







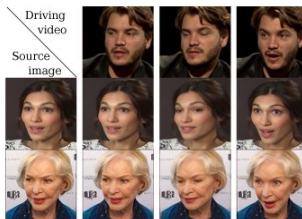
# Introduction

- AI and Deep learning has accelerated this trend
- Tasks such as generation of talking face videos, speech/text based lip synthesis and movie generation have witnessed a tremendous growth.



**Figure 2:** The actor Mark Hamill being de-aged using AI for the part of Luke Skywalker in the Mandalorian series.

## Some more Examples



**Figure 3:** A SOTA model (FOMM) used to animate a "source image" as per a "driving video"



**Figure 4:** An Editing framework (STIT) editing a video featuring Kamala Harris













g as  $\mathcal{M}$  is known, all we need is a *differentiable*  
*parameterization* of  $\mathcal{M}$  and then we can simply *progress*  
 that map our input frames to correct outputs using  
 function.

**do we have the differentiable parameterization?**

**Ans : NO!**

we use **Generative Modelling!**

- As long as  $\mathcal{M}$  is known, all we need is a *differentiable parameterization* of  $\mathcal{M}$  and then we can simply *progress* along those regions that map our input frames to correct outputs using a suitable loss function.
- **But do we have the differentiable parameterization?**
- **Answer : NO!**
- To solve this use **Generative Modelling!**







# StyleGAN

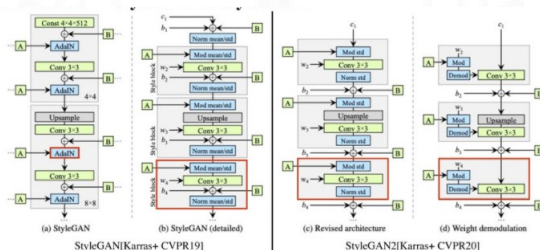
- Choose Karras et.al's pretrained Face StyleGAN trained on the FFHQ Dataset.
- This can generate large diversity of generated images including diverse faces in various poses and natural images such as animals, cars etc.



**Figure 6:** The diverse set of images capable of being generated through a StyleGAN

# StyleGAN Architecture

Architecture gives us *flexibility* in choice of Latent Space as well!



**Figure 7:** Architecture of a StyleGAN



# StyleGAN Architecture

The algorithm of the generator is as follows:

- ➊  $z \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), z \in \mathbb{R}^{512}$
- ➋ Get  $w = m(z), w \in \mathbb{R}^{512}$ .  $m(\cdot)$  is a learnt network termed as the **mapping network**
- ➌ Pass  $w$  through learnt affine transformations  $A_1, \dots, A_n$
- ➍ Add to convolution layer corresponding to successive resolution (8x8 – 1024x1024)
- ➎ Obtain HD Image



# How to traverse the Latent Space?

We need to perform curve traversal in our chosen latent space. There are broadly two kinds of methods employed for this purpose:

- 1 **Optimization** - Involves optimization for a single image or a set of images using latents initialized in the latent space. Unstable for larger sets
- 2 **Geometry Based** - Traversal along directions corresponding to a *local* or *global basis* of directions in the latent space. Stable for larger sets. We shall use this

# Geometry Aware Traversal

Two kinds of curve traversal methods.

- 1 **Local Basis:** Define for every  $w \in \mathcal{W}$  separately as the basis of the Tangent Space at  $w$  ( $\mathcal{T}_w$ ). As  $m(\cdot)$  is a black-box neural network, to find this we need  $z \in \mathcal{Z} \ni m(z) = w$ . Then  $dm_z$ , (**pushforward**), the jacobian of  $m$  at  $z$  can be computed.

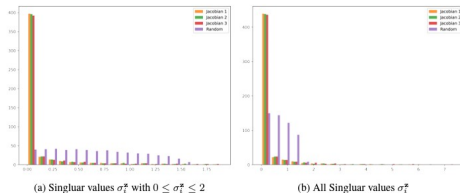
$$dm_z : \mathcal{T}_z \mapsto \mathcal{T}_w \quad (3)$$

To find the basis perform a dimensional reduction such as PCA/SVD on  $dm_z$  (**Manifold Hypothesis**)

भारतीय विज्ञान संस्थान

# Geometry Aware Traversal

Validity of the Manifold Hypothesis in this case been seen through Fig 7. which shows the distribution of singular values for 3 random  $z$ , a very significant chunk ( $> 400$  out of 512) is 0 or slightly larger)



**Figure 8:** Distribution of singular values of the jacobian of the mapping network

# Local Basis method in practice

- Very stable traversal method that has the advantages of preserving facial identity and smooth transitions
- It needs to know the  $z$  at every stage that maps to an intermediate  $w$  which greatly reduces its viability.

भारतीय विज्ञान संस्थान

# Geometry Aware Traversal - Global Basis

② **Global Traversal:** Assume  $\exists$  *global basis* of directions for the manifold  $\mathcal{W}$ ; algorithm to find this:

- ① Sample  $N$  vectors  $z_{1:N}$  from  $\mathcal{N}(\mathbf{0}, \mathbf{I})$  where  $N$  is necessarily very large
- ② Get  $w_{1:N}$  corresponding to these from the mapping network
- ③ Perform a dimensionality reduction algorithm such as PCA or SVD
- ④ Take the top  $k$  components to be the basis  $V$
- ⑤ Traversal is now given by the equation

$$w' = w + Vx \quad (4)$$

where  $x$  is a control vector used to control the extent of traversal along the basis directions.

भारतीय विज्ञान संस्थान

# Geometry Aware Traversal - Global Basis

- Global Traversal is less accurate but *more general*
- $x$  can be learnt as  $f(w)$  depending on task
- Use this method based on video generation framework - **MocoGANHD**



# MocoGANHD

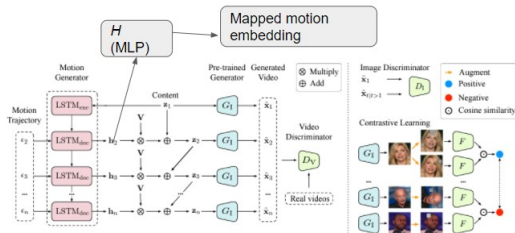
- MocoGANHD (Tian et.al , *ICLR 2021* generates synthetic videos on top of a pretrained StyleGAN
- Given  $w \in \mathcal{W}$  returns *ordered set* of latents,  $W_T = \{w_1, \dots, w_T\}$  (**Video with same identity**)
- Learns control vector per frame for global traversal
- Learns *temporal consistency* between frames
- Learns *realistic motions* i.e head pose and expressions

# MocoGANHD-Architectural Innovations

- Utilizes a conditional **LSTM** a sequential architecture defined for an ordered series of time steps  $1 : T$  (during inference can work for arbitrary  $T$ )
- Output of network at time  $t$ ;  $x_t = f_t(x_{t-1}, \epsilon_t)$ ;  $\epsilon_t$  is random vector corr. to motion at  $t$ ,  $f_t$  is learnable
- $t^{th}$  frame of video obtained as:

$$w_t = w + Vx_t \quad (5)$$

# MocoGANHD



**Figure 9:** Architecture of MocoGANHD



**Figure 10:** Sample videos generated by MocoGANHD

# Our innovations for the Audio Sync Setting

## Generalizing MocoGANHD!

- 1 Substitute Audio conditioning in-place of motion
- 2 Starting  $w$  needs to be known always. Finding  $w$  for arbitrary image is *inversion*. Figure shows how hard problem is (SOTA). **This problem is resolved as well.**



**Figure 11:** Inversion of sample Video Frames

# Solving the Inversion Problem

- Inversion is in fact a *separate* problem.
- For audio-sync by default datasets pair audio segments with frames belonging to different timestamps from the same video
- This can give a reconstruction prior for supervision
- **Not needed here** as generator can synthesize high quality images
- Therefore we can work with generated images!

# Generality is built from the ground up!

- Our method uses no ground truth image supervision
- If our method works for generated images we can solve for arbitrary real images!
- The onus is now on the inversion method - **modularity attained!**

भारतीय विज्ञान संस्थान

# Our Architecture

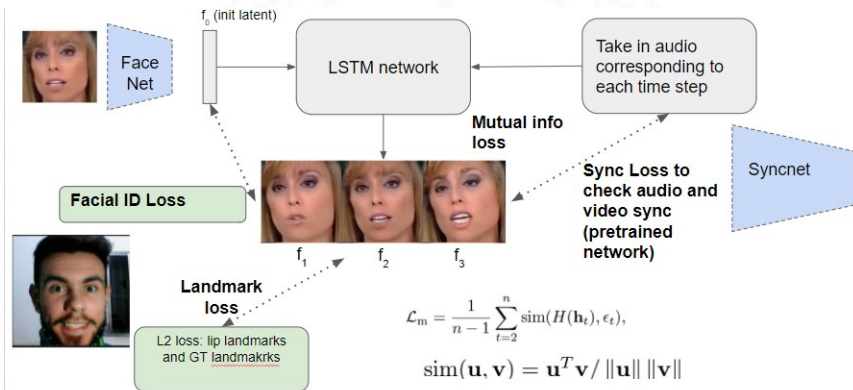
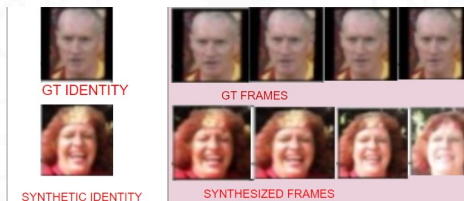


Figure 12: Our architecture

# Results



**Figure 13:** Expression transfer at low resolution  $128 \times 128$

Shows expression transfer from a given identity to a synthetic identity made from randomly sampling in the latent space. Done in low resolution to show the consistency.





# Unresolved problems

Problems yet to be resolved:

- ① *Identity Jitters* : Unstable with respect to facial identity for large  $T$ .
- ② *Motion Consistency* : Audio conditioning seems to be **insufficient** for motion. Explicit motion priors may be needed
- ③ *Improving Traversal*: More accuracy needed. Can utilize recent works utilizing Differential Geometry such as Riemannian CNFs (Mathieu et.al, Neurips 2020), Moserflow (Best paper, Neurips 2021)
- ④ *Inversion*: Open problem. Theoretically feasible for a large class of images but real world solutions are still limited.



THANKS!

