

# Hierarchical Video Generation from Orthogonal Information: Optical flow and Texture

---

**Katsunori Ohnishi**<sup>\*§1</sup>, Shohei Yamamoto<sup>\*1</sup>,  
Yoshitaka Ushiku<sup>1</sup>, and Tatsuya Harada<sup>12</sup>

<sup>1</sup> The University of Tokyo    <sup>2</sup> RIKEN

\* indicates equal contribution.

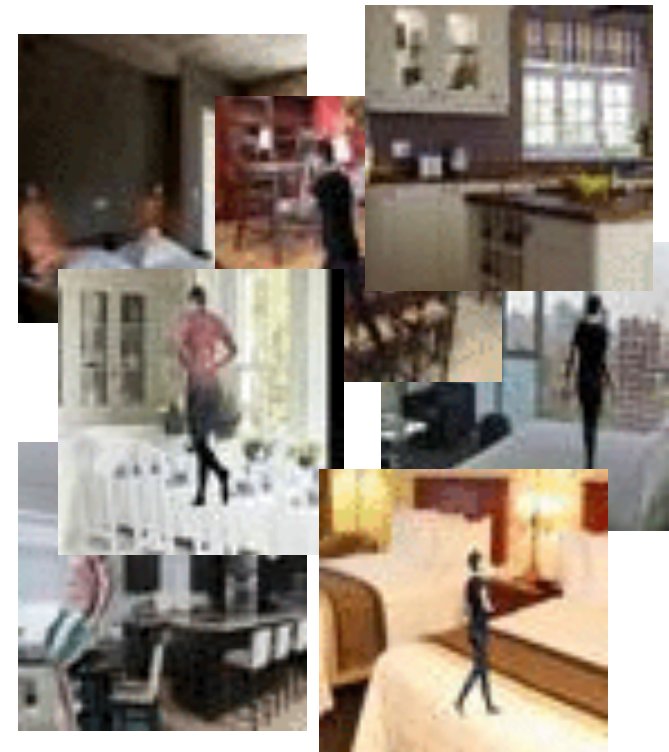
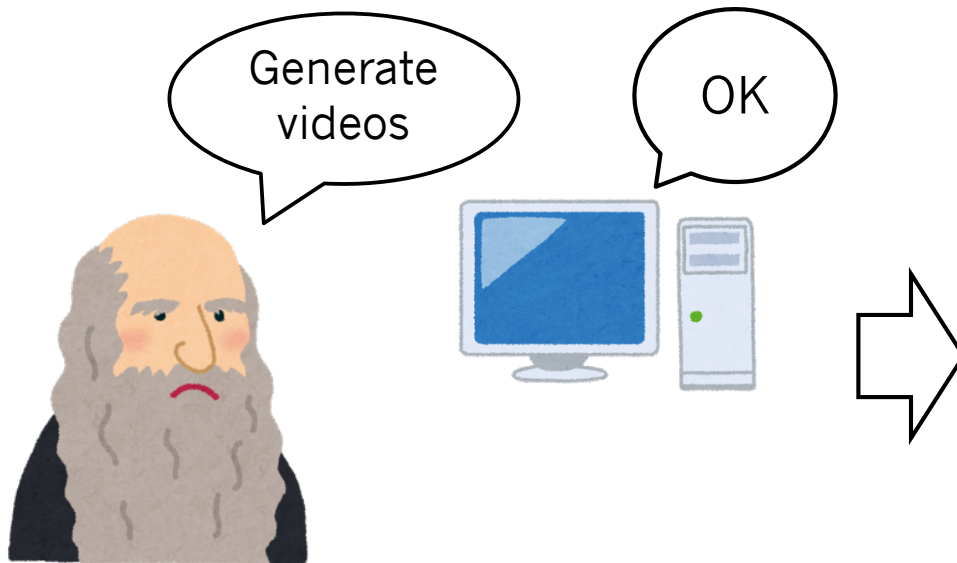
§ currently belongs to DeNA Co., Ltd.

Paper&Slides: <http://katsunoriohnishi.github.io/>

# Goal

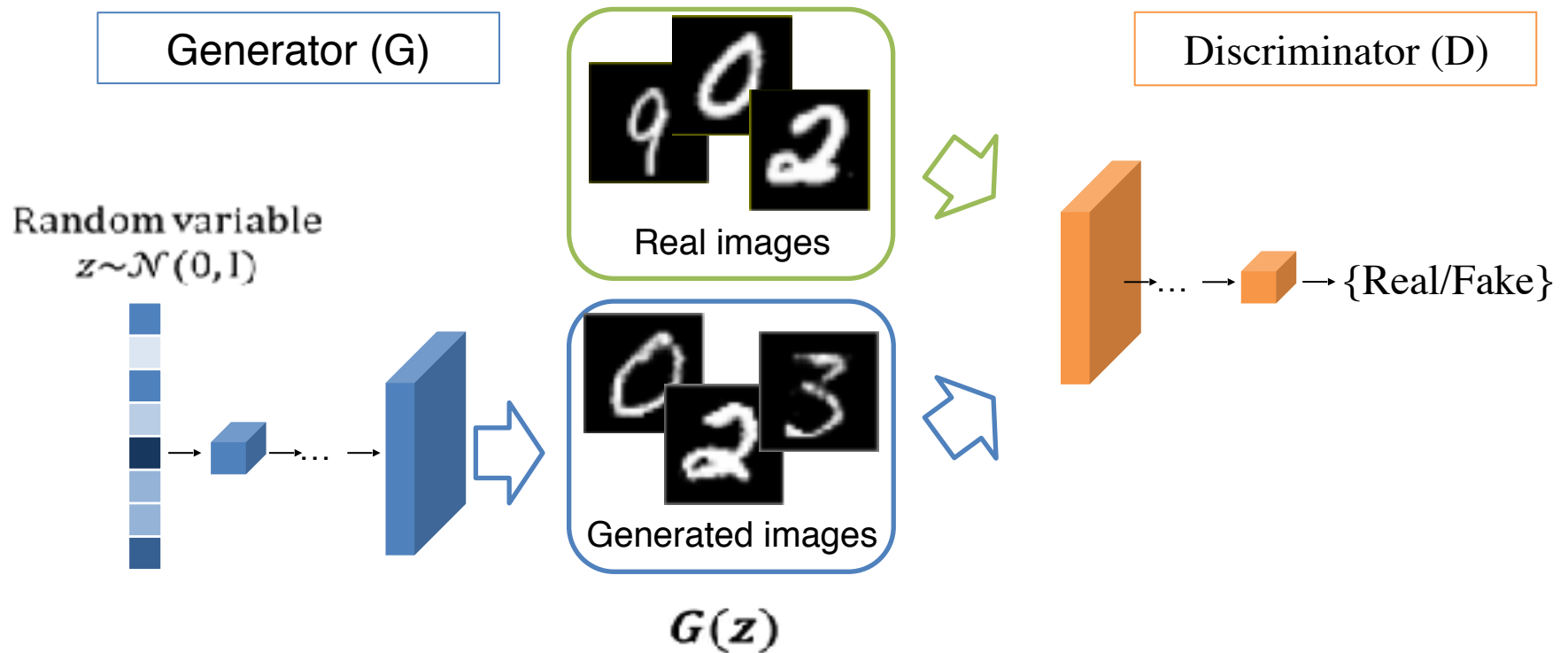
---

- Video generation
  - Applications)  
Human AI collaboration, dataset extension



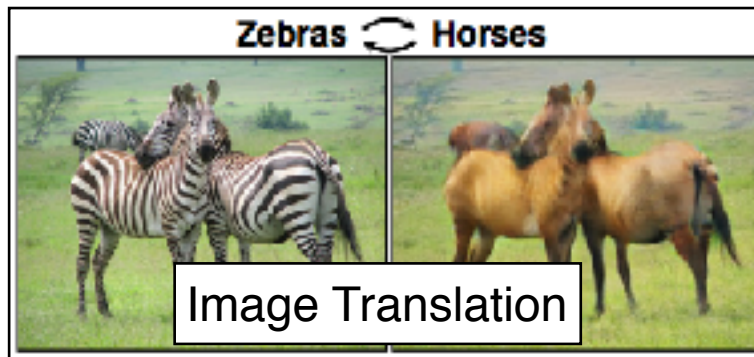
# Generative Adversarial Network

- Generative Adversarial Network (GAN)  
[I. Goodfellow+, NIPS14]

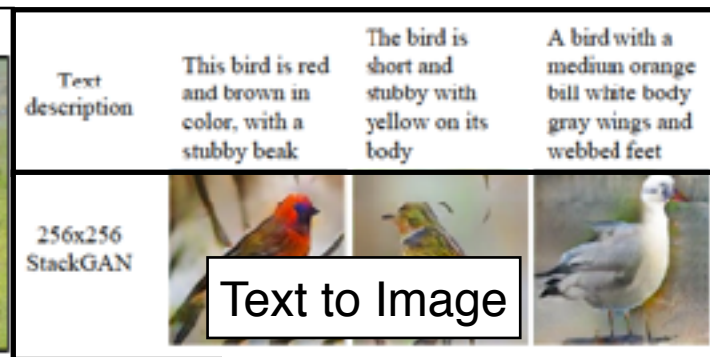


# Generative Adversarial Networks

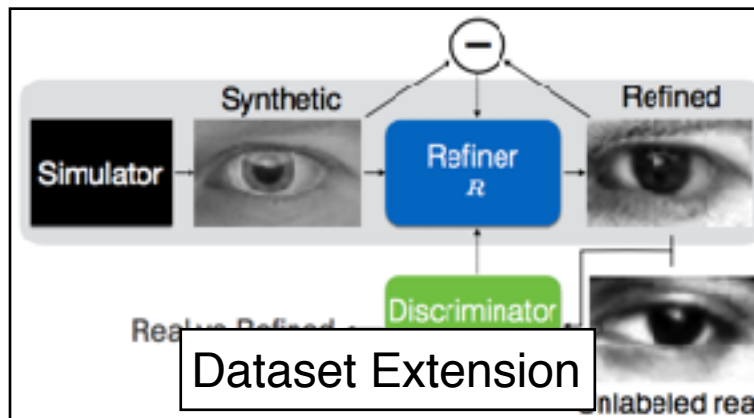
## Application



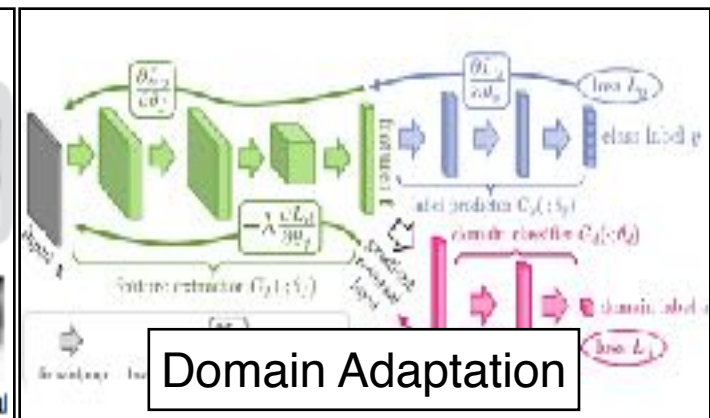
[JY Zhu, et al., ICCV17]



[H. Zhang, et al., ICCV17]



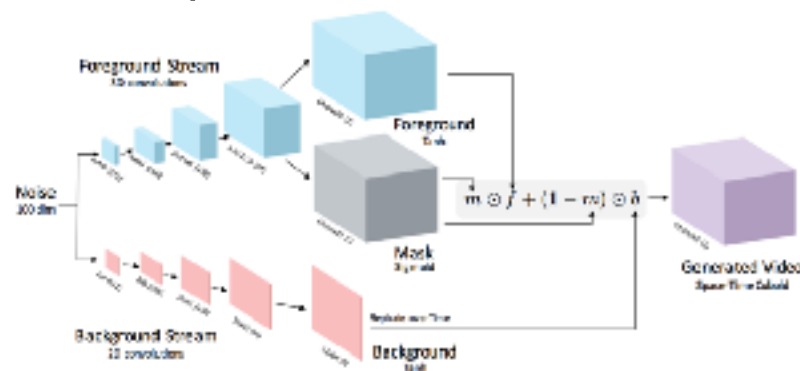
[A. Shrivastava, et al., CVPR17]



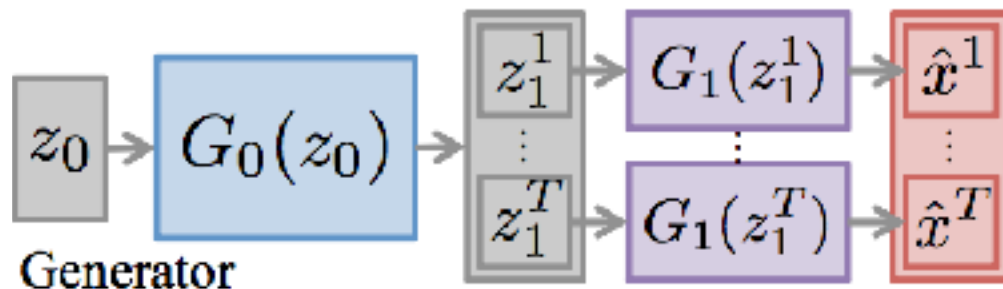
Y. Ganin, et al., ICML15]

# GANs for Video

- Previous works
  - Video GAN (VGAN) [C. Vondrick, et al., NIPS16]



- Temporal GAN (TGAN) [M. Saito, et al., ICCV17]\*



\* We refer their first arXiv version in the paper because ICCV17 papers were not published yet at our submission time.

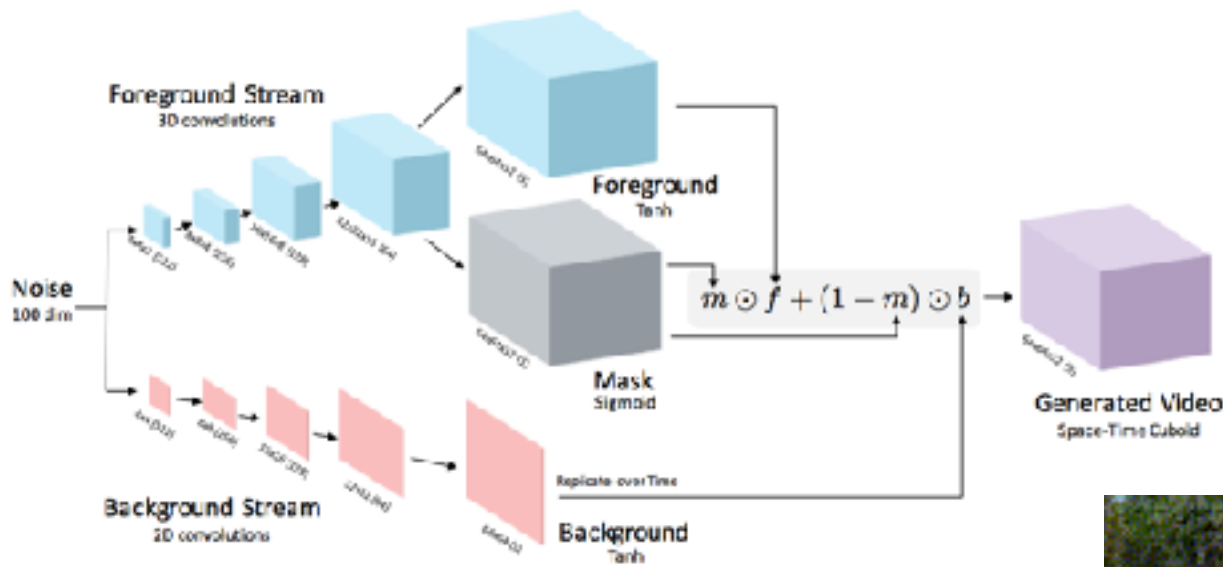
# Challenges in video generation

---

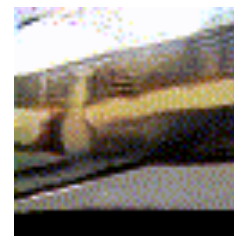
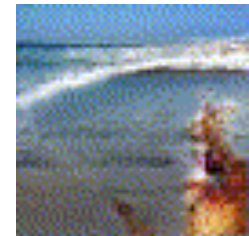
- Important factors for realistic video generation:
  1. Realistic frame
  2. Scene consistency
  3. Reasonable motion

# GANs for Video

- Video GAN (VGAN) [C. Vondrick, et al., NIPS16]

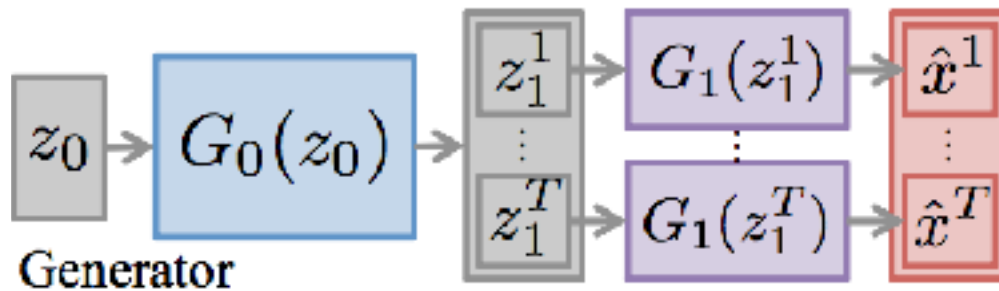


1. ~~Realistic frame~~
2. Scene consistency
3. ~~Reasonable motion~~



# GANs for Video

- Temporal GAN (TGAN) [M. Saito, et al., ICCV17]\*



1. Realistic frame
2. ~~Scene consistency~~
3. ~~Reasonable motion~~



\* We refer their first arXiv version in the paper because ICCV17 papers were not published yet at our submission time.



# Challenges in video generation

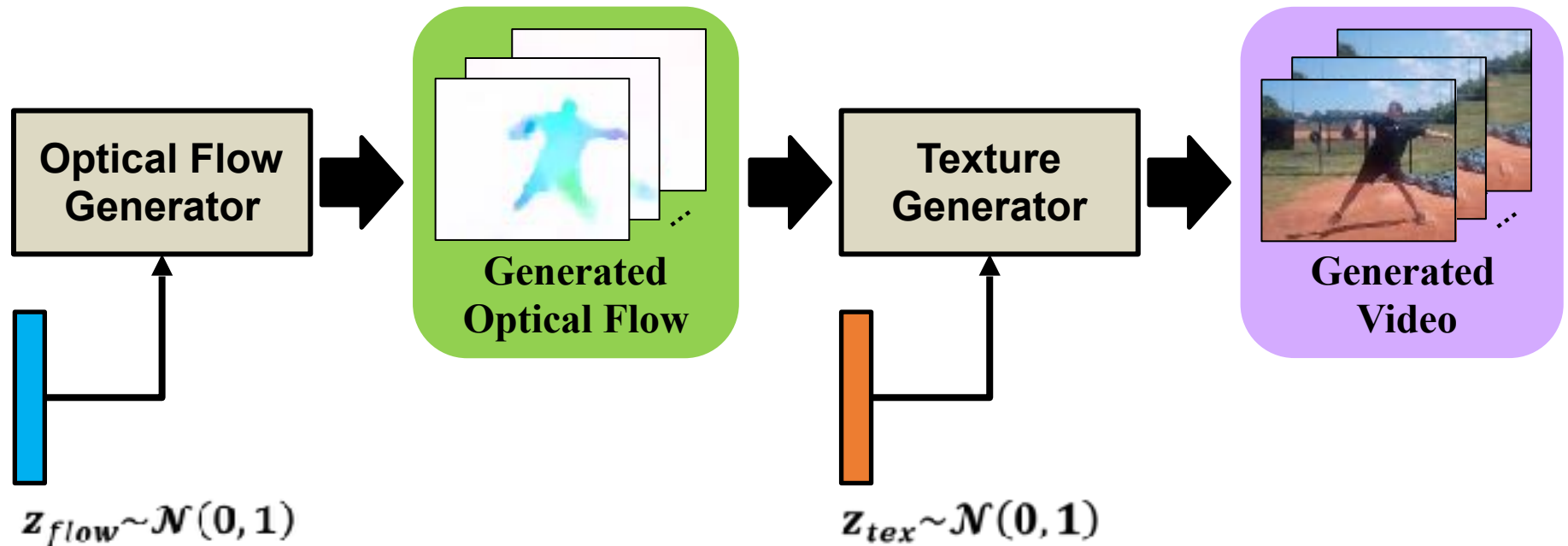
---

- Important factors for realistic video generation:
  1. Realistic frame
  2. Scene consistency
  3. Reasonable motion

It is important to consider structure of video and to make a video generation pipeline that can express the structure.

# Hierarchical Video Generation

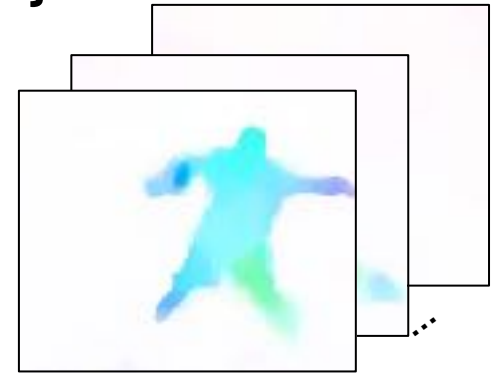
- Generating video via optical flow
  1. Generate optical flow as motion information
  2. Give texture to generated optical flow



# Features of Optical flow

---

- Extractable unsupervisedly
- Holding the contour of a moving object
- Continuity in the time direction
- No texture information

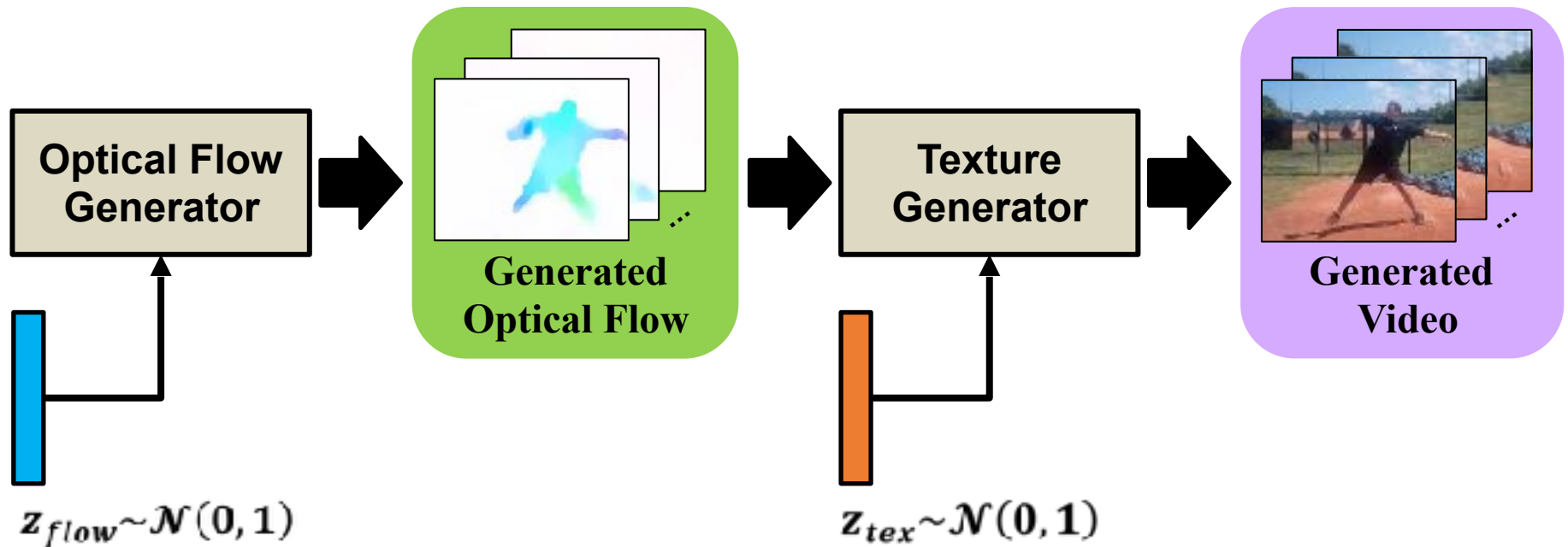


Generating optical flow first makes it ...

- possible to generate a video with reasonable motion  
&
- easier to generate a realistic video than without optical flow

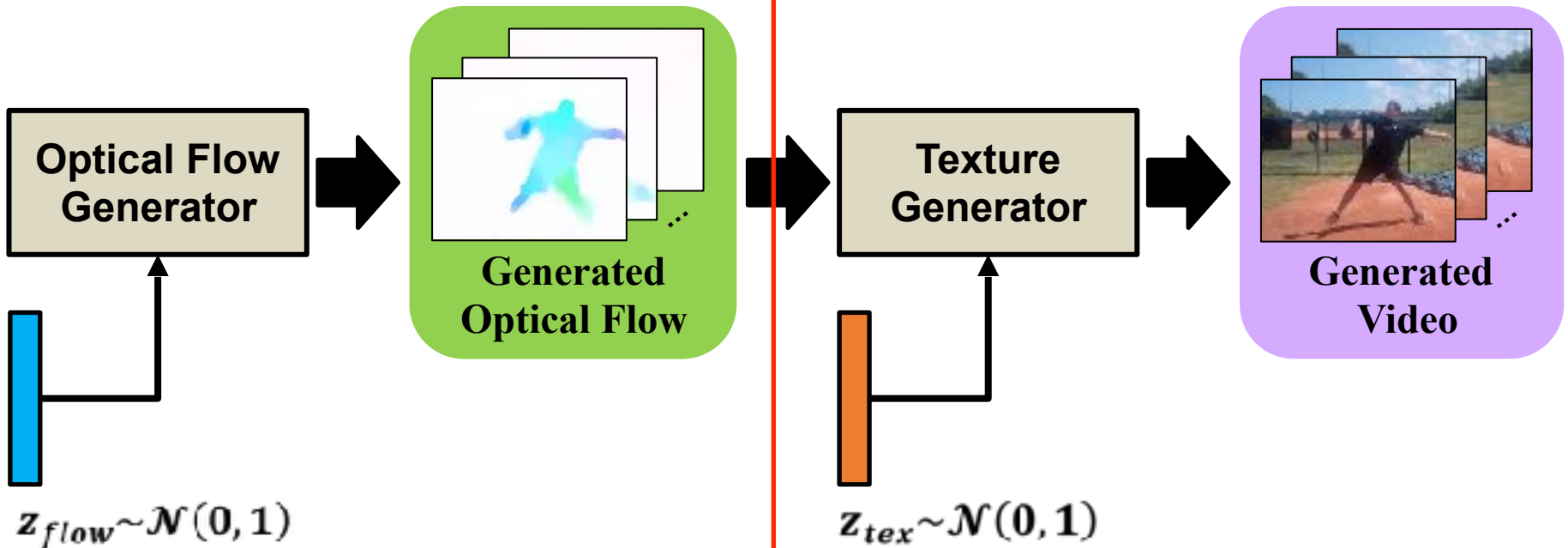
# Proposed Method

- Overview of generator



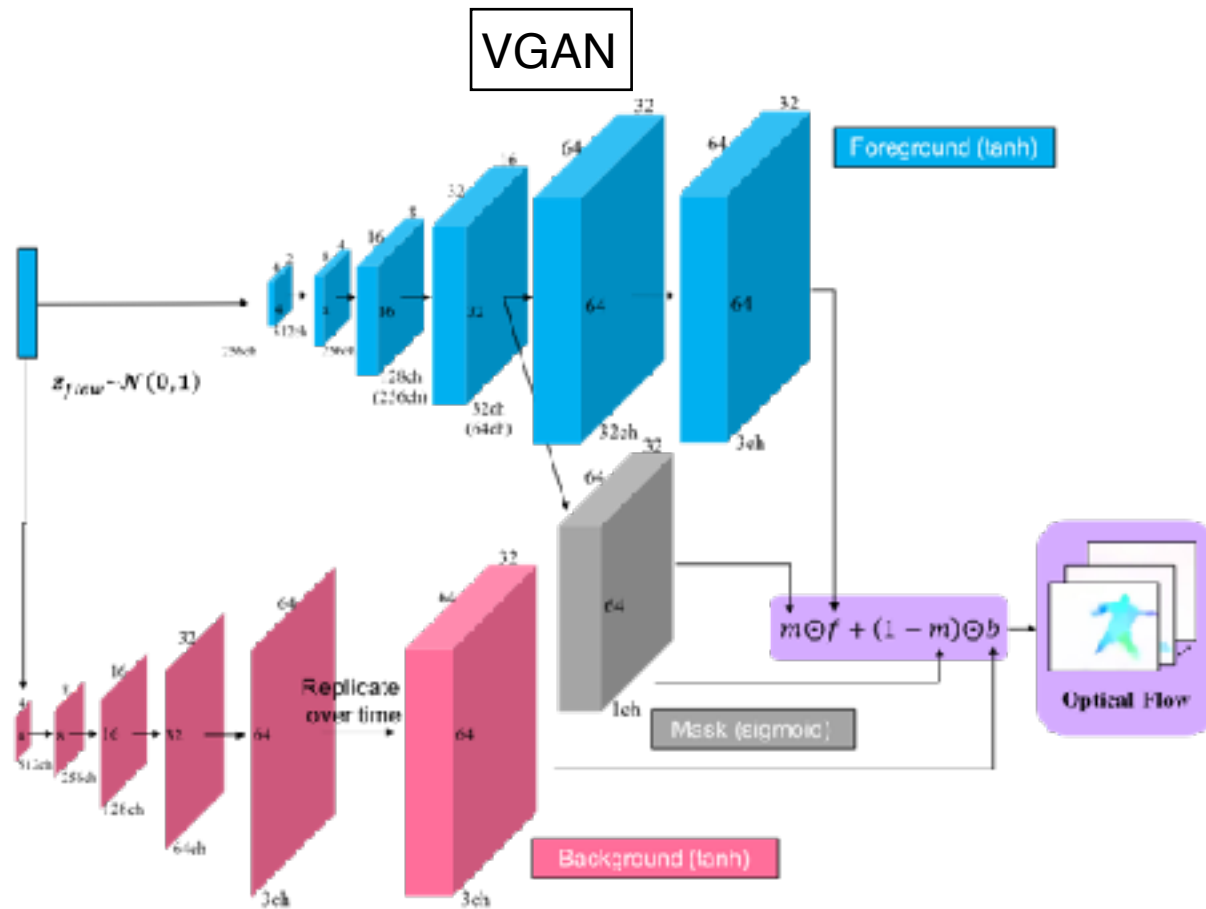
# Proposed Method

- Optical flow generator



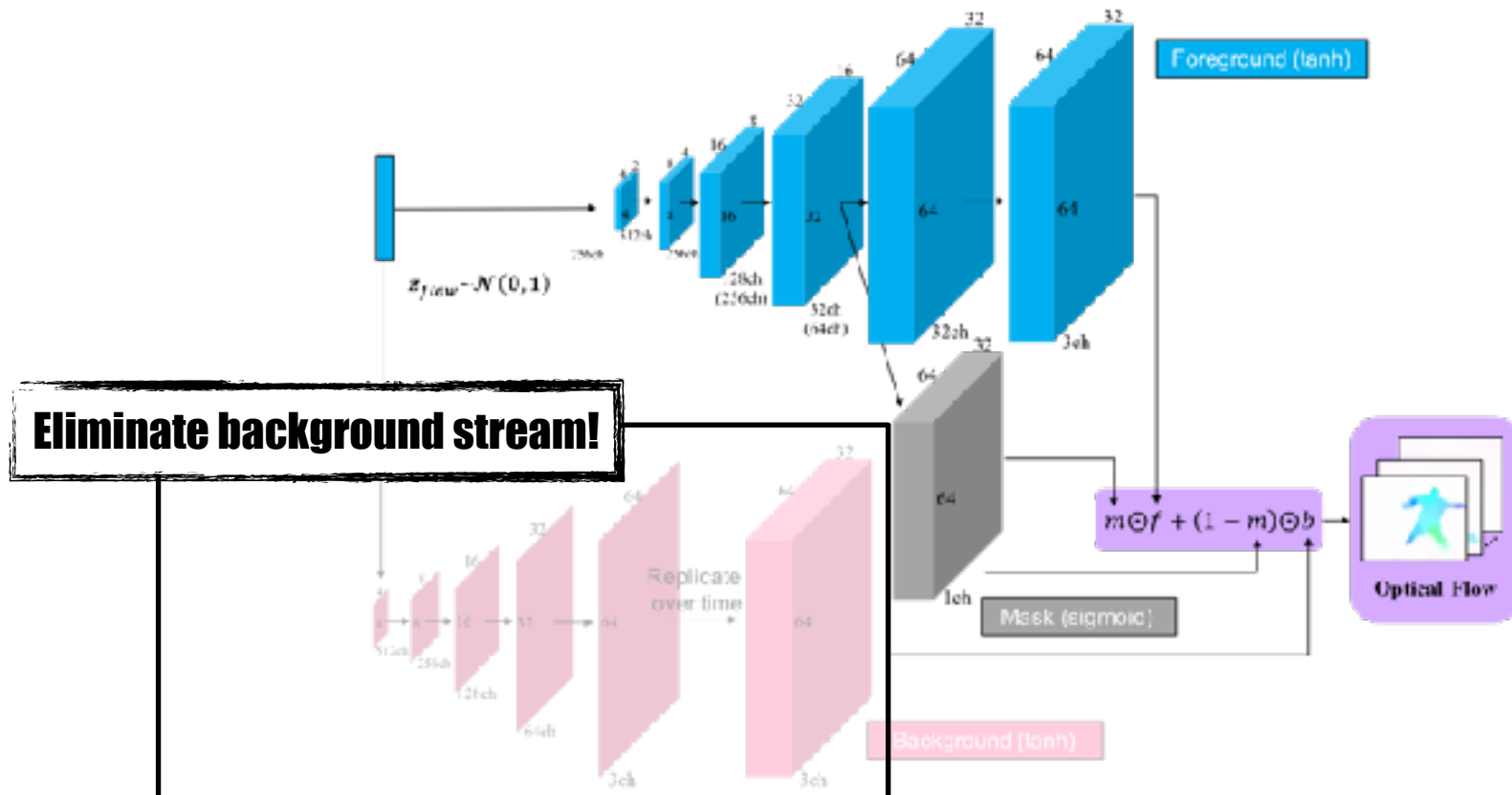
# Optical flow generator

- Optical flow generator is constructed based on the pipeline of VGAN [C. Vondrick+, NIPS16].



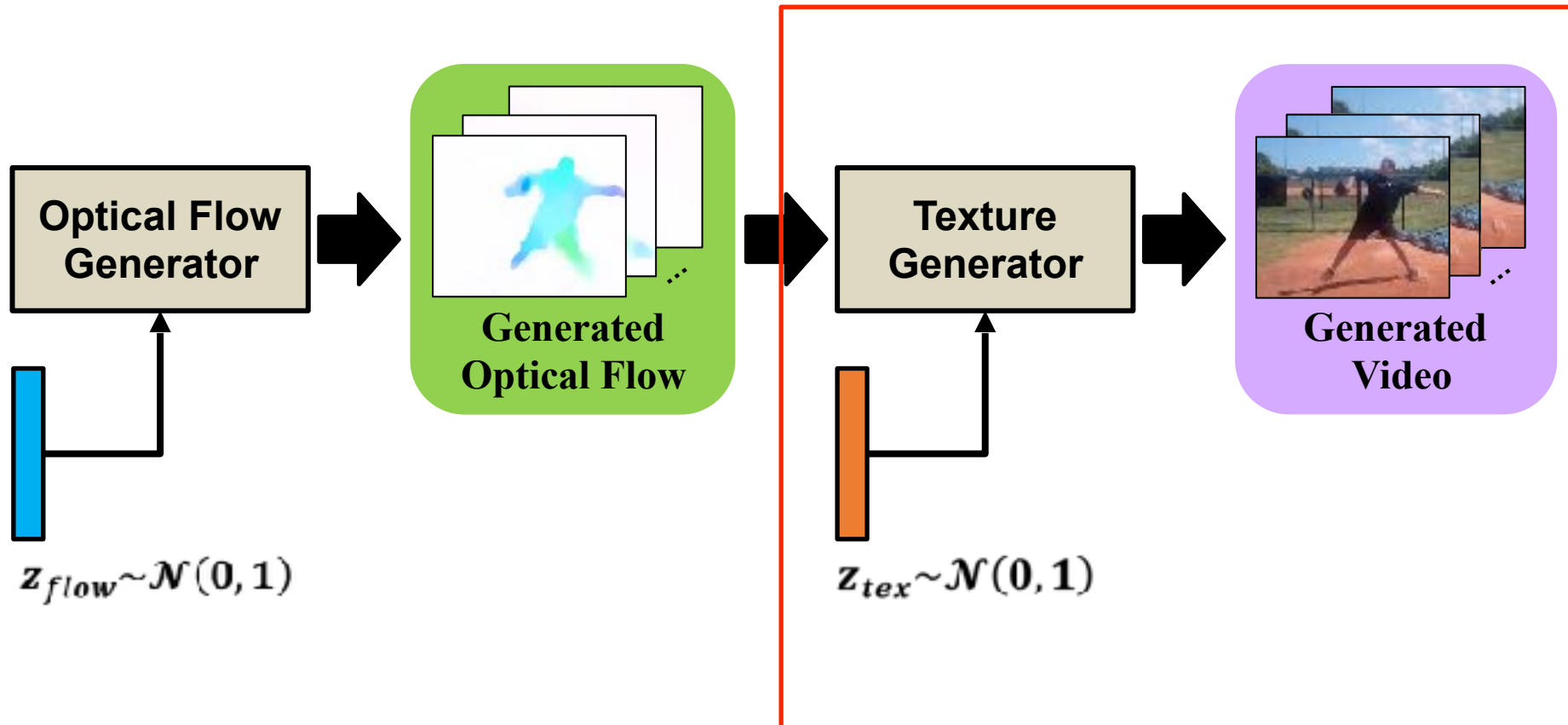
# Optical flow generator

- Background optical flow should be zero
  - If the camera is fixed.



# Proposed Method

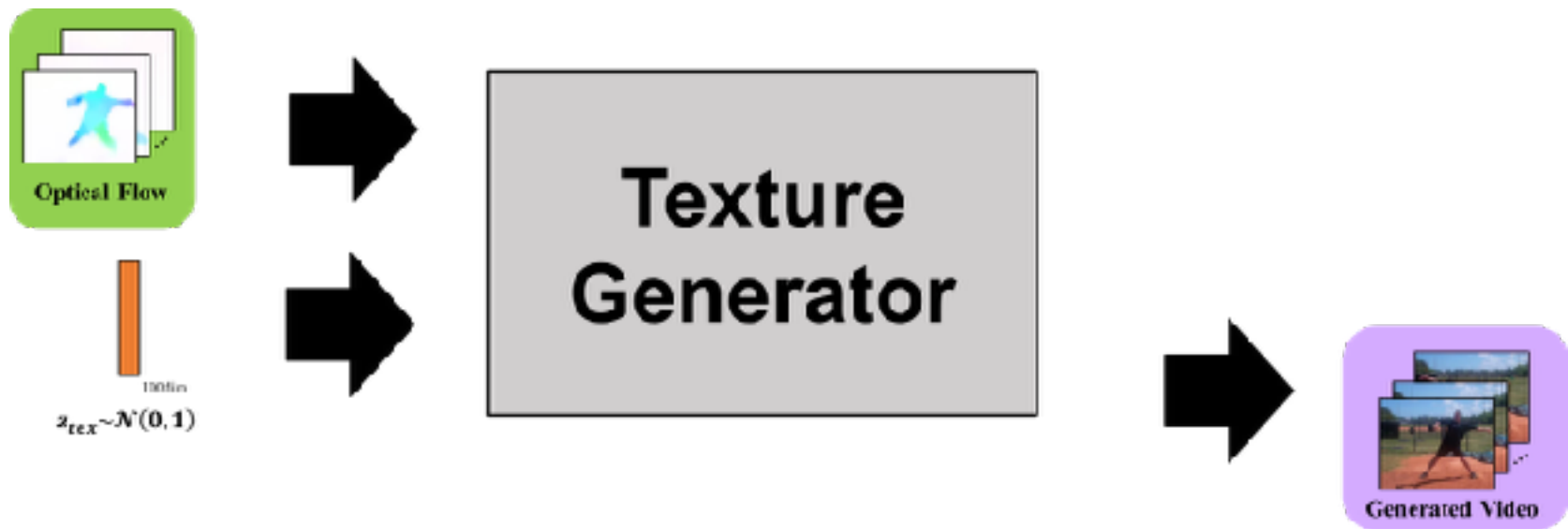
- Texture generator





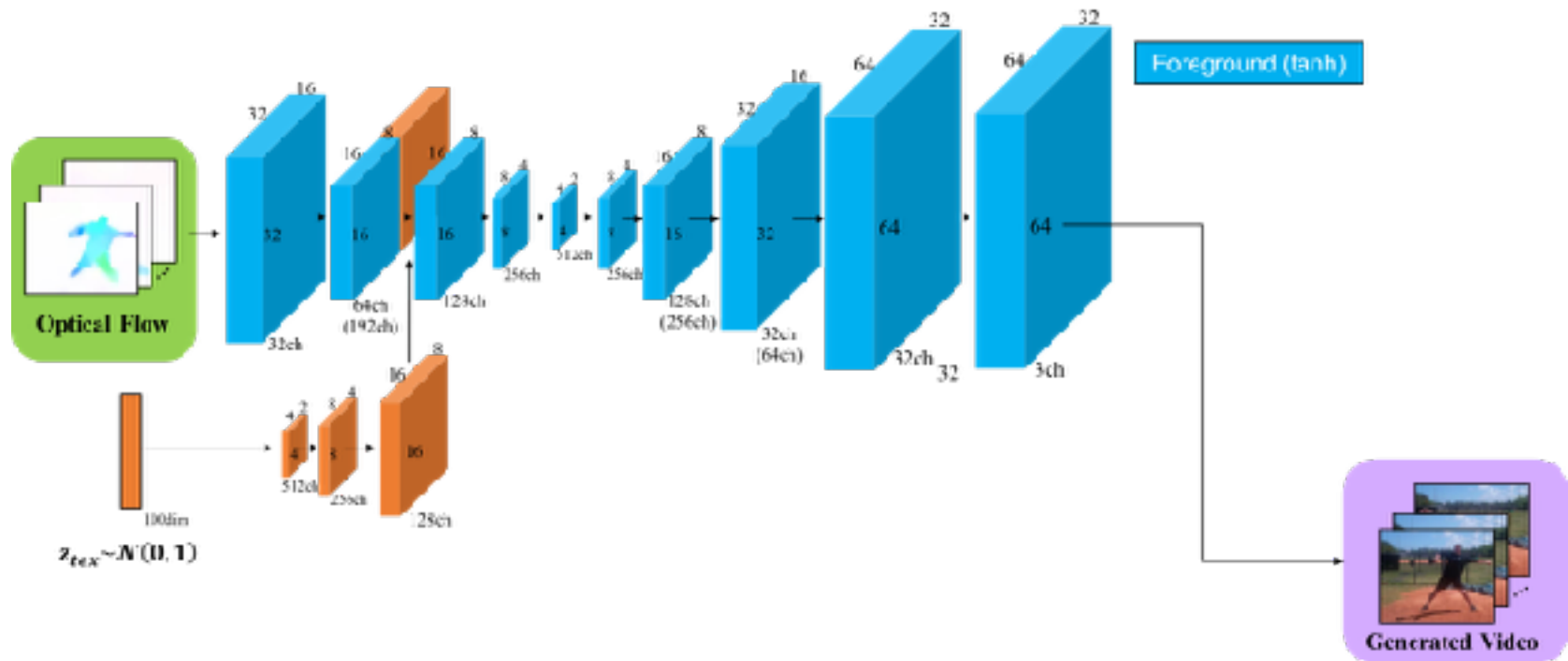
# Texture generator

- Generate RGB video from random noise and optical flow.



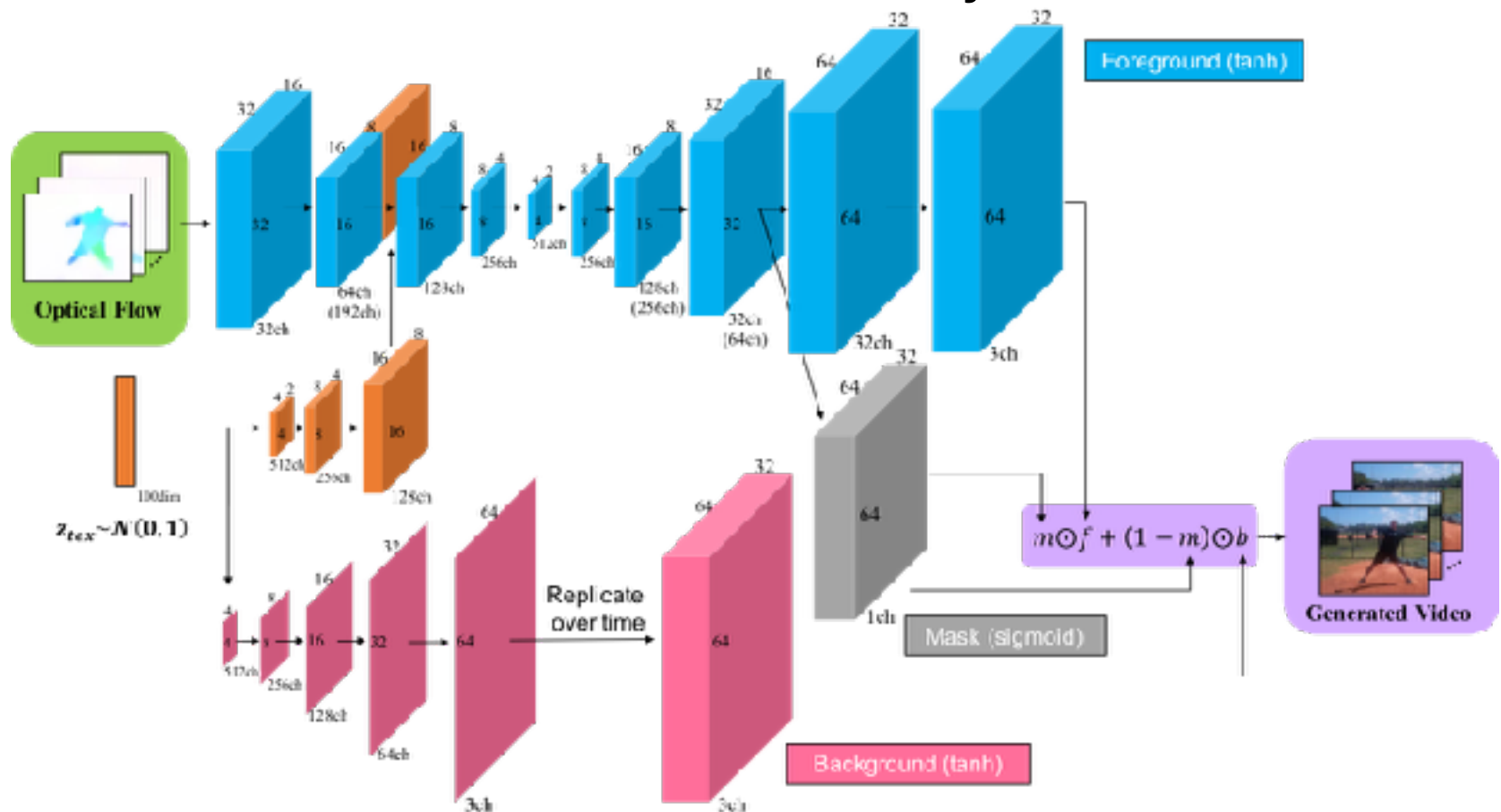
# Texture Generator

- Auto-encoder that converts optical flow to RGB video

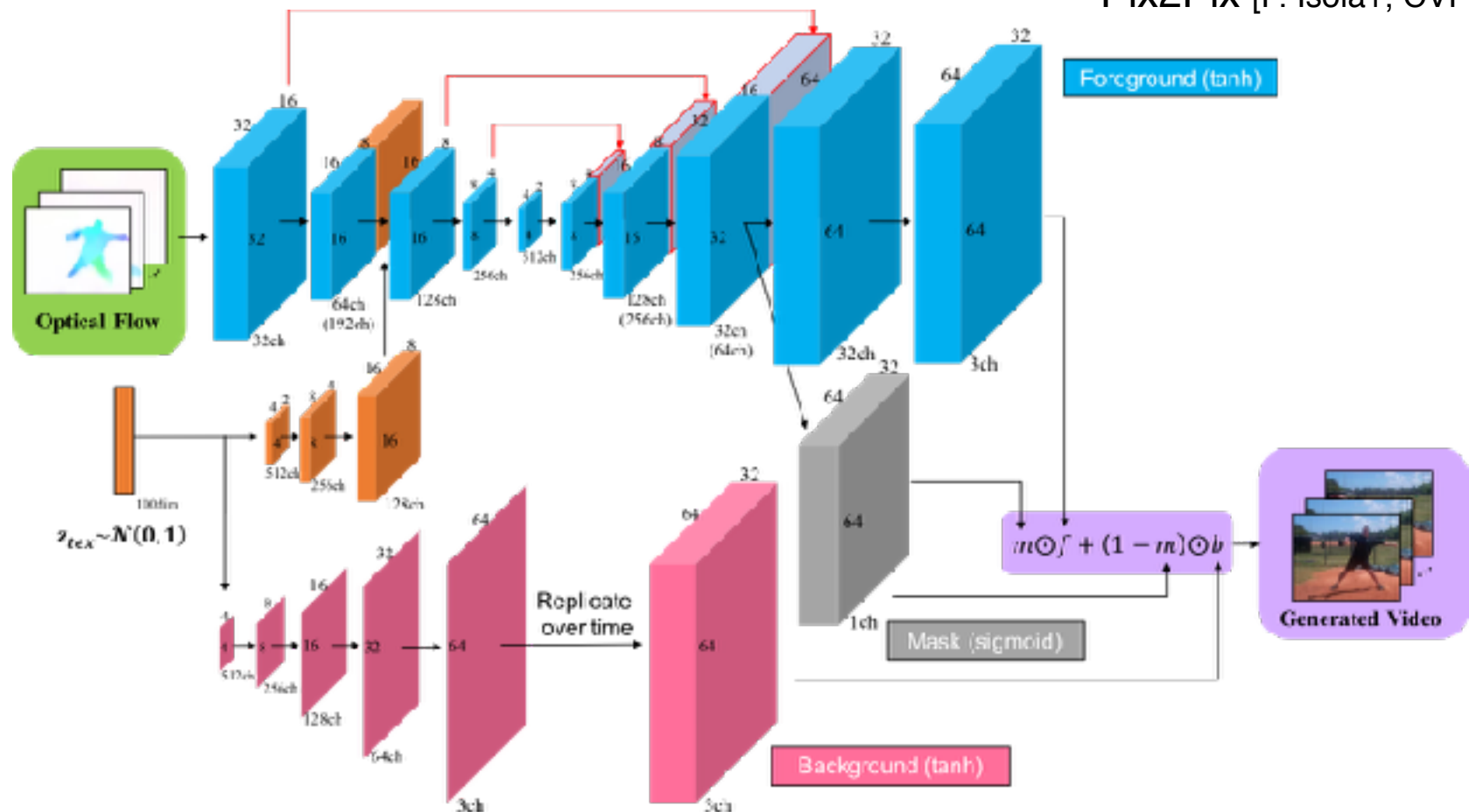


# Texture Generator

- Add background stream as VGAN
  - to obtain scene consistency

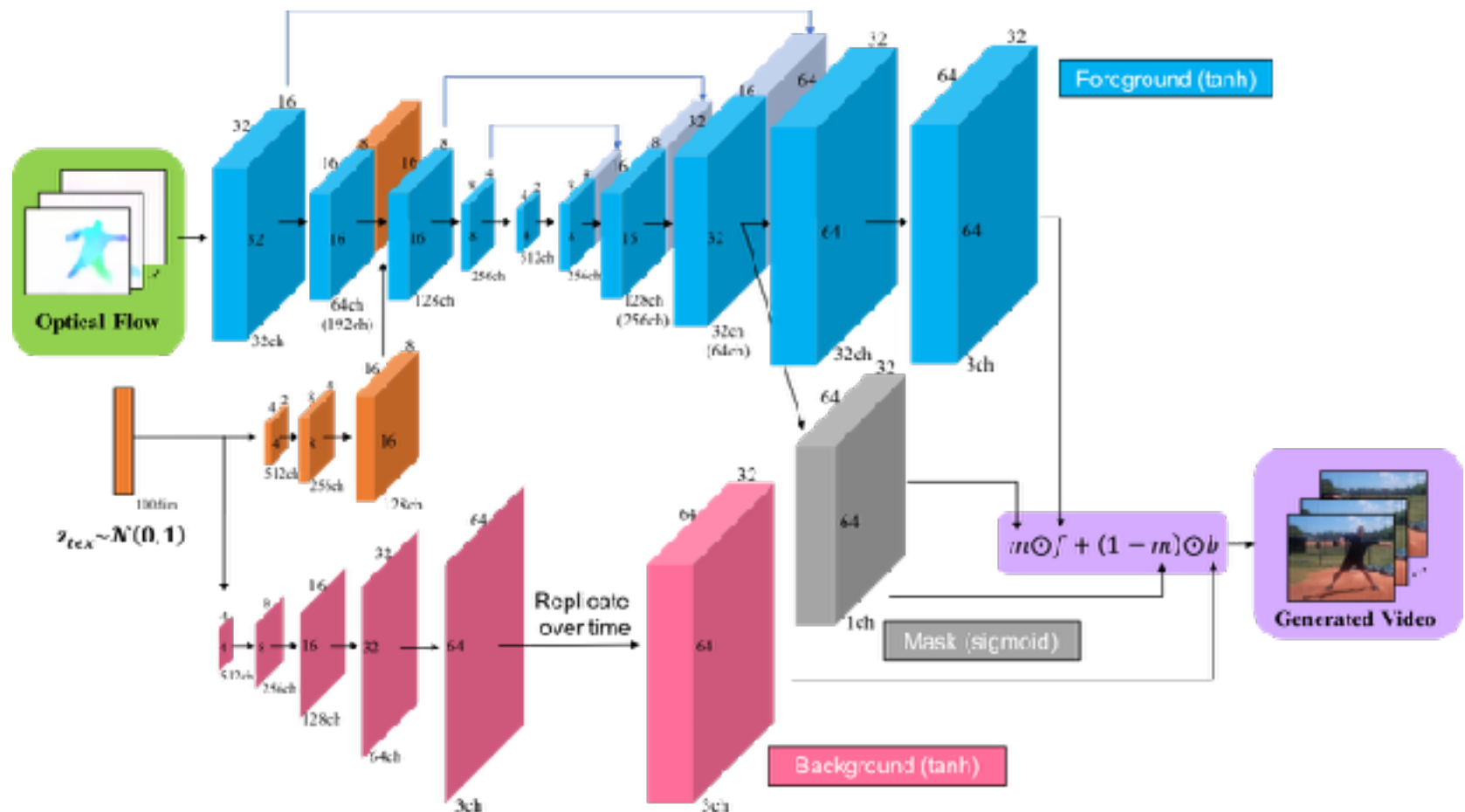


- cf.)  
Pix2Pix [P. Isola+, CVPR17]



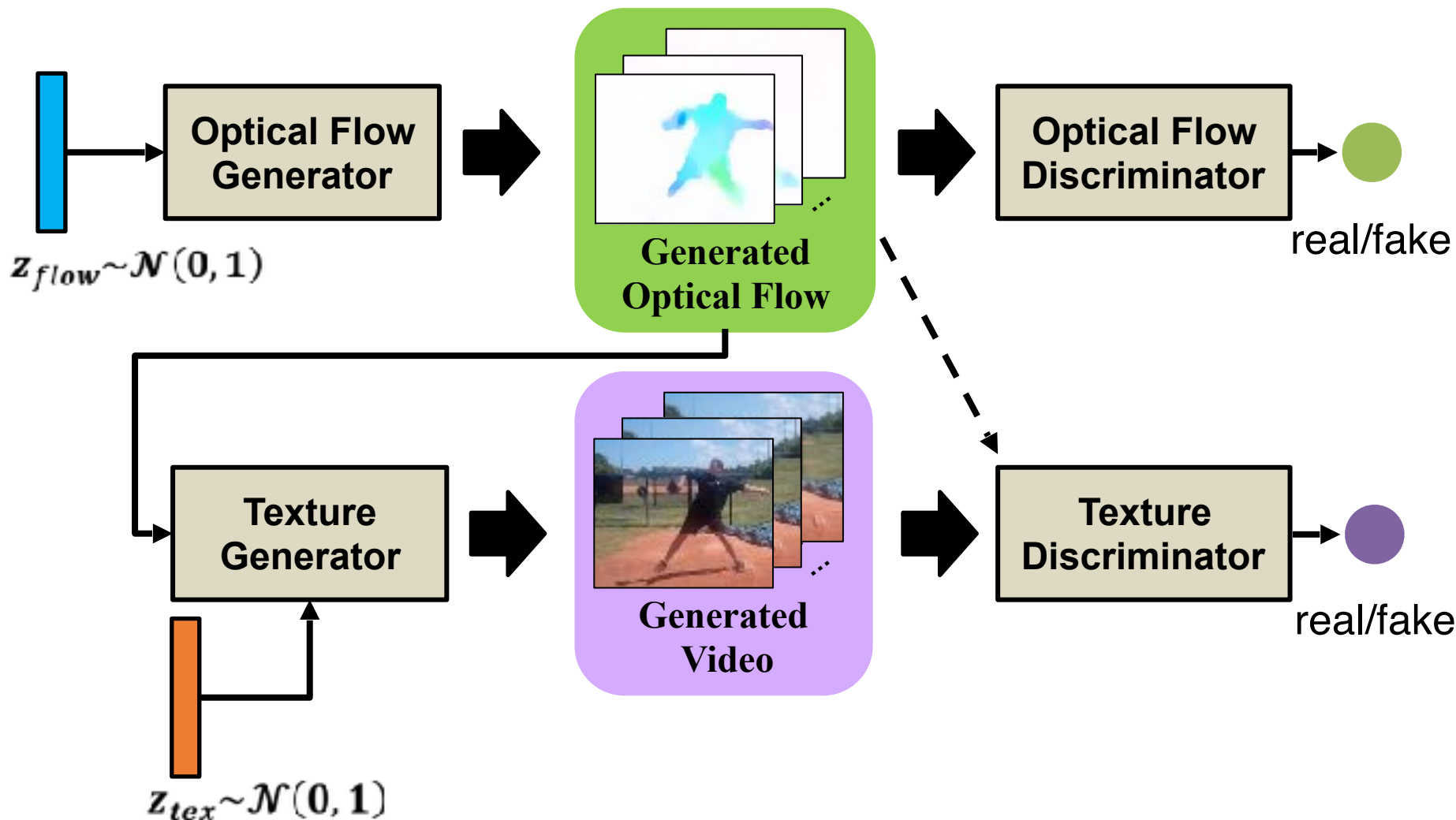
# Texture Generator

- The whole pipeline of our texture generator



# Overview of Proposed Method

- Hierarchical video generation via optical flow



# Experiments

---

- Experiment 1:
  - Examples of generated results
  - Qualitative comparison with baseline
  - Human evaluation
- Experiment 2:
  - Walk in dual  $z$
- Experiment 3:
  - Unsupervised action classification

# Experiments

---

- Experiment 1:
  - Examples of generated results
  - Qualitative comparison with baseline
  - Human evaluation
- Experiment 2:
  - Walk in dual  $z$
- Experiment 3:
  - Unsupervised action classification

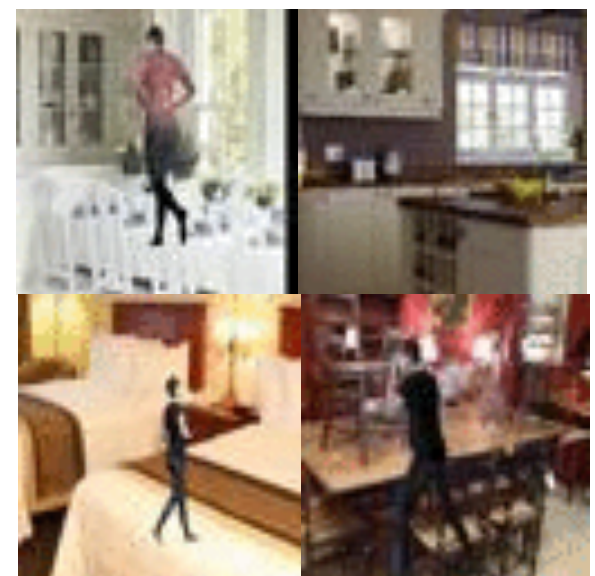
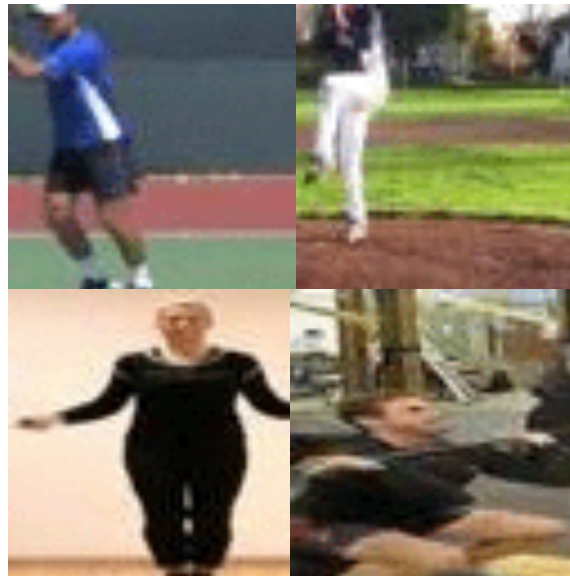


# Experiments

- Dataset
  - Resolution: 64x64
  - Time: 32frames ( $\approx$  1~2 seconds )

Penn Action  
[W. Zhang+, ICCV13] → bounding box Penn Action  
Cropped

SURREAL  
[G. Varol+, CVPR17]

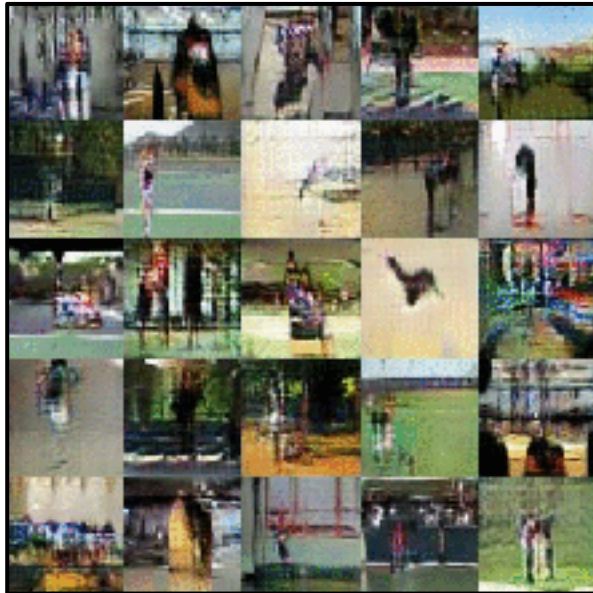


# Experiments

---

- Examples of generated results
  - Various videos

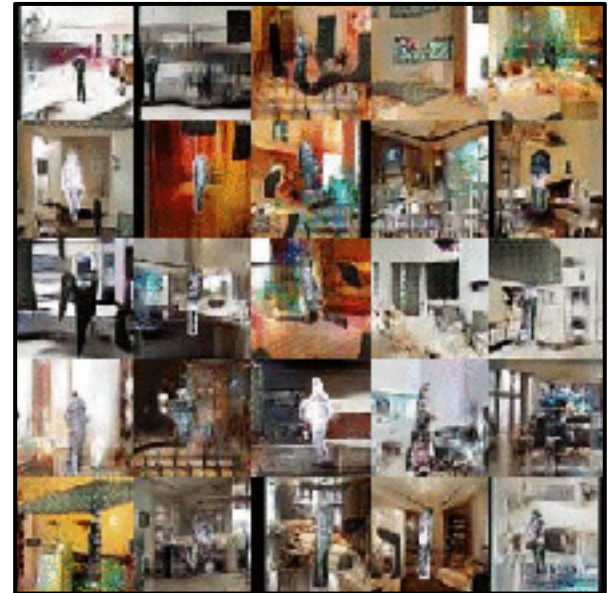
Penn Action



Penn Action  
Cropped



SURREAL



# Experiments

## □ Qualitative comparison with VGAN

Result on  
SURREAL

- FTGAN generates a video with reasonable motion

VGAN



- A person is walking without moving his legs.

FTGAN  
(ours)



- A person is walking by moving their left and right feet in turn.

# Experiments

---

- Qualitative comparison with VGAN
  - FTGAN generates a video with reasonable motion

Result on  
PennAction  
cropped

VGAN



- The outline and the axis of rotation are unclear.

FTGAN  
(ours)



Pull-up

- The outline and the axis of rotation are clear.



# Experiments

---

- AB test on Amazon Mechanical Turk
  - Q: In which video is it easier to figure out what action is being performed?
  - 200 videos
  - 9 votes on each video

Number of videos with better evaluation

	Penn Action	Penn Action Cropped	SURREAL
VGAN	76	91	95
<b>FTGAN (ours)</b>	<b>124</b>	<b>109</b>	<b>105</b>

←  
As the complexity of the dataset increases,  
the proposed method becomes effective

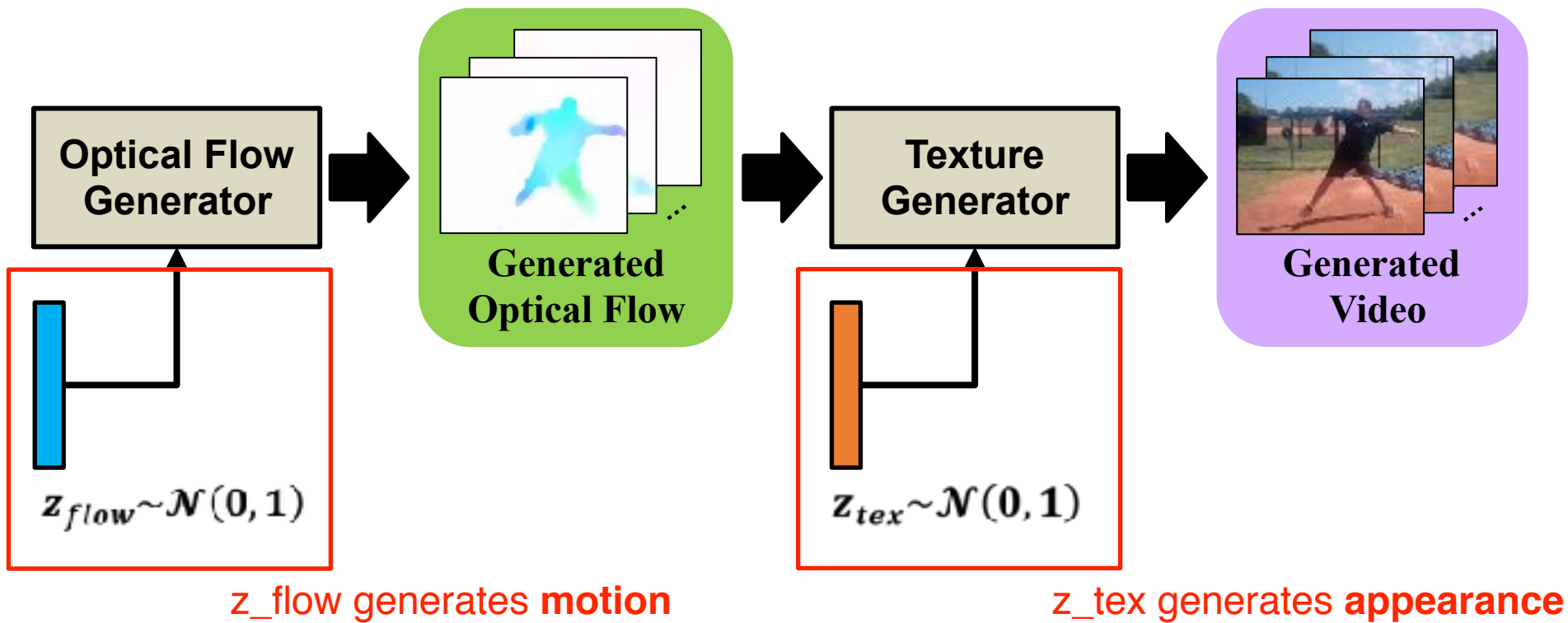
# Experiments

---

- Experiment 1:
  - Examples of generated results
  - Qualitative comparison with baseline
  - Human evaluation
- Experiment 2:
  - Walk in dual  $z$
- Experiment 3:
  - Unsupervised action classification

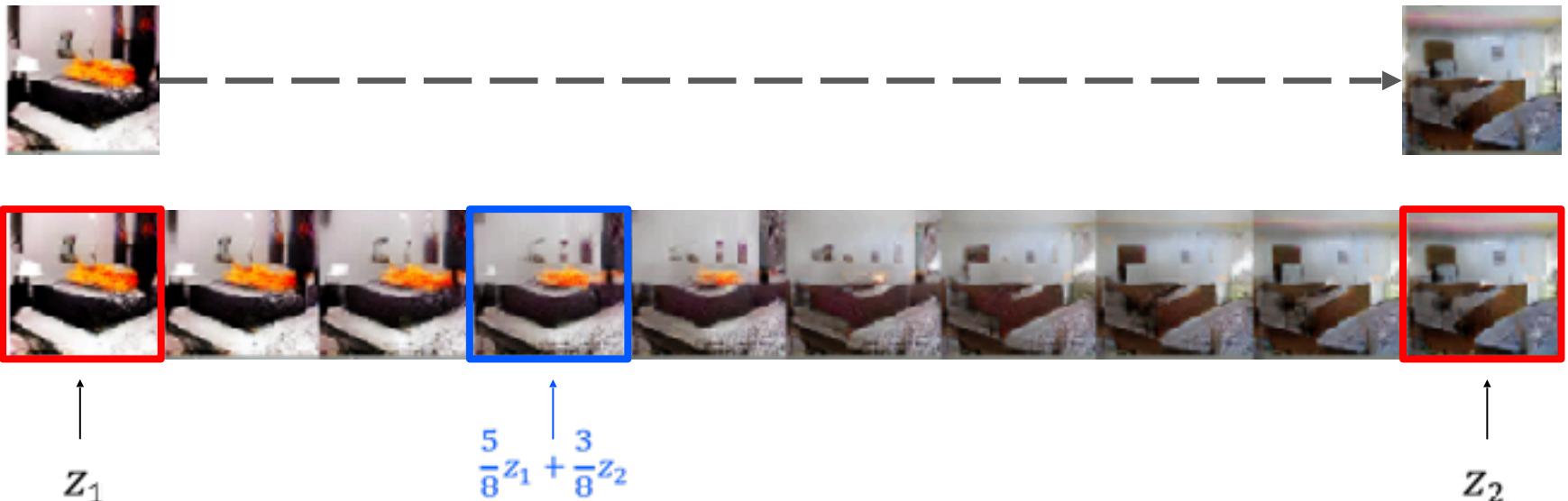
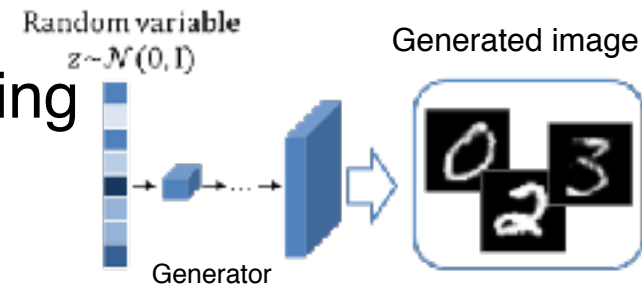
# Experiments

- Overview of generator



# Experiments

- Walk in  $z$ 
  - Conforming mode-collapse avoiding
  - Mixing two images at any ratio



[A. Radford+, ICLR16]



# Experiments

- Dataset
  - SURREAL [G. Varol+, CVPR17] cropped

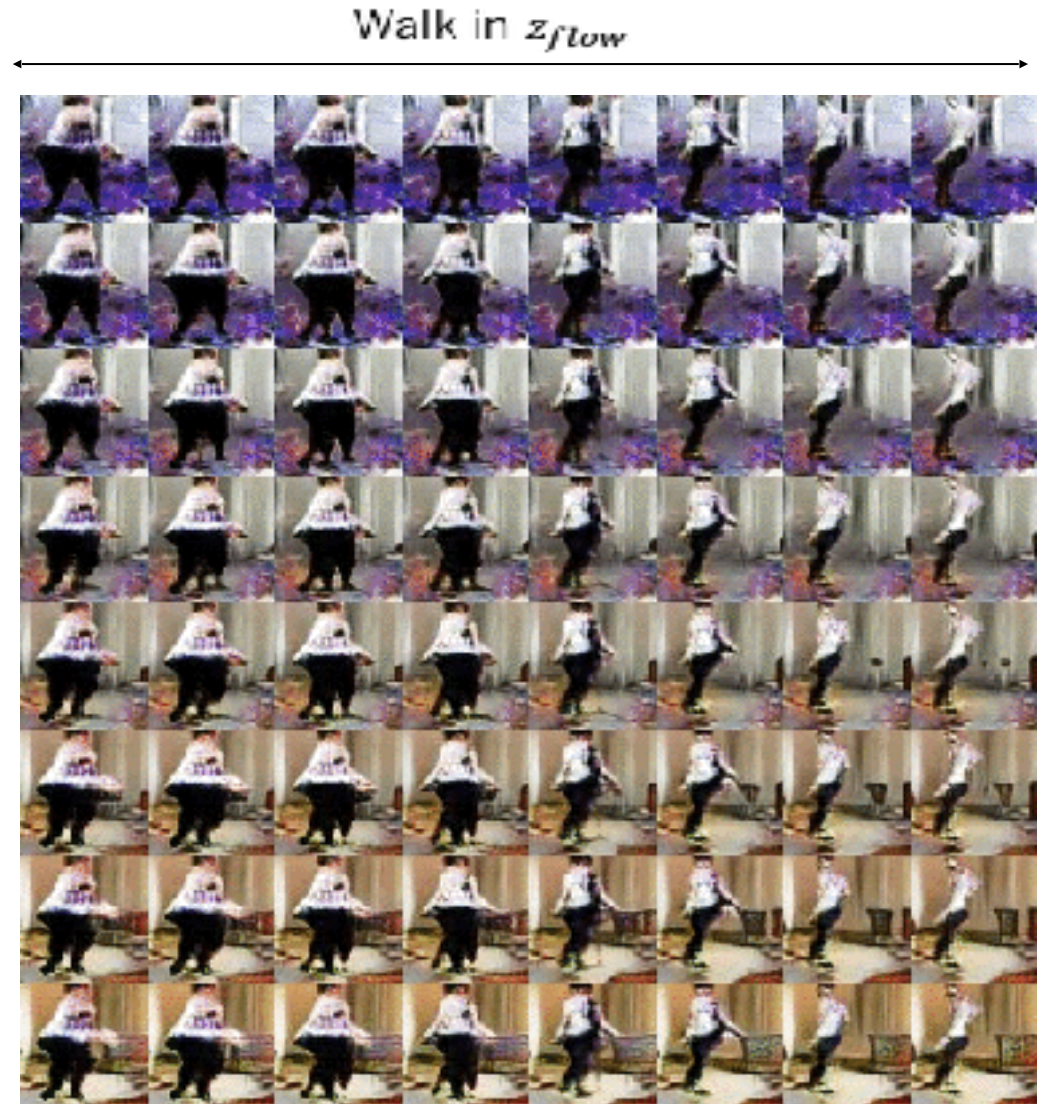


Crop videos with bounding box

# Experiments

- Walk in  $z_{flow}$ 
  - The same motion
  - The appearance changes gradually
- Walk in  $z_{tex}$ 
  - The same appearance
  - The motion changes gradually

Walk in  $z_{tex}$





# Experiments

- Walk in  $z_{flow}$ 
  - The same motion
  - The appearance changes gradually
- Walk in  $z_{tex}$ 
  - The same appearance
  - The motion changes gradually

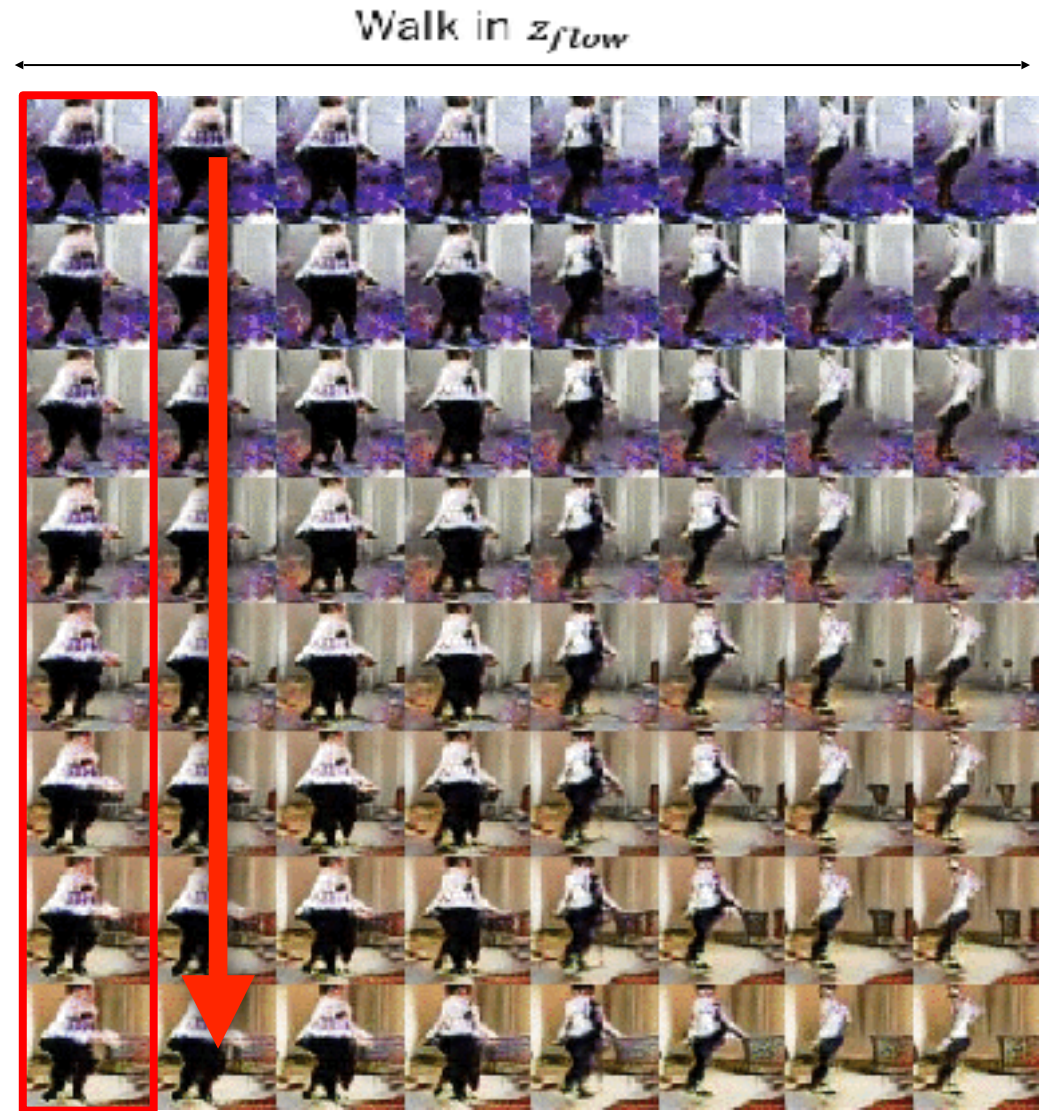
Walk in  $z_{tex}$



# Experiments

- Walk in  $z_{flow}$ 
  - The same motion
  - The appearance changes gradually
- Walk in  $z_{tex}$ 
  - The same appearance
  - The motion changes gradually

Walk in  $z_{tex}$



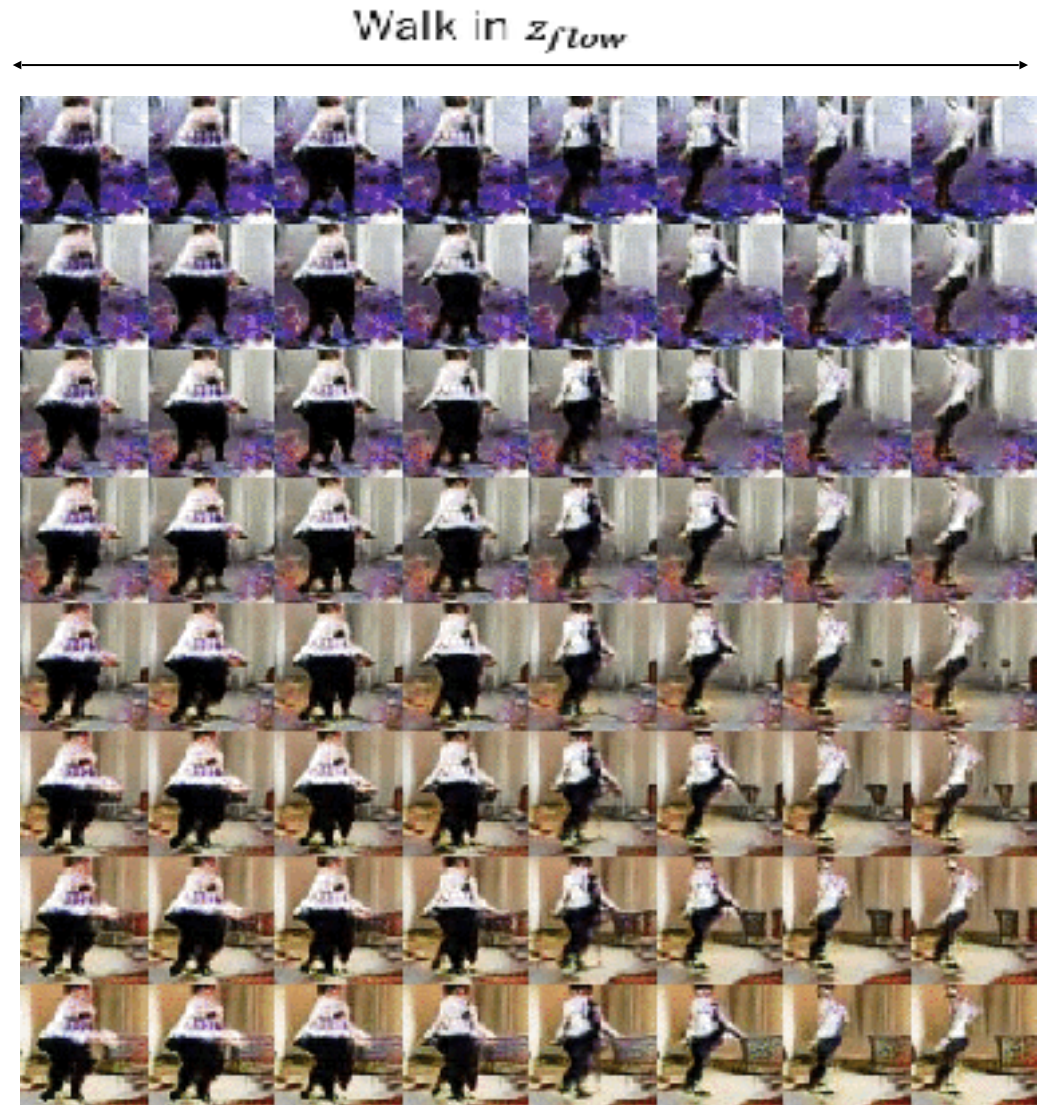


# Experiments

- Walk in  $z_{flow}$ 
  - The same motion
  - The appearance changes gradually
- Walk in  $z_{tex}$ 
  - The same appearance
  - The motion changes gradually

Walk in  $z_{tex}$

Our method can generate videos by independently controlling motion and appearance



# Experiments

---

- Experiment 1:
  - Examples of generated results
  - Qualitative comparison with baseline
  - Human evaluation
- Experiment 2:
  - Walk in dual  $z$
- Experiment 3:
  - Unsupervised action classification

# Experiments

---

## Purpose

- Investigate the unsupervised feature expression learning capability as the same way with previous works

## Method

- Extract the last layer in discriminator as feature

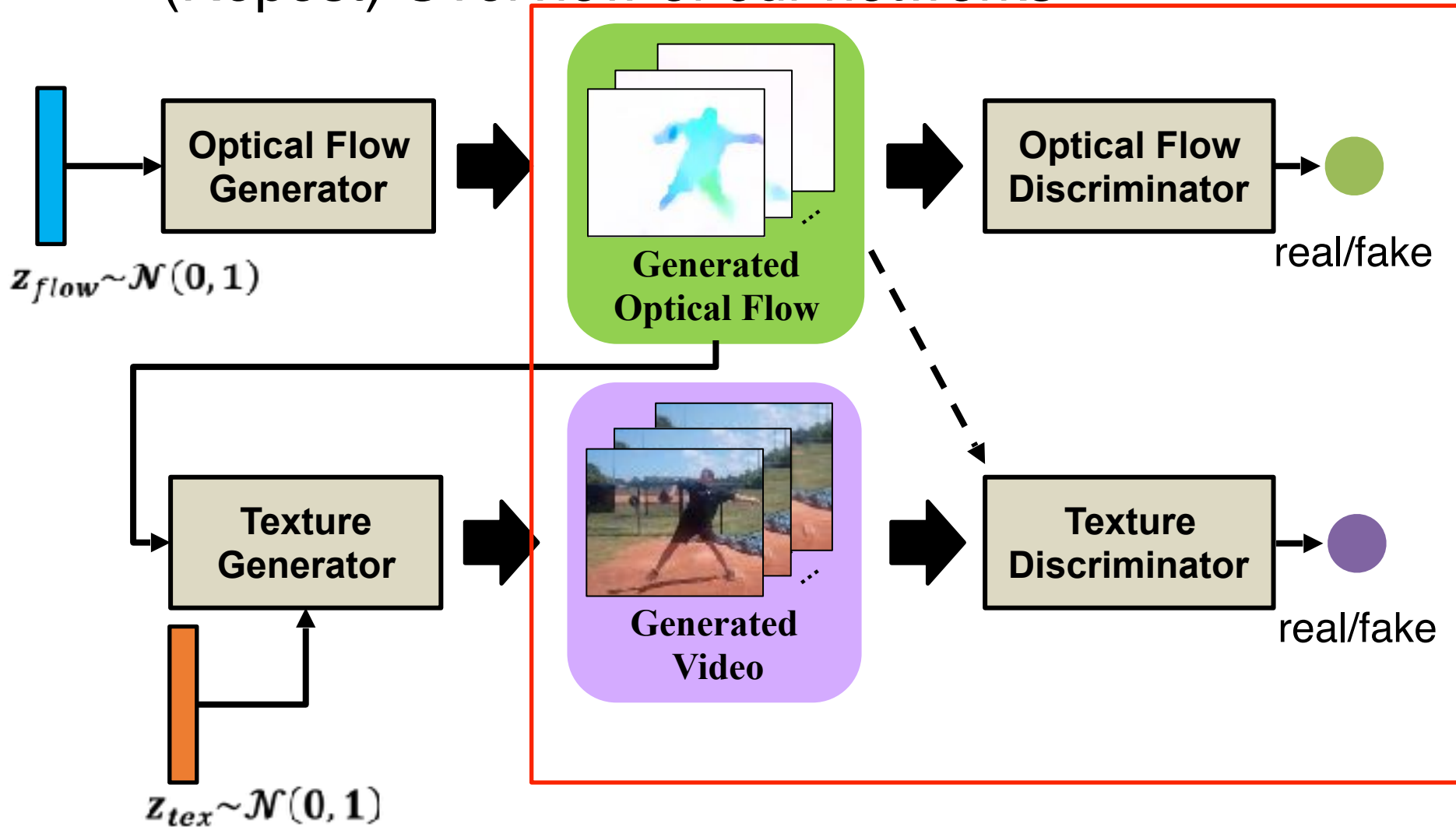
## Setting

- Dataset: UCF101 [K. Soomro, et al., arXiv 2012]
  - 101 classes
  - 13320 videos
- Classifier: SVM

All of these settings are following previous works (VGAN and TGAN).

# Experiments

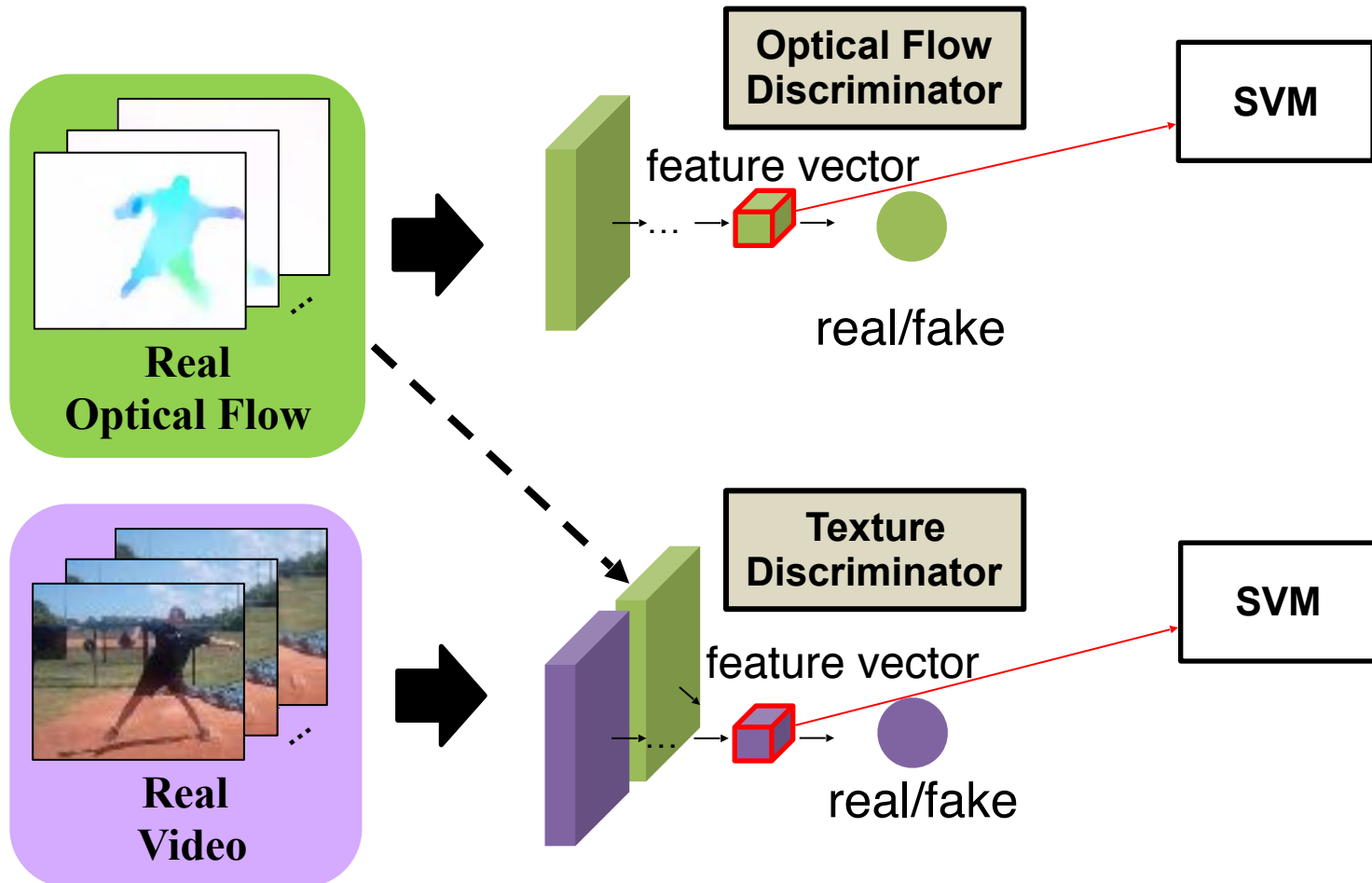
- (Repost) Overview of our networks.





# Experiments

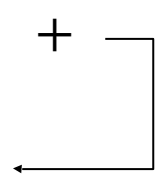
- Extract the last layer of discriminator as feature vector



# Experiments

- Late fusion of flow-discriminator and texture-discriminator improves recognition accuracy.
  - Features learned by each network is complementary, which means...
    - Flow-discriminator learns **motion** information
    - Texture-discriminator learns **appearance** information

Method	Accuracy
Chance	0.9%
(a) Flow-discriminator + Linear SVM ( <b>ours</b> )	<b>48.0%</b>
(b) Texture-discriminator + Linear SVM ( <b>ours</b> )	<b>50.3%</b>
(a) + (b) FTGAN (fusion by Linear SVM) ( <b>ours</b> )	<b>59.7%</b>



A diagram on the right side of the table shows a plus sign (+) with a line connecting it to the '59.7%' accuracy value in the final row, indicating that the final result is a combination of the two methods above it.

# Experiments

- FTGAN outperforms VGAN and TGAN
  - Separating information ensures the capture of much richer video characteristics

Method	Accuracy
Chance	0.9%
VGAN + Random Init [C. Vondrick+, NIPS16]	36.7%
TGAN: Image-discriminator + Linear SVM [M. Saito et al, arXiv]	38.6%
TGAN: Temporal-discriminator + Linear SVM [M. Saito et al, arXiv]	23.3%
(a) Flow-discriminator + Linear SVM ( <b>ours</b> )	48.0%
(b) Texture-discriminator + Linear SVM ( <b>ours</b> )	50.3%
(a) + (b) FTGAN (fusion by Linear SVM) ( <b>ours</b> )	59.7%

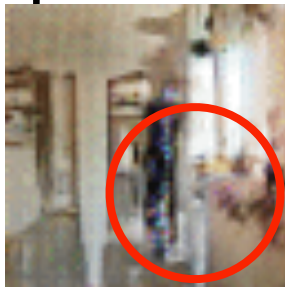
outperform!

# Summary

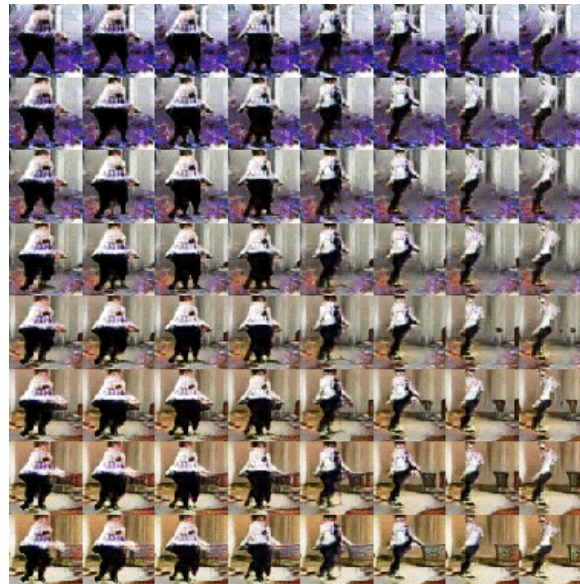
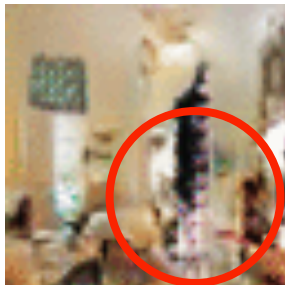
*Thank you!*

- We propose a hierarchical video generative model via optical flow: FTGAN.
- Experiments:

VGAN



FTGAN



UCF101

Method	Accuracy
VGAN	36.7%
TGAN	38.6%
<b>FTGAN</b>	<b>59.7%</b>

It is important to consider structure of video  
and to make a video generation pipeline  
that can express the structure.

---

- Fin

- 終

# Experiments

- Does  $z_{flow}$  generates motion and  $z_{tex}$  generates appearance independently?
  - vertical: generated from the same  $z_{flow}$
  - Horizontal: generated from the same  $z_{tex}$

Penn Action

$z_{flow}^1$        $z_{flow}^2$



Penn Action  
Cropped

$z_{flow}^1$        $z_{flow}^2$



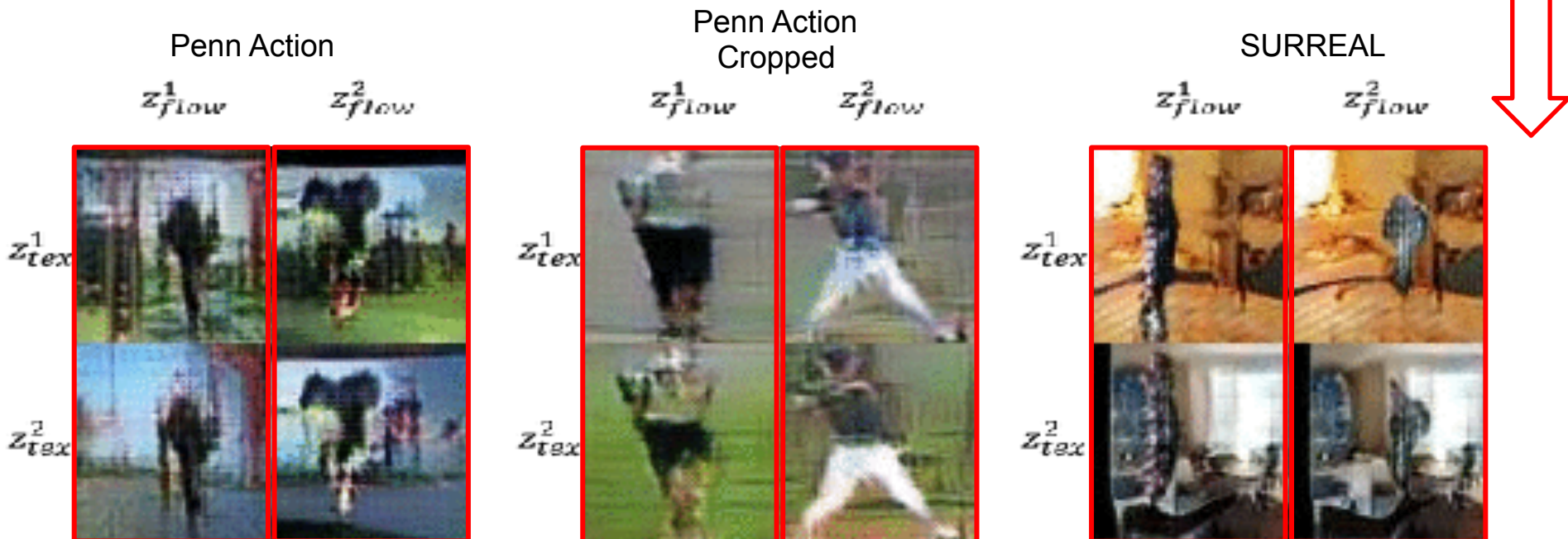
SURREAL

$z_{flow}^1$        $z_{flow}^2$



# Experiments

- Does  $z_{\text{flow}}$  generate motion and  $z_{\text{tex}}$  generate appearance independently?
  - vertical: generated from the same  $z_{\text{flow}}$
  - Horizontal: generated from the same  $z_{\text{tex}}$

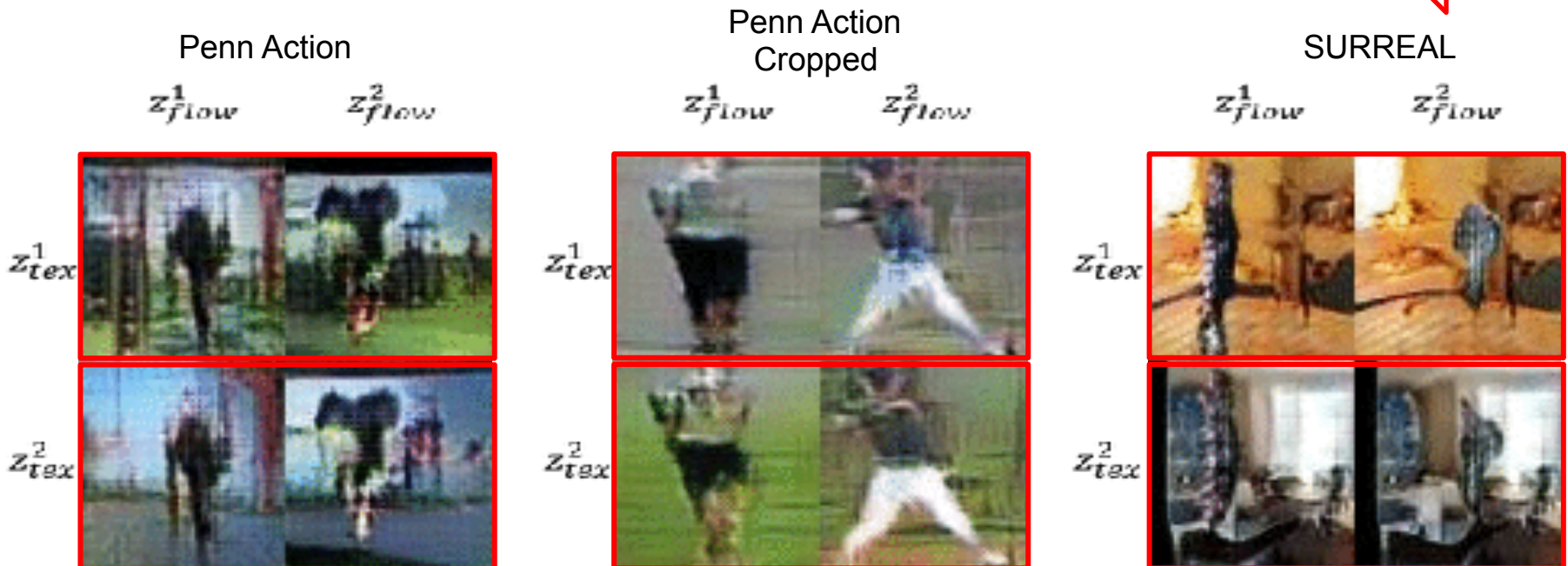


The same movements and different appearance



# Experiments

- Does  $z_{flow}$  generate motion and  $z_{tex}$  generate appearance independently?
  - vertical: generated from the same  $z_{flow}$
  - Horizontal: generated from the same  $z_{tex}$



Different movements and the same appearance







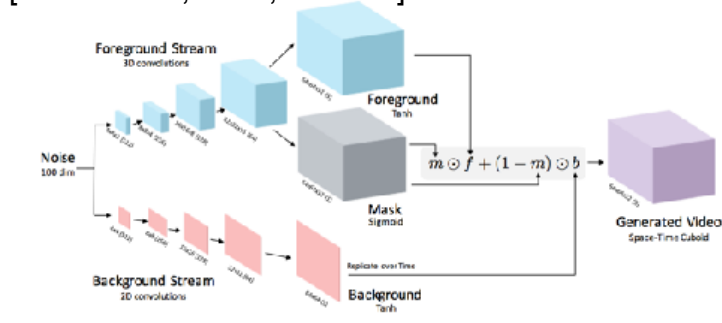
# 動画生成の難しさ

- 生成された動画が本物らしくあるためには、以下の3つの条件を満たしている必要がある
  - a. 各フレームがきれいな画像になっている
  - b. 動画内でのシーンの一貫性が保たれている
  - c. 動きが妥当なものになっている

## - GANを動画生成に拡張した手法

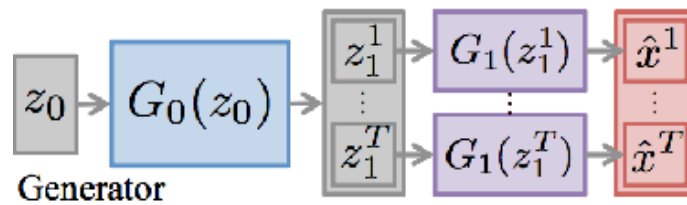
### Video GAN (VGAN)

[C. Vondrick, et al., NIPS16]



### Temporal GAN (TGAN)

[M. Saito et al, arxiv]

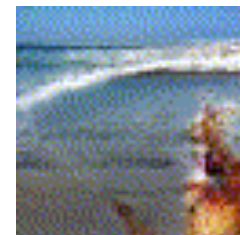
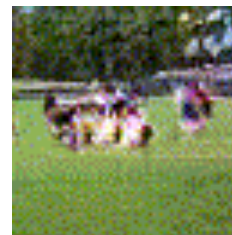


# 関連研究 -動画生成-

- 生成された動画が本物らしくあるためには、以下の3つの条件を満たしている必要がある
  - a. 各フレームがきれいな画像になっている
  - b. 動画内でのシーンの一貫性が保たれている
  - c. 動きが妥当なものになっている

手法: **VGAN** [C. Vondrick, et al., NIPS16]

- 動くものを前景、動かないものを背景として生成
  - 動画内で同じシーンが現れる
- 3D convolutionを使用
  - 見た目と動きを同時に学習



# 関連研究 -動画生成-

- 生成された動画が本物らしくあるためには、以下の3つの条件を満たしている必要がある
  - a. 各フレームがきれいな画像になっている
  - b. 動画内でのシーンの一貫性が保たれている
  - c. 動きが妥当なものになっている

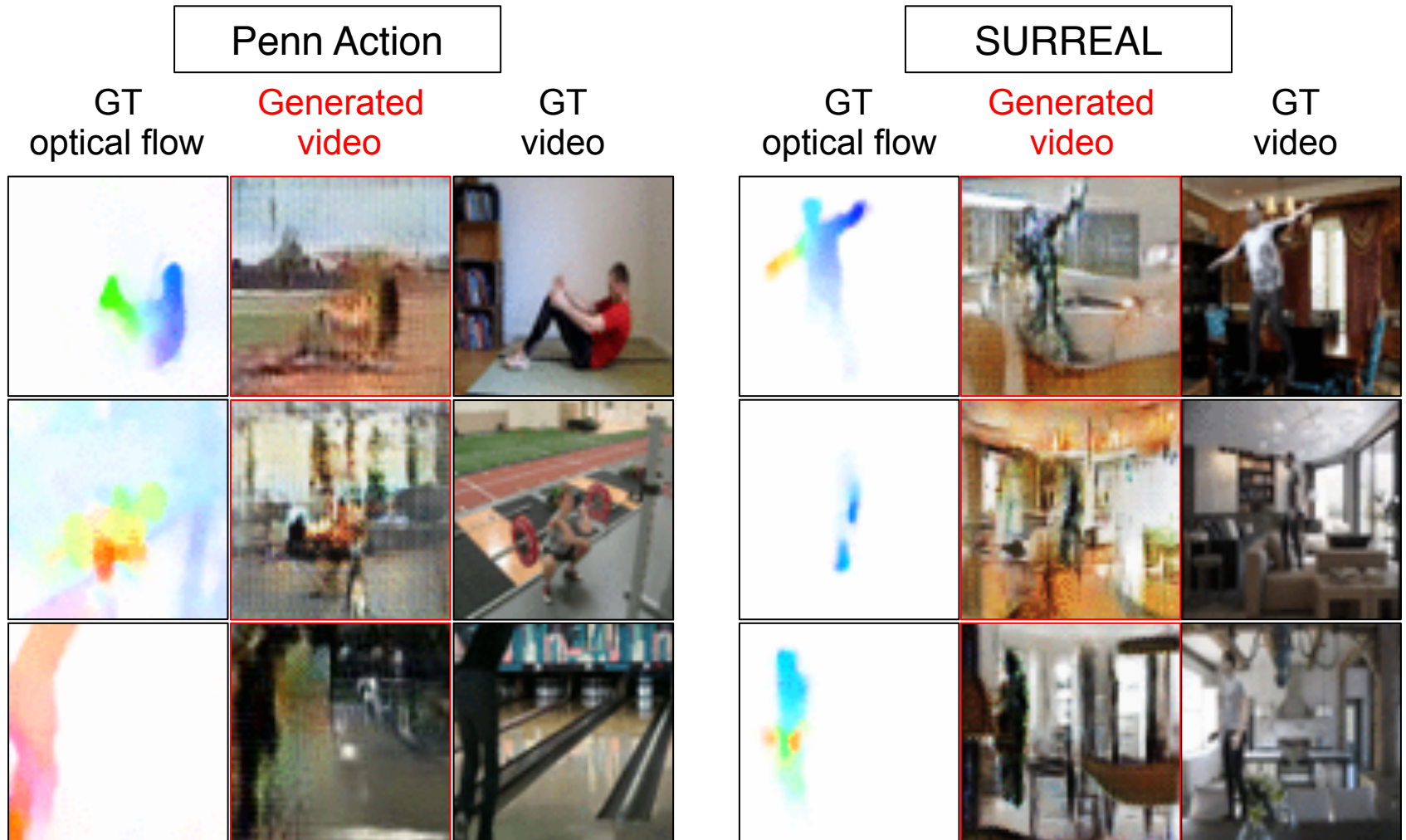
手法: TGAN [M. Saito et al, arxiv]

- 2D convolutionをXY方向にかけた後、1D convolutionをT方向に
  - 動き情報が見た目情報を抽象化された状態でしか取れてい



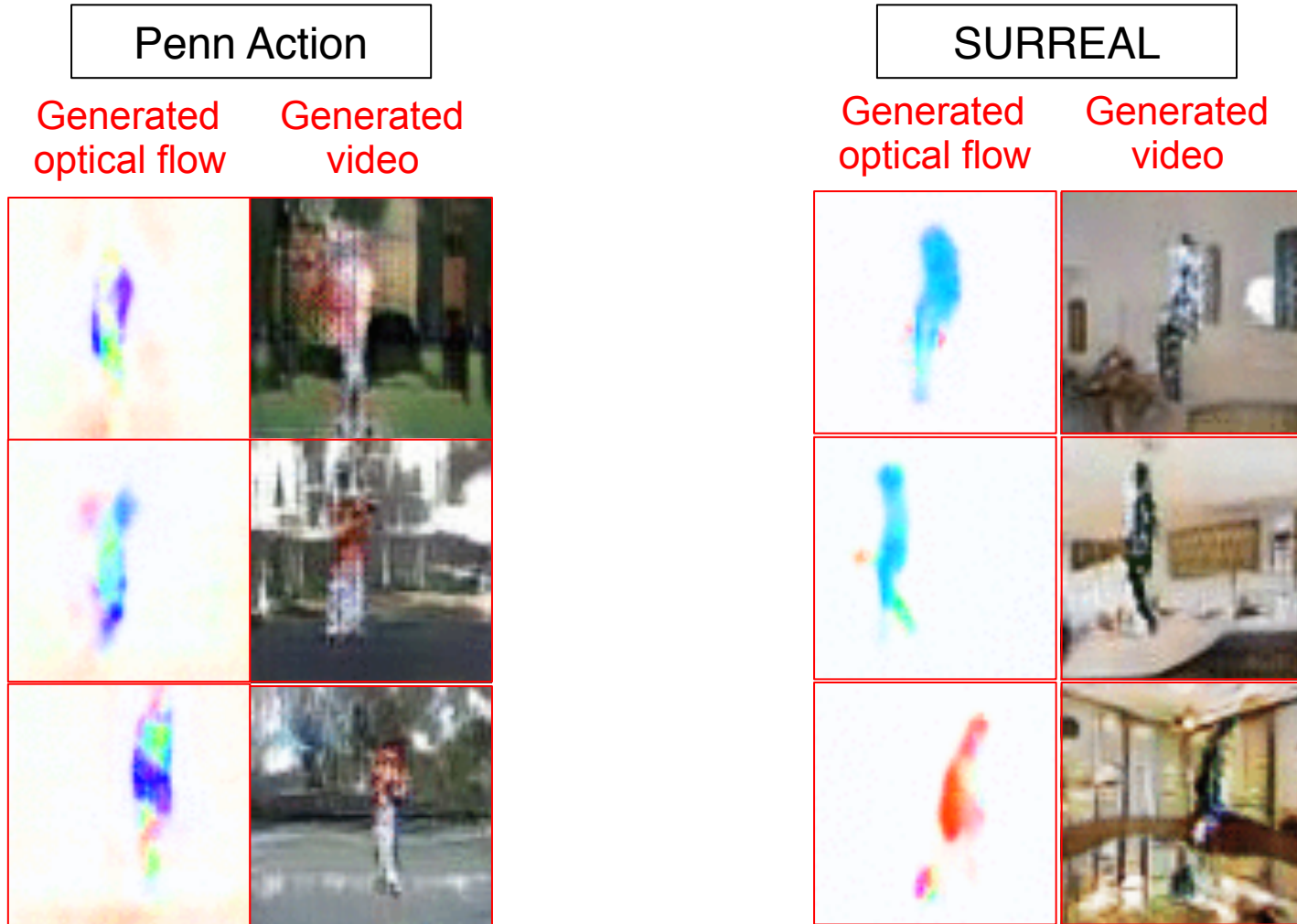
# 実験

- 生成結果: TextureGAN (Optical flowを与えた場合)



# 実験

- 生成結果: FTGAN (Optical flowも生成)





# 実験

未完成

## □ VGANとの比較

### ■ Penn Action

VGAN



TextureGAN



FTGAN





# 実験

未完成

## □ VGANとの比較

### ■ SURREAL

VGAN

TextureGAN

FTGAN

