

# Latent Image Animator: Learning to Animate Images via Latent Space Navigation

ICLR 2022

Yaohui Wang, Di Yang, Francois Bremond and Antitza Dantcheva

Inria

**Presenter: Jaeseong Lee**

*Computer Vision Seminar*  
*21 Feb 2022*

# Contents

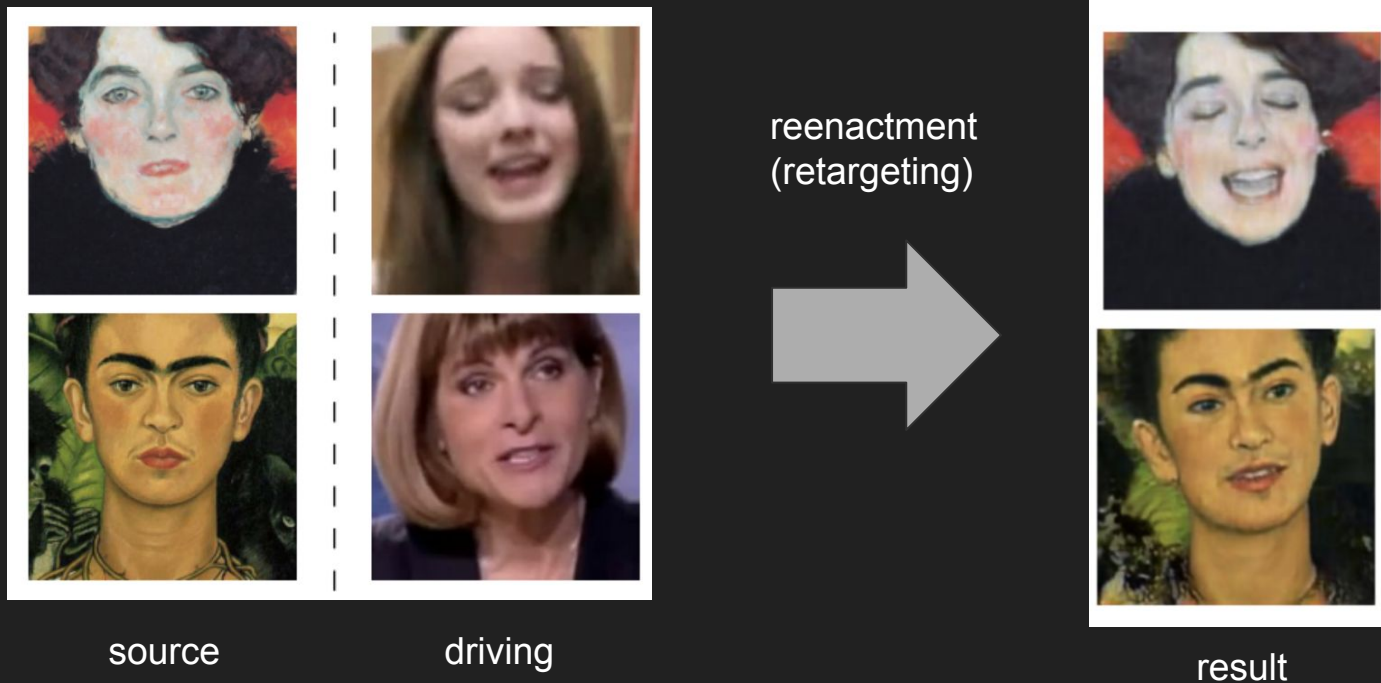
## 1. Preliminary for Neural Talking Heads

## 2. Main paper

- a. Introduction
- b. Method**
- c. Experiments
- d. Additional analysis
- e. Conclusion

# Preliminary(cont'd)

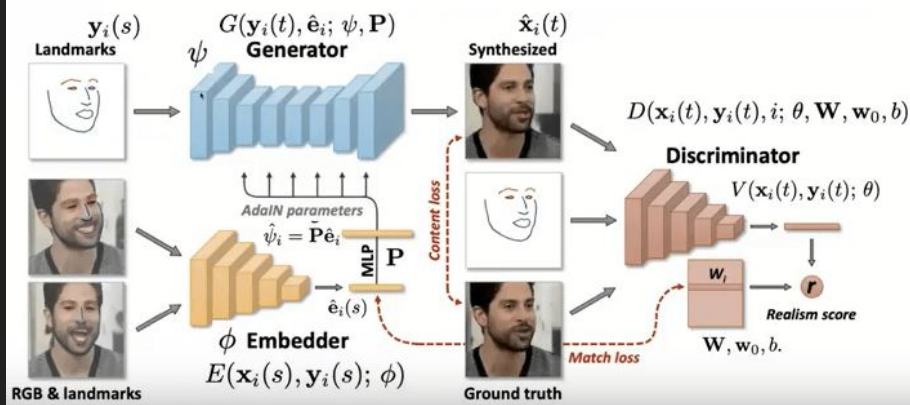
Neural Talking Head: Reenacting driving images(or a video)' **Facial Expression & Head Rotation** to a source image



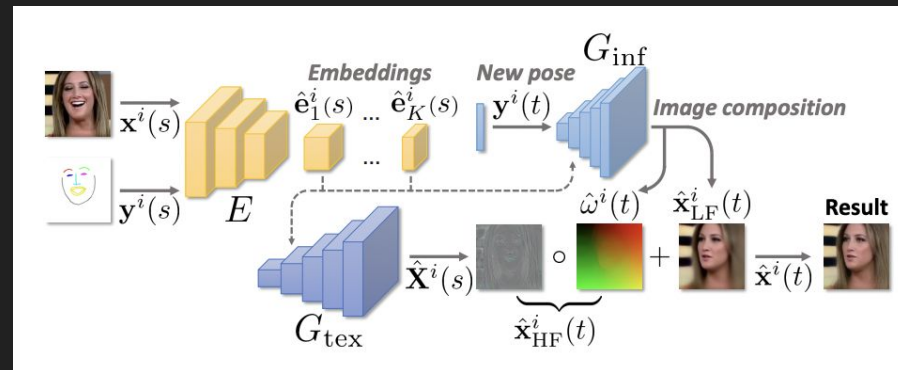
# Preliminary(cont'd)

off-the-shelf  
keypoints-based

## Architecture and notation



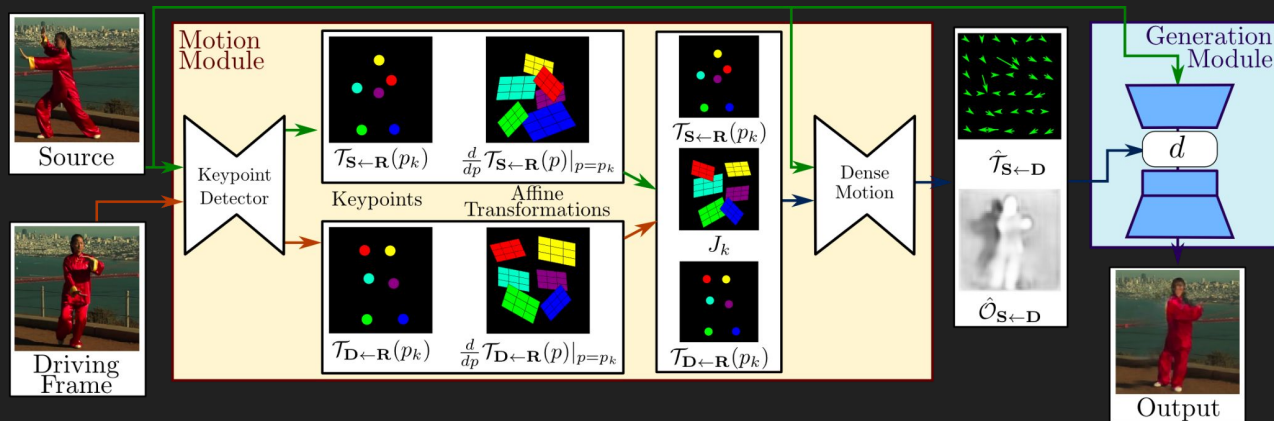
Few-Shot Adversarial Learning of  
Realistic Neural Talking Head Models  
(ICCV 2019)  
**Samsung AI, Moscow**



Fast Bi-layer Neural Synthesis of  
One-Shot Realistic Head Avatars  
(ECCV 2020)  
**Samsung AI, Moscow**

# Preliminary(cont'd)

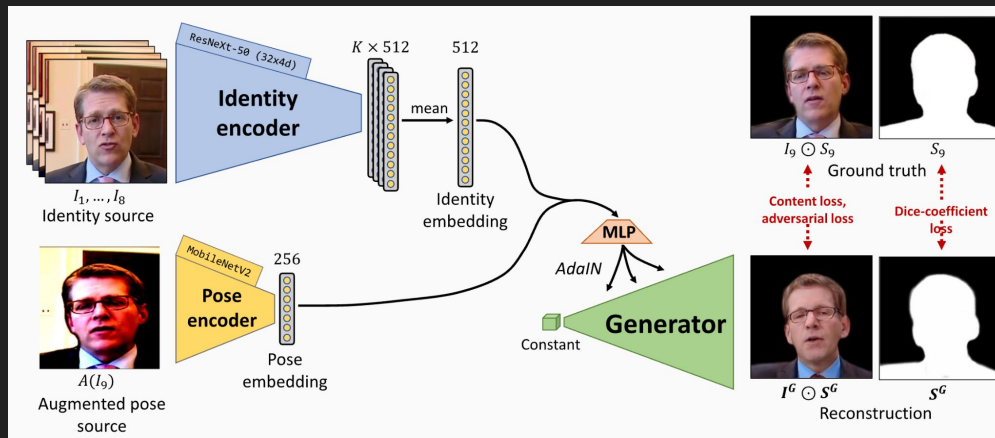
self-sup  
keypoints-based



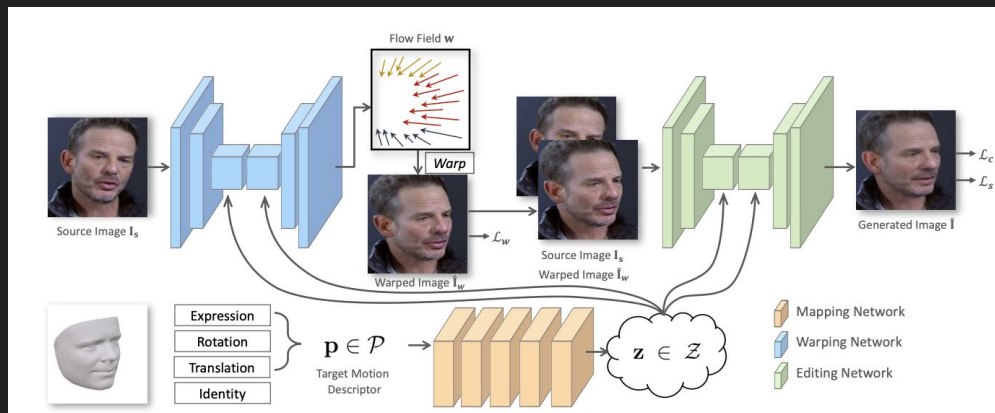
First Order Motion Model for Image  
Animation  
(NeurIPS 2019)  
**University of Trento**

# Preliminary

pose  
descriptor-based  
(AdaIN)



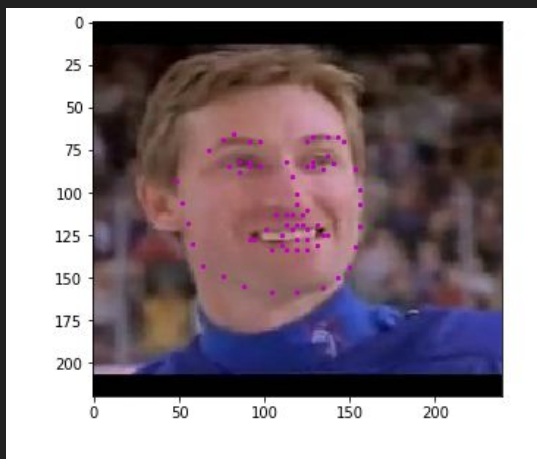
Neural Head Reenactment with Latent  
Pose Descriptors  
(CVPR 2020)  
Samsung AI, Moscow



PIRenderer: Controllable Portrait Image  
Generation via Semantic Neural  
Rendering  
(ICCV 2021)  
Peking University

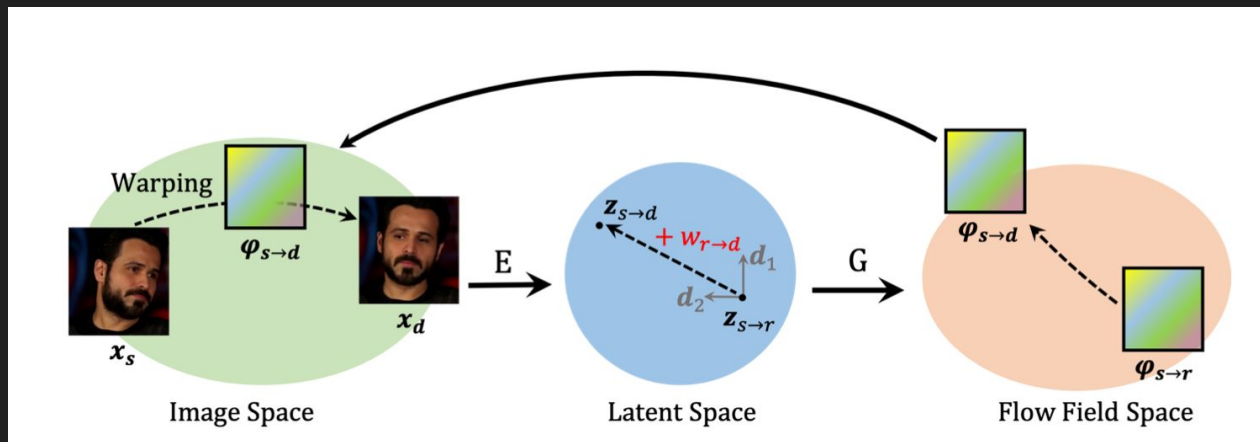
# Introduction(cont'd)

- The poses constraints(e.g., keypoint, 3DMM parameters) on applications, where such representations of unseen testing images might be **fragmentary** or **missing**.



# Introduction(cont'd)

- Existing methods using structural information(e.g., keypoints or descriptors) are hard to perform well when the source and the driving have **large appearance variation**. ▶▶ Latent Space Navigating



General pipeline of LIA

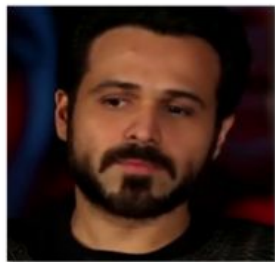


# Introduction

source



driving



Warp

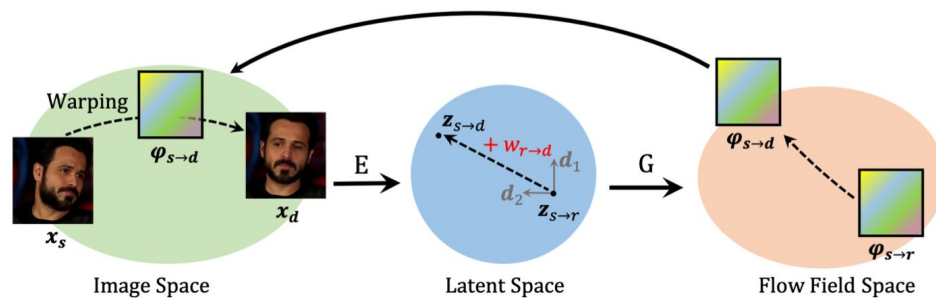
$w(r \rightarrow d)$

$z(s \rightarrow r)$

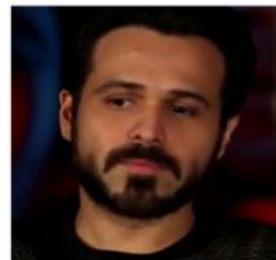
Latent Space

$z(s \rightarrow d)$

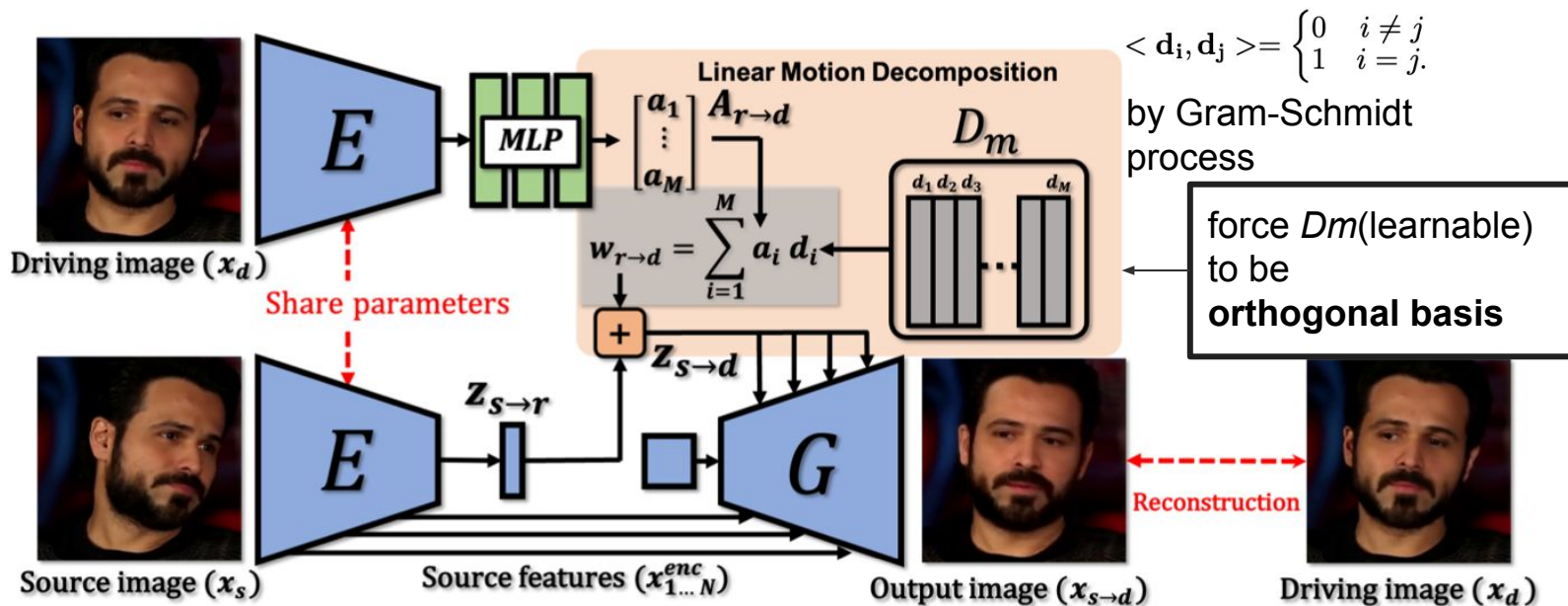
$\varphi_{s \rightarrow d}$



result



# Method(cont'd)



Objectives:

$$\mathcal{L}_{recon}(x_{s \rightarrow d}, x_d) = \mathbb{E}[\|x_d - x_{s \rightarrow d}\|_1].$$

$$\mathcal{L}_{vgg}(x_{s \rightarrow d}, x_d) = \mathbb{E}\left[\sum_n^N \|F_n(x_d) - F_n(x_{s \rightarrow d})\|_1\right],$$

$$\mathcal{L}_{adv}(x_{s \rightarrow d}) = \mathbb{E}_{x_{s \rightarrow d} \sim p_{rec}}[-\log(D(x_{s \rightarrow d}))],$$

# Method

## Goal

$$z_{s \rightarrow t} = (z_{s \rightarrow r} + w_{r \rightarrow s}) + (w_{r \rightarrow t} - w_{r \rightarrow 1})$$

when,  $s \neq 1$

@ Inference time → 'relative transfer'

$$= z_{s \rightarrow s} + (w_{r \rightarrow t} - w_{r \rightarrow 1}), \quad t \in \{1, \dots, T\}.$$

- $z(s \rightarrow s)$ : reconstruction
- $w(r \rightarrow t) - w(r \rightarrow 1)$ : motion from 1 to  $t$

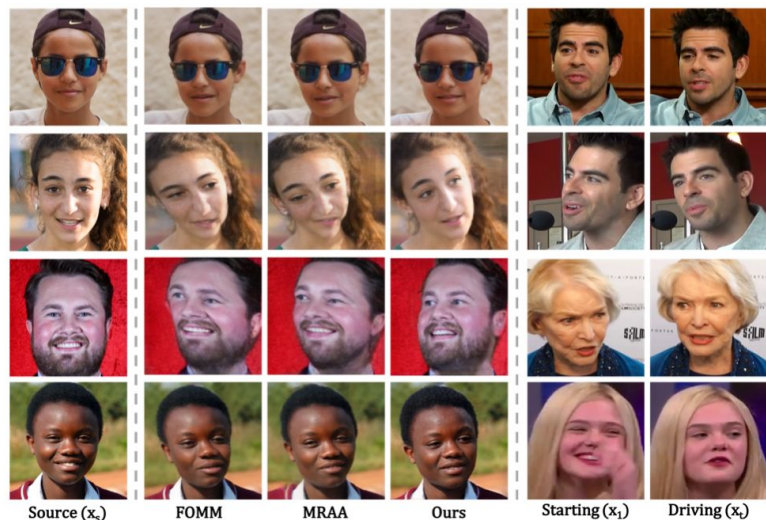
\* The more similar btw  $x_s$  &  $x_1$ , the more replication

when,  $s = 1$

cf) @ Training time → 'absolute transfer'

$$= z_{s \rightarrow r} + w_{r \rightarrow t}, \quad t \in \{1, \dots, T\}.$$

# Experiments



out-of-domain source  
images(FFHQ)  
inference



same-identity

Method	$\mathcal{L}_1$	VoxCeleb		
		AKD	AED	LPIPS
X2Face	0.078	7.687	0.405	-
Monkey-Net	0.049	1.878	0.199	-
FOMM	0.046	1.395	0.141	0.136
MRAA w/o bg	0.043	<b>1.307</b>	0.140	0.127
Ours	<b>0.041</b>	1.353	<b>0.138</b>	<b>0.123</b>

cross-identity

	VoxCeleb	GermanAudio
FOMM	0.323	0.456
MRAA	0.308	0.454
Ours	<b>0.161</b>	<b>0.406</b>

# Additional analysis(cont'd)

1. Literally, xr represents what?

A) Canonical(frontal/neutral) pose of xs, **regardless of original poses.**

source

reference



# Additional analysis(cont'd)

## 2. Manipulation on $D_m$



Figure 10: Manipulation of motion dictionary.

# Conclusion

‘LIA opens a new door in design of interpretable generative models for video generation.’