|--------|------|
| A | Spatial Generation | 0.2568 | 15.8561 | 18.1414 |
| B | Temporal Generation | 0.2627 | 15.3411 | 18.1914 |
| C | Spatial + Temporal | 0.3098 | 17.0236 | 18.6043 |
| D | C + Bi-Direction | 0.4287 | 20.2891 | 19.2389 |
| E | D + Dual-Discriminator | 0.4956 | **23.4734** | 19.2526 |
| F | E + Attention Module | **0.5383** | 22.5816 | **19.2895** |

and temporal information should be explored for video generation task. Baseline D by incorporating both bi-directional downstream and upstream generation further improves the performance to 0.4287/20.2891/19.2389, which demonstrates that generation are reversible in video generation task. Baseline E outperforms D showing the importance of using the proposed dual-discriminator, i.e., temporal discriminator and spatial discriminator. Baseline F by adopting the attention fusion strategy further increases the SSIM, PSNR and SD scores, which demonstrates the effectiveness of attention fusion.

## 5 CONCLUSION

In this paper, we propose a novel STA-GAN framework to address a novel cross-view exocentric to egocentric video synthesize problem by exploiting the temporal and spatial information in videos. Based on the property of videos, we propose a bi-directional strategy which generates video sequences in both downstream and upstream directions. Meanwhile, a novel temporal and spatial dual-discriminator is proposed for better network training. Moreover, we propose a novel attention fusion method which targets to refine the generation results. Extensive experimental results on the Top2Ego and Side2Ego datasets demonstrate that our method outperforms state-of-the-art approaches for the challenging cross-view exocentric to egocentric video synthesis.

# REFERENCES

[1] [n.d.]. https://gopro.com/en/us/.
[2] Omid Aghazadeh, Josephine Sullivan, and Stefan Carlsson. 2011. Novelty detection from an ego-centric perspective. In *CVPR*.
[3] Shervin Ardeshir and Ali Borji. 2016. Ego2top: Matching viewers in egocentric and top-view videos. In *ECCV*.
[4] Aayush Bansal, Shugao Ma, Deva Ramanan, and Yaser Sheikh. 2018. Recycle-gan: Unsupervised video retargeting. In *ECCV*.
[5] Andrew Brock, Jeff Donahue, and Karen Simonyan. 2019. Large scale GAN training for high fidelity natural image synthesis. In *ICLR*.
[6] Caroline Chan, Shiry Ginosar, Tinghui Zhou, and Alexei A Efros. 2019. Everybody dance now. In *ICCV*.
[7] Kai Chen, Jiaqi Wang, Shuo Yang, Xingcheng Zhang, Yuanjun Xiong, Chen Change Loy, and Dahua Lin. 2018. Optimizing video object detection via a scale-time lattice. In *CVPR*.
[8] Yunjey Choi, Youngjung Uh, Jaejun Yoo, and Jung-Woo Ha. 2020. Stargan v2: Diverse image synthesis for multiple domains. In *CVPR*.
[9] Emily Denton and Rob Fergus. 2018. Stochastic video generation with a learned prior. In *ICML*.
[10] Bin Duan, Wei Wang, Hao Tang, Hugo Latapie, and Yan Yan. 2021. Cascade attention guided residue learning gan for cross-modal translation. In *ICPR*.
[11] Mohamed Elfeki, Krishna Regmi, Shervin Ardeshir, and Ali Borji. 2019. From third person to first person: Dataset and baselines for synthesis and retrieval. In *CVPR*.
[12] Chenyou Fan, Jangwon Lee, Mingze Xu, Krishna Kumar Singh, Yong Jae Lee, David J Crandall, and Michael S Ryoo. 2017. Identifying first-person camera wearers in third-person videos. In *CVPR*.
[13] Alireza Fathi, Ali Farhadi, and James M Rehg. 2011. Understanding egocentric activities. In *ICCV*.
[14] Alircza Fathi, Jessica K Hodgins, and James M Rehg. 2012. Social interactions: A first-person perspective. In *CVPR*.
[15] Alireza Fathi, Yin Li, and James M Rehg. 2012. Learning to recognize daily actions using gaze. In *ECCV*.
[16] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. 2014. Generative adversarial nets. (2014).
[17] Minyoung Huh, Shao-Hua Sun, and Ning Zhang. 2019. Feedback adversarial learning: Spatial feedback for improving generative adversarial networks. In *CVPR*.
[18] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. 2017. Image-to-image translation with conditional adversarial networks. In *CVPR*.
[19] Takeo Kanade and Martial Hebert. 2012. First-person vision. *Proc. IEEE* 100, 8 (2012), 2442–2453.
[20] Tero Karras, Samuli Laine, and Timo Aila. 2019. A style-based generator architecture for generative adversarial networks. In *CVPR*.
[21] Guosheng Lin, Fayao Liu, Anton Milan, Chunhua Shen, and Ian Reid. 2019. Refinenet: Multi-path refinement networks for dense prediction. *IEEE TPAMI* 42, 5 (2019), 1228–1242.
[22] Guosheng Lin, Anton Milan, Chunhua Shen, and Ian Reid. 2017. Refinenet: Multi-path refinement networks for high-resolution semantic segmentation. In *CVPR*.
[23] Gaowen Liu, Hao Tang, Hugo Latapie, and Yan Yan. 2020. Exocentric to egocentric image generation via parallel generative adversarial network. In *ICASSP*.
[24] Ming-Yu Liu, Xun Huang, Arun Mallya, Tero Karras, Timo Aila, Jaakko Lehtinen, and Jan Kautz. 2019. Few-shot unsupervised image-to-image translation. In *ICCV*.
[25] Ming-Yu Liu and Oncel Tuzel. 2016. Coupled generative adversarial networks. *NeurIPS* (2016).
[26] Michael Mathieu, Camille Couprie, and Yann LeCun. 2016. Deep multi-scale video prediction beyond mean square error. In *ICLR*.
[27] Keisuke Ogaki, Kris M Kitani, Yusuke Sugano, and Yoichi Sato. 2012. Coupling eye-motion and ego-motion features for first-person activity recognition. In *CVPR Workshops*.
[28] Junting Pan, Chengyu Wang, Xu Jia, Jing Shao, Lu Sheng, Junjie Yan, and Xiaogang Wang. 2019. Video generation from single semantic label map. In *CVPR*.
[29] Eunbyung Park, Jimei Yang, Ersin Yumer, Duygu Ceylan, and Alexander C Berg. 2017. Transformation-grounded image generation network for novel 3d view synthesis. In *CVPR*.
[30] Hamed Pirsiavash and Deva Ramanan. 2012. Detecting activities of daily living in first-person camera views. In *CVPR*.
[31] Yair Poleg, Chetan Arora, and Shmuel Peleg. 2014. Temporal segmentation of egocentric videos. In *CVPR*.
[32] Krishna Regmi and Ali Borji. 2018. Cross-view image synthesis using conditional gans. In *CVPR*.
[33] Krishna Regmi and Ali Borji. 2019. Cross-view image synthesis using geometry-guided conditional gans. *Elsevier CVIU* 187 (2019), 102788.

[34] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. 2015. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*.
[35] Michael S Ryoo and Larry Matthies. 2013. First-person activity recognition: What are they doing to me?. In *CVPR*.
[36] Masaki Saito, Eiichi Matsumoto, and Shunta Saito. 2017. Temporal generative adversarial nets with singular value clipping. In *ICCV*.
[37] Tamar Rott Shaham, Tali Dekel, and Tomer Michaeli. 2019. Singan: Learning a generative model from a single natural image. In *ICCV*.
[38] Firas Shama, Roey Mechrez, Alon Shoshan, and Lihi Zelnik-Manor. 2019. Adversarial feedback loop. In *ICCV*.
[39] Hao Tang, Song Bai, and Nicu Sebe. 2020. Dual attention gans for semantic image synthesis. In *ACM MM*.
[40] Hao Tang, Song Bai, Philip HS Torr, and Nicu Sebe. 2020. Bipartite graph reasoning gans for person image generation. In *BMVC*.
[41] Hao Tang, Song Bai, Li Zhang, Philip HS Torr, and Nicu Sebe. 2020. Xinggan for person image generation. In *ECCV*.
[42] Hao Tang, Hong Liu, and Nicu Sebe. 2020. Unified generative adversarial networks for controllable image-to-image translation. *IEEE TIP* 29 (2020), 8916–8929.
[43] Hao Tang, Wei Wang, Dan Xu, Yan Yan, and Nicu Sebe. 2018. Gesturegan for hand gesture-to-gesture translation in the wild. In *ACM MM*.
[44] Hao Tang, Dan Xu, Gaowen Liu, Wei Wang, Nicu Sebe, and Yan Yan. 2019. Cycle in cycle generative adversarial networks for keypoint-guided image generation. In *ACM MM*.
[45] Hao Tang, Dan Xu, Nicu Sebe, Yanzhi Wang, Jason J Corso, and Yan Yan. 2019. Multi-channel attention selection gan with cascaded semantic guidance for cross-view image translation. In *CVPR*.
[46] Hao Tang, Dan Xu, Yan Yan, Philip HS Torr, and Nicu Sebe. 2020. Local class-specific and global image-level generative adversarial networks for semantic-guided scene generation. In *CVPR*.
[47] Ekaterina Taralova, Fernando De la Torre, and Martial Hebert. 2011. Source constrained clustering. In *ICCV*.
[48] Maxim Tatarchenko, Alexey Dosovitskiy, and Thomas Brox. 2016. Multi-view 3d models from single images with a convolutional network. In *ECCV*.
[49] Sergey Tulyakov, Ming-Yu Liu, Xiaodong Yang, and Jan Kautz. 2018. Mocogan: Decomposing motion and content for video generation. In *CVPR*.
[50] Ting-Chun Wang, Ming-Yu Liu, Jun-Yan Zhu, Guilin Liu, Andrew Tao, Jan Kautz, and Bryan Catanzaro. 2018. Video-to-video synthesis. In *NeurIPS*.
[51] Xiaolong Wang, Allan Jabri, and Alexei A Efros. 2019. Learning correspondence from the cycle-consistency of time. In *CVPR*.
[52] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. 2004. Image quality assessment: from error visibility to structural similarity. *IEEE TIP* 13, 4 (2004), 600–612.
[53] Yan Yan, Elisa Ricci, Gaowen Liu, and Nicu Sebe. 2015. Egocentric daily activity recognition via multitask clustering. *IEEE TIP* 24, 10 (2015), 2984–2995.
[54] Ryo Yonetani, Kris M Kitani, and Yoichi Sato. 2016. Visual motif discovery via first-person vision. In *ECCV*.
[55] Han Zhang, Ian Goodfellow, Dimitris Metaxas, and Augustus Odena. 2019. Self-attention generative adversarial networks. In *ICML*.
[56] Tinghui Zhou, Shubham Tulsiani, Weilun Sun, Jitendra Malik, and Alexei A Efros. 2016. View synthesis by appearance flow. In *ECCV*.
[57] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. 2017. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *ICCV*.