

# Neural Implicit Morphing of Face Images

Guilherme Schardong, Tiago Novello, Hallison Paz, Iurii Medvedev, Vinícius da Silva, Luiz Velho and Nuno Gonçalves

ISR-UC, IMPA, PUC-Rio, INCM

16/04/2024

# Outline

Introduction

Problem Statement and Methodology

Results

Conclusions

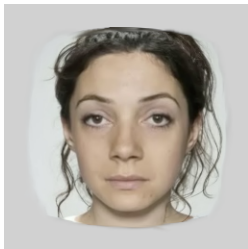
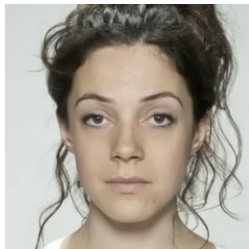
# Image Warping: What is it and how can we use it?

- Morphing = Warping + blending
- Warping maps the pixels of an image to a different geometry
  - **Perspective corrections**, lens distortion corrections, geometric transformations, morphing
- We “move” pixels in an image to a different location in another (or the same) image



# Image Warping: What is it and how can we use it?

- Morphing = Warping + blending
- Warping maps the pixels of an image to a different geometry
  - Perspective corrections, **lens distortion corrections**, geometric transformations, morphing
- We “move” pixels in an image to a different location in another (or the same) image





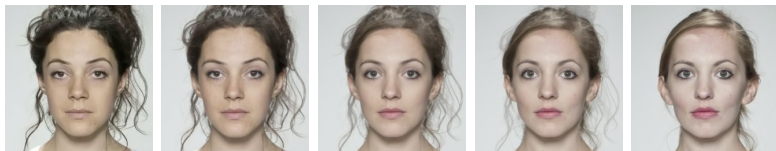
# Image Warping: What is it and how can we use it?

- Morphing = Warping + blending
- Warping maps the pixels of an image to a different geometry
  - Perspective corrections, lens distortion corrections, **geometric transformations**, morphing
- We “move” pixels in an image to a different location in another (or the same) image



# Image Warping: What is it and how can we use it?

- Morphing = Warping + blending
- Warping maps the pixels of an image to a different geometry
  - Perspective corrections, lens distortion corrections, geometric transformations, **morphing**
- We “move” pixels in an image to a different location in another (or the same) image



## Image Warping: How to apply it to faces?

- For faces, warping is used to align features (eyes, nose, mouth, ears, ...)
  - Even under perfect studio conditions, faces are still “misaligned” by nature



- How can we align them?

## Image Warping: How to apply it to faces?

- For faces, warping is used to align features (eyes, nose, mouth, ears, ...)
  - Even under perfect studio conditions, faces are still “misaligned” by nature



- How can we align them?

Find correspondences between every pixel in the images and warp them

# Face Warping: Issues and Motivation

- Matching every single pixel does not scale well
- Plus: How to find the correspondences between every single pixel?
  - Some have no correspondences at all (signs, scars, tattoos, ...)
- Additionally, we don't really need to match all pixels, only a sufficient subset of them
  - Then we simply transform the others based on this subset
- How to find/define a subset of pixels?

# Face Warping: Affine Transformation

- We can use facial landmarks as a viable subset of pixels:
  - use landmark detectors<sup>1,2</sup> or;
  - manually mark them
  - or both
- The matching part is trickier
- We have a correspondence between landmarks, but what about the other pixels?
  - How do we move them to their new locations?

---

<sup>1</sup>Kazemi and Sullivan, One Millisecond Face Alignment with an Ensemble of Regression Trees, CVPR 2014

<sup>2</sup>Lugaresi *et al.*, MediaPipe: A Framework for Building Perception Pipelines, arXiv, 2019

# Face Warping: Affine Transformation

- We can use facial landmarks as a viable subset of pixels:
  - use landmark detectors<sup>1,2</sup> or;
  - manually mark them
  - or both
- The matching part is trickier
- We have a correspondence between landmarks, but what about the other pixels?
  - How do we move them to their new locations?
- We can think of each landmark as influencing a region around it (thoughts of Voronoi)
- When we move a landmark, any pixels in this region must be moved as well

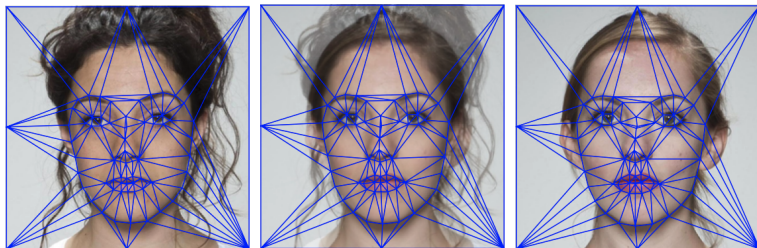
---

<sup>1</sup>Kazemi and Sullivan, One Millisecond Face Alignment with an Ensemble of Regression Trees, CVPR 2014

<sup>2</sup>Lugaresi *et al.*, MediaPipe: A Framework for Building Perception Pipelines, arXiv, 2019

# Face Warping: Affine Transformation

- It make sense to think that neighboring landmarks will influence each other
- Thus, we can triangulate landmarks to define the regions
- As the landmarks move, so will the triangles (and any pixels in them)





# Face Warping: Affine Transformation

- It make sense to think that neighboring landmarks will influence each other
- Thus, we can triangulate landmarks to define the regions
- As the landmarks move, so will the triangles (and any pixels in them)



# Good and Bad of warping by affine transformation

- Simple to implement
  - Delaunay triangulation is done once and shared for all faces
  - Landmark detection is done through 3rd party libraries
  - Image blending is done by linear interpolation as well
- Affine transforms may not be enough to align faces at all intermediate times
  - May lead to weird face proportions
  - Ghosting artifacts due to pose differences
  - It is, by nature, discrete

# Motivation

- Warping by affine transforms is not that bad, but can we improve it?
  1. More specifically: Can we improve it by performing a non-linear warping?
  2. Additionally, if the warping is **smooth**, we can exploit energy functionals proposed in the literature. Is it possible to produce a **continuous** and **smooth** warping?
- Neural networks may be used as a non-linear transformation
- More specifically: Sinusoidal neural networks

# SIRENs

- Sinusoidal Representation Networks (SIRENs) are simply multilayer-perceptrons with sine activation functions
- Sine was the first non-linearity proposed to be used with neural networks<sup>3</sup>
  - But were plagued by convergence issues
- Recent renaissance due to more in-depth studies<sup>4</sup>
- Eventually, better initialization methods made them feasible<sup>5</sup>
- Sines are  $C^\infty$ , thus **smooth**
  - Thus, we can calculate the (high-order) derivatives of output w.r.t the input

---

<sup>3</sup>Lapedes and Farber, Nonlinear signal processing using neural networks: Prediction and system modelling, 1987

<sup>4</sup>Parascandolo *et al.*, Taming the waves: sine as activation function in deep neural networks, 2016

<sup>5</sup>Sitzmann *et al.*, Implicit neural representations with periodic activation functions, 2020

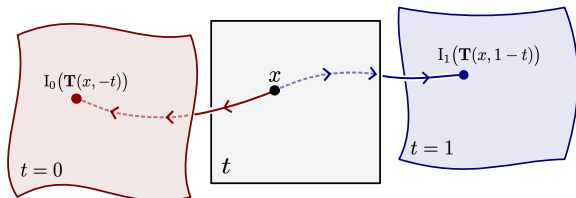
## On the definition and intuition of $\mathbf{T}$

We define a transform  $\mathbf{T} : \mathbb{R}^2 \times \mathbb{R} \rightarrow \mathbb{R}^2$ ,  $x \in \mathbb{R}^2$  and  $t \in [-1, 1]$ , such that:

- $\mathbf{T}(x, 0) = x$  (identity property)
- $\mathbf{T}(\mathbf{T}(x, t), -t) = x$  (inverse property)

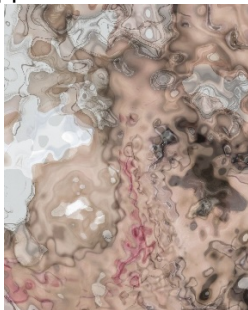
Given the landmarks  $p_i$  in  $I_0$  and their correspondences  $q_i$  in  $I_1$ :

- $\mathbf{T}(p_i, t) = \mathbf{T}(q_i, 1 - t)$  (landmark matching property)



## On the need for regularization

We want  $\mathbf{T}$  to be non-rigid, but such “non-rigidity” must be under control, or this may happen:



A way to achieve this is by introducing the **Thin-plate** property for regularization:

- $\min \|\mathbf{Hess}(\mathbf{T})(x, t)\|_F^2$

## Formal definition: Alignment Loss

$$\mathcal{L}(\theta) = \lambda_1 \mathcal{W}(\theta) + \lambda_2 \mathcal{D}(\theta) + \lambda_3 \mathcal{T}(\theta). \quad (1)$$

$\mathcal{W}(\theta)$ ,  $\mathcal{D}(\theta)$ ,  $\mathcal{T}(\theta)$  are the *warping*, *data*, *thin-plate* constraints

$$\mathcal{W}(\theta) = \underbrace{\int_{\mathbb{R}^2} \|\mathbf{T}(x, 0) - x\|^2 dx}_{\text{Identity constraint}} + \underbrace{\int_{\mathbb{R}^2 \times \mathbb{R}} \left\| \mathbf{T}(\mathbf{T}(x, t), -t) - x \right\|^2 dx dt}_{\text{Inverse constraint}}. \quad (2)$$

$$\mathcal{D}(\theta) = \int_{[0,1]} \|\mathbf{T}(p_i, t) - \mathbf{T}(q_i, 1 - t)\|^2 dt \quad (3)$$

$$\mathcal{T}(\theta) = \int_{\mathbb{R}^2 \times \mathbb{R}} \|\mathbf{Hess}(\mathbf{T})(x, t)\|_F^2 dx dt \quad (4)$$

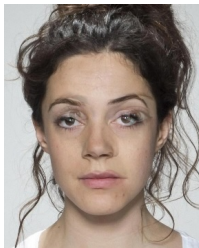
## Effects of each loss term



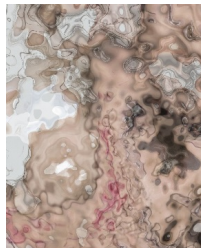
no inv



no id



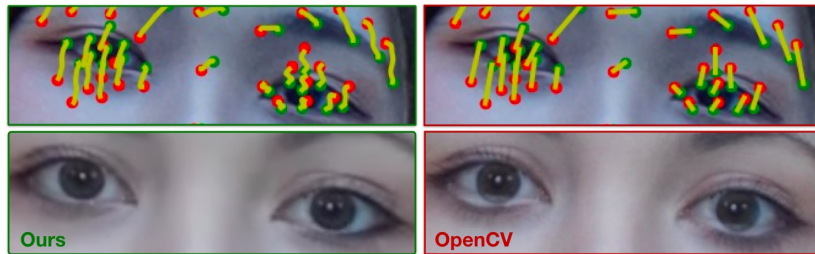
no  $\mathcal{D}(\theta)$



no  $\mathcal{I}(\theta)$

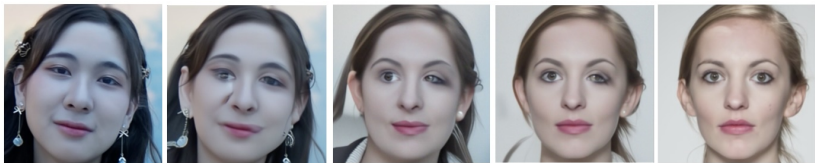


## Results: Differences in warping paths

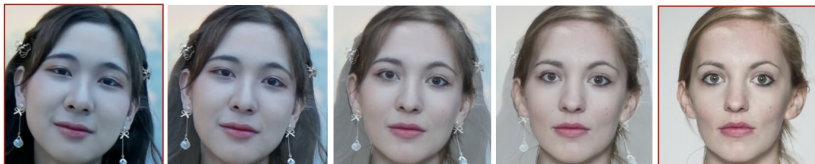


# Results: Dealing with rotations

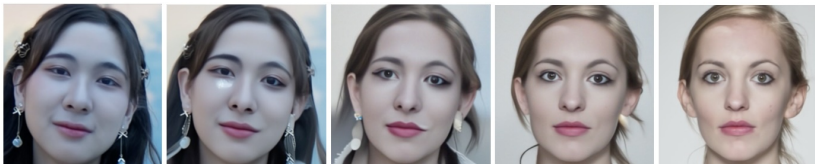
Pure diffAE



Neural warping + linear blending [Ours]



Neural warping + diffAE [Ours]



# Results: Dealing with poses, expressions and occlusions

## Different poses



## Different expressions



## Eyes occlusion + faces in the wild



# Results: Morphings between different ethnicities and genders



# Results: Morphings between different ethnicities and genders



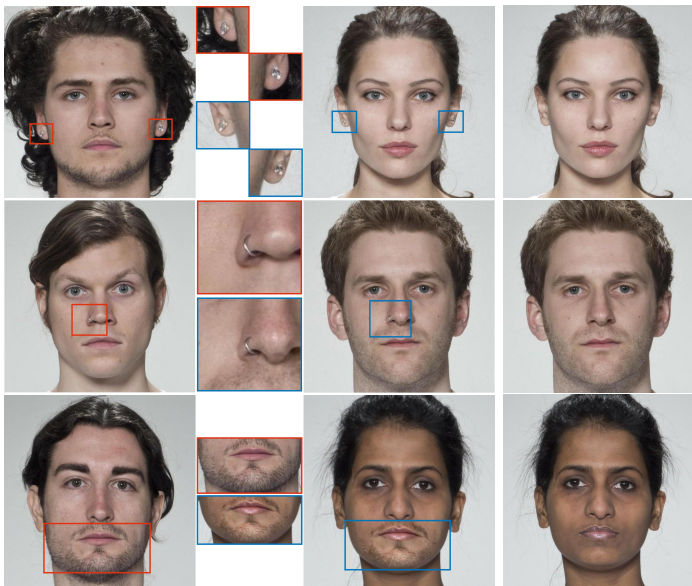
# Results: Morphings between different ethnicities and genders



# Results: Morphings between different ethnicities and genders

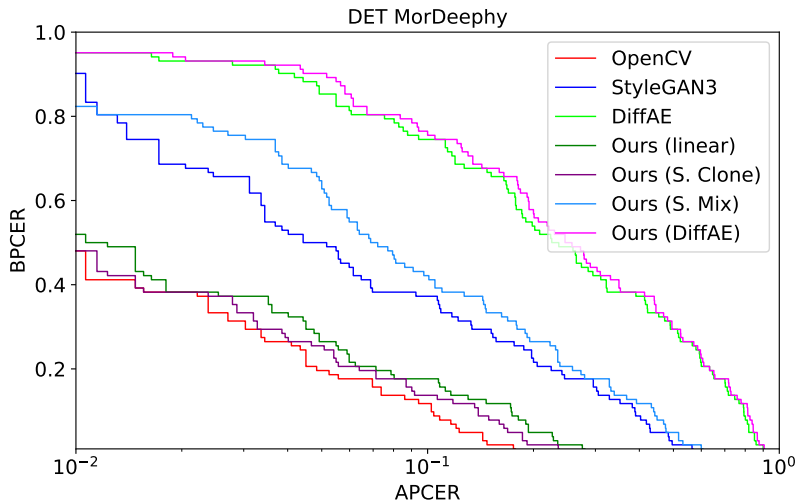


# Results: Feature transference





# Results: Morphing Attack Detection Effectiveness



## Results: Morphing failures



- Needs more time to converge. But why?
  - Maybe the initialization is not good enough

# Conclusions and Future Directions

- Introduction of classic energy functionals to use with neural networks is fairly straightforward
  - They are still very powerful and yield good results
- We can exploit the network smoothness to use its derivatives on the loss
- Separation of problems adds flexibility, since we can “mix and match” various techniques
- Investigate why some cases take so long to converge

# Thank you

