

MoCoGAN: Decomposing Motion and Content for Video Generation

CVPR 2018

Sergey Tulyakov,
Snap Research

`stulyakov@snap.com`

Ming-Yu Liu, Xiaodong Yang, Jan Kautz
NVIDIA

`{mingyul,xiaodongy,jkautz}@nvidia.com`

Motivation

- deep feature representation of data in a unsupervised manner
- generate novel data for various applications
- very good image generation models
- weather prediction, autonomous driving



<https://thispersondoesnotexist.com/>

Challenges

- both the appearance model and the motion model
- the time dimension brings in a huge amount of variations
- human beings are more sensitive to motion

Previous approach

VGAN, TGAN: video as a point in latent space

Cons:

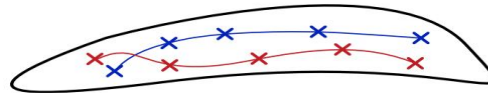
- Complexity
- Fixed length video

Approach

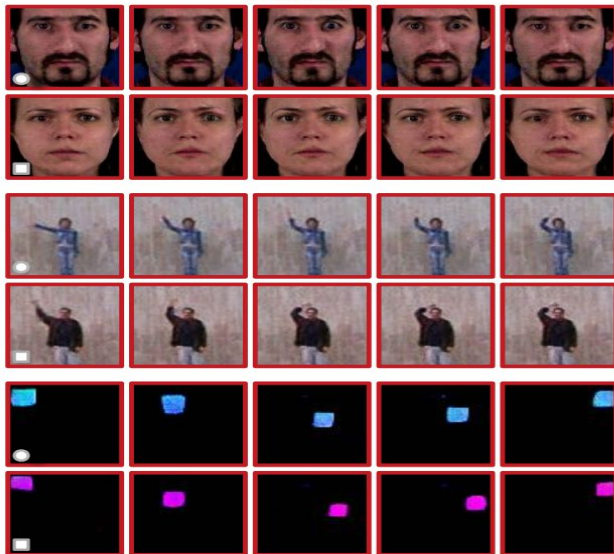
Content subspace



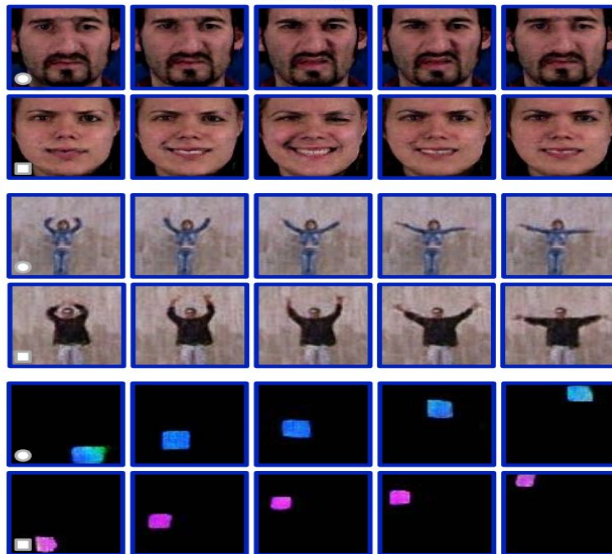
Motion subspace



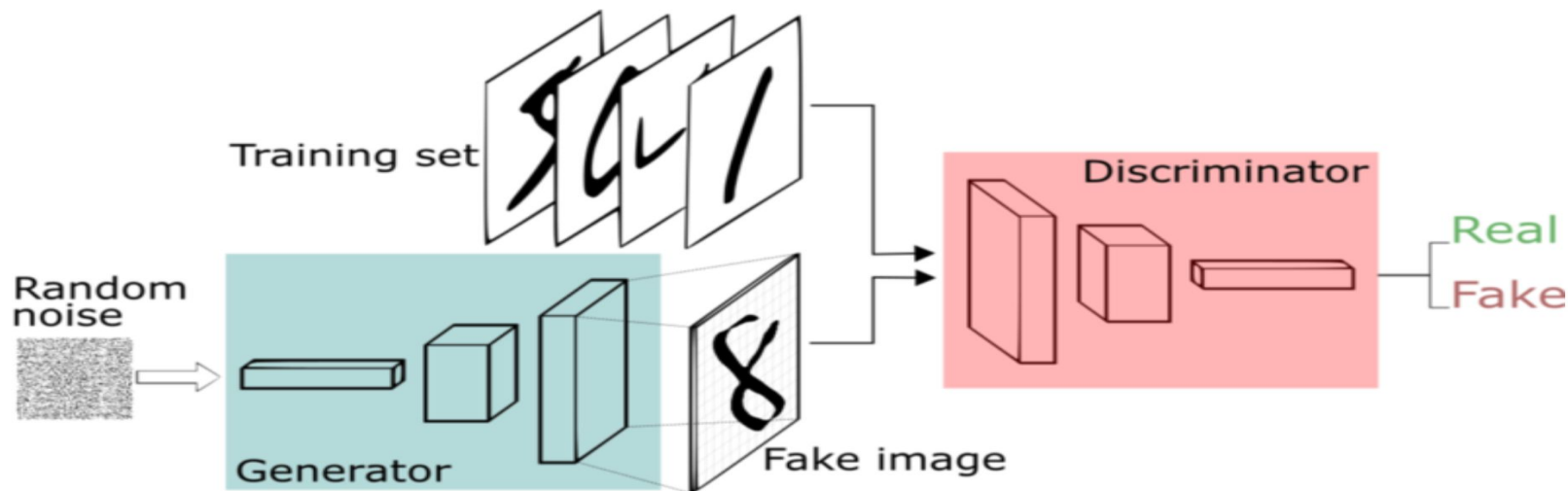
Motion 1



Motion 2



GAN (Image)



$$\max_{G_I} \min_{D_I} \mathcal{F}_I(D_I, G_I)$$

$$\mathcal{F}_I(D_I, G_I) = \mathbb{E}_{\mathbf{x} \sim p_X} [-\log D_I(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_{Z_I}} [-\log(1 - D_I(G_I(\mathbf{z})))]$$

MoCoGAN

Latent space: $Z_I = Z_C \times Z_M$ $Z_C = \mathbb{R}^{d_C}$, $Z_M = \mathbb{R}^{d_M}$ $d = d_C + d_M$

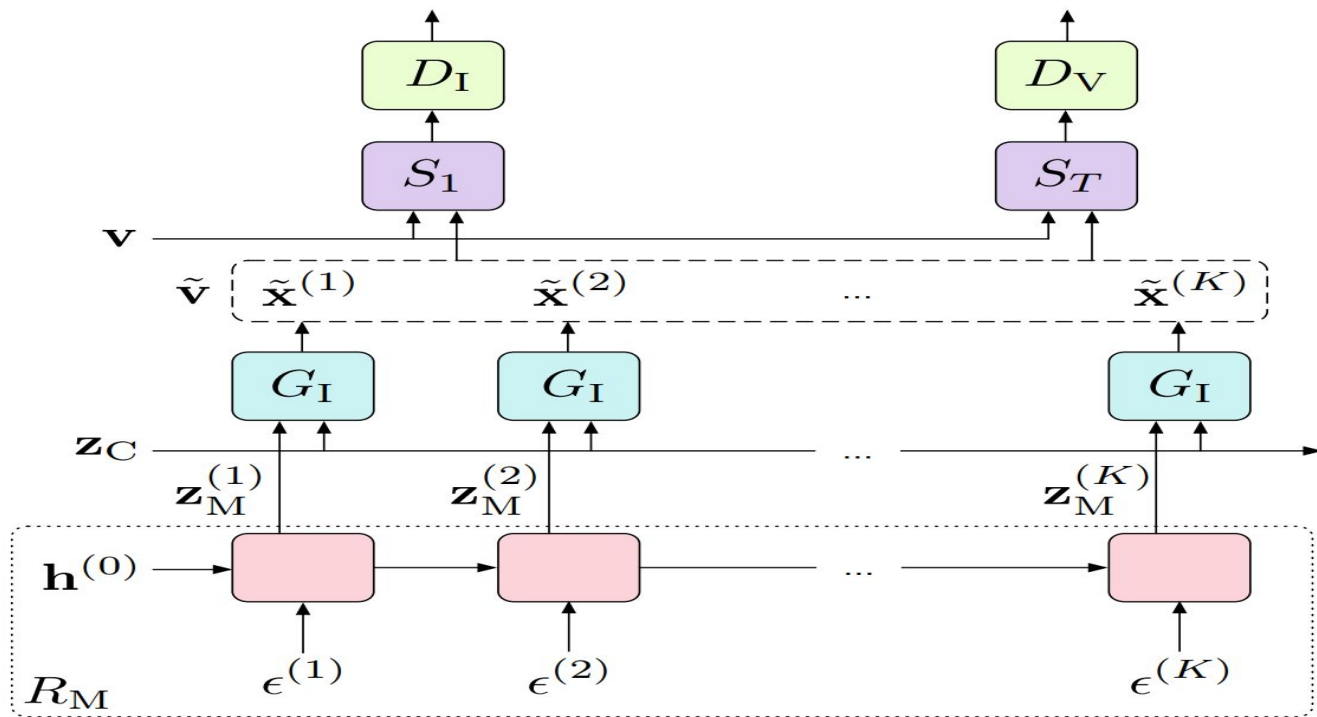
$$[\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(K)}] = \left[\begin{bmatrix} \mathbf{z}_C \\ \mathbf{z}_M^{(1)} \end{bmatrix}, \dots, \begin{bmatrix} \mathbf{z}_C \\ \mathbf{z}_M^{(K)} \end{bmatrix} \right]$$

Content subspace: $\mathbf{z}_C \sim p_{Z_C} \equiv \mathcal{N}(\mathbf{z}|0, I_{d_C})$

Motion subspace: $[\epsilon^{(1)}, \dots, \epsilon^{(K)}] \xrightarrow{\text{RNN}} Z_M^{(1)}, \dots, Z_M^{(K)}$

$$\epsilon^{(k)} \sim p_E \equiv \mathcal{N}(\epsilon|0, I_{d_E})$$

MoCoGAN



MoCoGAN

G_I	Configuration
Input	$[\mathbf{z}_a \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \mathbf{z}_m \sim R_M]$
0	D CONV-(N512, K6, S0, P0), BN, LeakyReLU
1	D CONV-(N256, K4, S2, P1), BN, LeakyReLU
2	D CONV-(N128, K4, S2, P1), BN, LeakyReLU
3	D CONV-(N64, K4, S2, P1), BN, LeakyReLU
4	D CONV-(N3, K4, S2, P1), BN, LeakyReLU

D_I	Configuration
Input	height \times width \times 3
0	CONV-(N64, K4, S2, P1), BN, LeakyReLU
1	CONV-(N128, K4, S2, P1), BN, LeakyReLU
2	CONV-(N256, K4, S2, P1), BN, LeakyReLU
3	CONV-(N1, K4, S2, P1), Sigmoid

D_V	Configuration
Input	$16 \times$ height \times width \times 3
0	CONV3D-(N64, K4, S1, P0), BN, LeakyReLU
1	CONV3D-(N128, K4, S1, P0), BN, LeakyReLU
2	CONV3D-(N256, K4, S1, P0), BN, LeakyReLU
3	CONV3D-(N1, K4, S1, P0), Sigmoid

MoCoGAN

$$\max_{G_I, R_M} \min_{D_I, D_V} \mathcal{F}_V(D_I, D_V, G_I, R_M)$$

$$\mathbb{E}_{\mathbf{v}}[-\log D_I(S_1(\mathbf{v}))] + \mathbb{E}_{\tilde{\mathbf{v}}}[-\log(1 - D_I(S_1(\tilde{\mathbf{v}})))] + \\ \mathbb{E}_{\mathbf{v}}[-\log D_V(S_T(\mathbf{v}))] + \mathbb{E}_{\tilde{\mathbf{v}}}[-\log(1 - D_V(S_T(\tilde{\mathbf{v}})))] ,$$

Categorical dynamics: $\left[\left[\begin{array}{c} \mathbf{z}_A \\ \epsilon^{(1)} \end{array} \right], \dots, \left[\begin{array}{c} \mathbf{z}_A \\ \epsilon^{(K)} \end{array} \right] \right]$

$$\mathcal{F}_V(D_I, D_V, G_I, R_M) + \lambda L_I(G_I, Q)$$

Dataset

Name	Number of videos
Shape motion	4000
MUG Facial Expression	1254
Tai-Chi	4500
Weizmann Human Actions	81
UCF101	13 220

Video Generation Performance

ACD	Shape Motion	Facial Expressions
Reference	0	0.116
VGAN [40]	5.02	0.322
TGAN [30]	2.08	0.305
MoCoGAN	1.79	0.201

ACD: Average Content Distance

L2 distance between average color vectors (Shape motion)

L2 distance between OpenFace feature vector (Facial expression)

Video Generation Performance

Inception score				Facial expression and Tai-Chi datasets		
	VGAN	TGAN	MoCoGAN	User preference, %	Facial Exp.	Tai-Chi
UCF101	$8.18 \pm .05$	$11.85 \pm .07$	$12.42 \pm .03$	MoCoGAN / VGAN	84.2 / 15.8	75.4 / 24.6
				MoCoGAN / TGAN	54.7 / 45.3	68.0 / 32.0

Inception score:

- Images have variety
- Each image distinctly looks like something

Qualitative evaluation



Categorical Video Generation

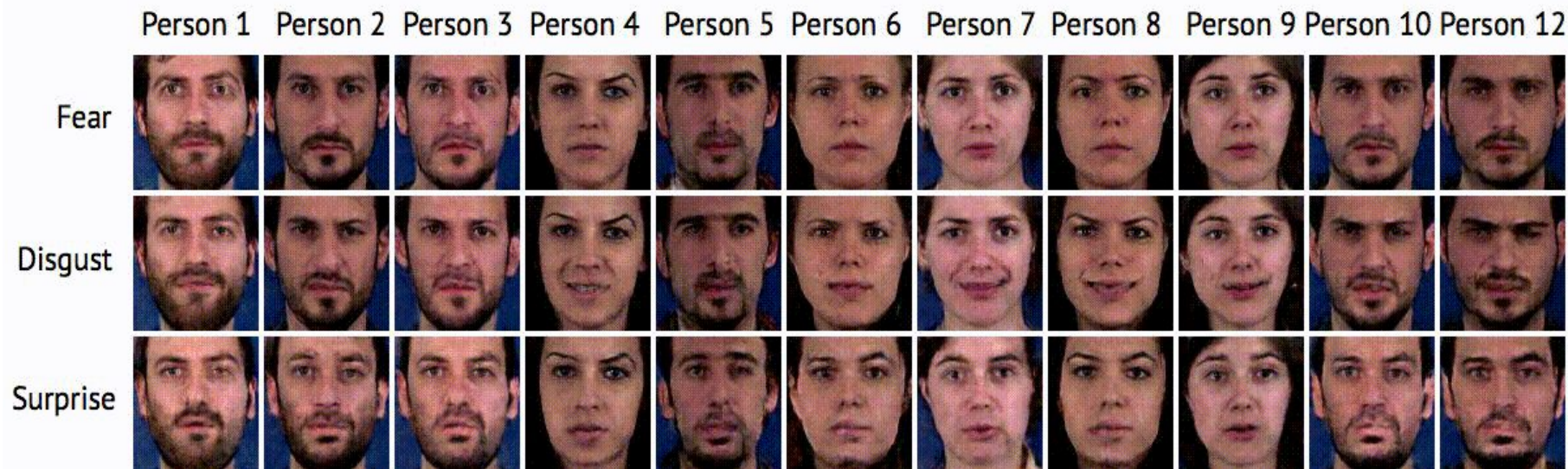


Image-to-video Translation



User preference, %	Tai-Chi
MoCoGAN / C-VGAN	66.9 / 33.1
MoCoGAN / MCNET	65.6 / 34.4

Strengths

- A Novel GAN framework for video generation.
- Can control content and motion in video generation.
- Several experiments with multiple datasets.

Weaknesses

- Small and medium sized dataset.
- Assuming there is a fixed content for the whole video.
- Does not work well with bigger dataset (Kinetics or even UCF101).

Discussion questions

- Why do we need image discriminator?
- Do you think the assumption of fixed content for the video is reasonable? Can we learn everything from data?

Thank you
Questions?