# Al-Driven Video Synthesis and its Implications





Visual Computing Group Prof. Matthias Nießner

## Visual Computing Group @ TUM

**Al-Driven Video** 



### Photo-realistic Image Synthesis

The Rendering Equation [Kajiya 86]

$$L_{\mathrm{o}}(\mathbf{x},\,\omega_{\mathrm{o}},\,\lambda,\,t)\,=\,L_{e}(\mathbf{x},\,\omega_{\mathrm{o}},\,\lambda,\,t)\,+\,\int_{\Omega}f_{r}(\mathbf{x},\,\omega_{\mathrm{i}},\,\omega_{\mathrm{o}},\,\lambda,\,t)\,L_{\mathrm{i}}(\mathbf{x},\,\omega_{\mathrm{i}},\,\lambda,\,t)\,(\omega_{\mathrm{i}}\,\cdot\,\mathbf{n})\,\,\mathrm{d}\,\omega_{\mathrm{i}}$$



## Need 3D Content for Rendering



Geometry

Textures

Material & Lighting

#### Computer Vision as Inverse Graphics

 $L_{\mathrm{o}}(\mathbf{x},\,\omega_{\mathrm{o}},\,\lambda,\,t)^{-1}$ 

#### Can we invert the Rendering Equation?



#### Computer Vision as Inverse Graphics

E(P) =





#### Priors: Parametric Face Model



[Blanz and Vetter 99] BlendShapes [Alexander et al. 09/10] Digital Emily [Chen et al. 14] FaceWarehouse

...







|P| = 6 + 80



Material / Reflection

$$\boldsymbol{P} = \begin{pmatrix} \boldsymbol{\Phi} \\ \boldsymbol{\alpha} \\ \boldsymbol{\beta} \end{pmatrix}$$

|P| = 6 + 80 + 80



#### **Expression Parameters**

$$\boldsymbol{P} = \begin{pmatrix} \boldsymbol{\Phi} \\ \boldsymbol{\alpha} \\ \boldsymbol{\beta} \\ \boldsymbol{\delta} \end{pmatrix}$$

|P| = 6 + 80 + 80 + 76



**Lighting Parameters** 

$$\boldsymbol{P} = \begin{pmatrix} \boldsymbol{\Phi} \\ \boldsymbol{\alpha} \\ \boldsymbol{\beta} \\ \boldsymbol{\delta} \\ \boldsymbol{\gamma} \end{pmatrix}$$

|P| = 6 + 80 + 80 + 76 + 27





![](_page_13_Picture_1.jpeg)

### Fitting Parametric Model to RGB Image

E(P) =

![](_page_14_Picture_2.jpeg)

![](_page_14_Picture_3.jpeg)

### Analysis-by-Synthesis

Given: Parametric Model *M*(*P*)

- 1. Render *M* with parameters  $P_k$
- 2. Compute diff. between rendering and target; i.e., E(P)
- 3. Update  $P_k \rightarrow P_{k+1}$ ; e.g., using differentiable renderer
- 4. If (diff > thresh) GOTO 1

E(P) =

![](_page_16_Picture_2.jpeg)

![](_page_16_Picture_3.jpeg)

$$E(P) = E_{col}(P)$$

Color Consistency

![](_page_17_Picture_3.jpeg)

![](_page_17_Picture_4.jpeg)

Distance in RGB Color Space

![](_page_17_Figure_6.jpeg)

$$E(P) = E_{col}(P) + E_{mrk}(P)$$

Color Consistency Feature Similarity

![](_page_18_Picture_4.jpeg)

![](_page_18_Picture_5.jpeg)

![](_page_18_Figure_6.jpeg)

![](_page_19_Figure_1.jpeg)

$$E(P) = E_{col}(P) + E_{mrk}(P) + E_{reg}(P)$$
Color Feature Regularization
Consistency Similarity

- Coarse-to-fine Gauss-Newton optimization (IRLS)
- Gradients through differentiable rendering

![](_page_21_Picture_1.jpeg)

### 3D Model + Image-based Rendering

![](_page_22_Picture_1.jpeg)

## 3D Model + Image-based Rendering

![](_page_23_Picture_1.jpeg)

Image-based mouth retrieval

![](_page_23_Picture_3.jpeg)

## 3D Model + Image-based Rendering

![](_page_24_Picture_1.jpeg)

## Facial Expression Transfer

![](_page_25_Figure_1.jpeg)

#### Face2Face

![](_page_26_Picture_1.jpeg)

![](_page_27_Picture_1.jpeg)

Source Actor

![](_page_27_Picture_3.jpeg)

**Reenacted Proxy** 

![](_page_27_Picture_5.jpeg)

#### **Reenacted Output**

![](_page_28_Picture_1.jpeg)

Source Actor

![](_page_28_Picture_3.jpeg)

**Reenacted Proxy** 

![](_page_28_Picture_5.jpeg)

#### **Reenacted Output**

![](_page_29_Picture_1.jpeg)

Source Actor

![](_page_29_Picture_3.jpeg)

**Reenacted Proxy** 

![](_page_29_Picture_5.jpeg)

#### **Reenacted Output**

![](_page_30_Picture_1.jpeg)

Source Actor

![](_page_30_Picture_3.jpeg)

**Reenacted Proxy** 

![](_page_30_Picture_5.jpeg)

#### **Reenacted Output**

#### Analysis-by-Synthesis

Parametric model needs to be flexible -> there needs to be a P that re-creates captured RGB input

Optimizable -> Must be able to find good optimum in energy landscape E(P)

Incompleteness

-> Image-based tricks to fix 3D artifacts are unsatisfactory

Over-parameterized models -> can re-create input

#### Over-parameterized models -> can re-create input

![](_page_34_Picture_2.jpeg)

Generator loss

 $J^{(G)} = -J^{(D)}$ 

GANs [Goodfellow et al. 14], Pix2Pix [Isola et al. 17], ProGAN [Karras et al. 18], ...

#### Over-parameterized models -> can re-create input

![](_page_35_Picture_2.jpeg)

No explicit no control -> struggle with videos

Discriminator loss  $J^{(D)} = -\frac{1}{2} \mathbb{E}_{\boldsymbol{x} \sim p_{\text{data}}} \log D(\boldsymbol{x}) - \frac{1}{2} \mathbb{E}_{\boldsymbol{z}} \log (1 - D(G(\boldsymbol{z})))$ Generator loss  $J^{(G)} = -J^{(D)}$ 

GANs [Goodfellow et al. 14], Pix2Pix [Isola et al. 17], ProGAN [Karras et al. 18], ...
### Conditional GANs



**Conditioning Input** 



Fully Controlled Output Video



### Conditional GANs



#### Conditioning Input



Fully Controlled Output Video



# Conditioning on Face Reconstruction



#### Source Sequence

Conditioning Images

Result

#### Neural Network converts synthetic data to realistic video

# Conditioning on Face Reconstruction



# Conditioning on Face Reconstruction



# Video Editing



- cGANs work with different input
- Requires consistent input
  i.e., accurate tracking

- cGANs work with different input
- Requires consistent input
  i.e., accurate tracking



[Chan et al. 18] Everybody Dance Now

- cGANs work with different input
- Requires consistent input
  i.e., accurate tracking



[Chan et al. 18] Everybody Dance Now

- cGANs work with different input
- Requires consistent input
  i.e., accurate tracking
- Network has no explicit 3D notion



[Chan et al. 18] Everybody Dance Now

# Videos still challenging for cGANs...



#### DeepVoxels: Explicit 3D Features



Simplified overview for novel view synthesis

CVPR'19 (Oral) [Sitzmann et al.]: DeepVoxels

#### DeepVoxels: Explicit 3D Features





#### **3D Geometry**



**Neural Texture** 





#### Novel View-Point Synthesis





#### Novel View-Point Synthesis





# **Geometry** Editing

Sequenc

Input











# Scene Editing





# Scene Editing





#### **Animation Synthesis**



#### **Animation Synthesis**



#### **Animation Synthesis**



#### **Animation Synthesis**



#### **Animation Synthesis**



# Deferred Neural Rendering

#### **Animation Synthesis**



# Deferred Neural Rendering

#### **Animation Synthesis**



# Conditioning on Audio: Neural Voice Puppetry



#### Audio to Video

#### **German News Video**



#### **English Audio**











Person-specific Blendshape Expression Model



#### Audio2Expression Training
#### Neural Voice Puppetry



[Thies et al. 19]: Neural Voice Puppetry

#### Neural Voice Puppetry: Audio to Video





[Thies et al. 19]: Neural Voice Puppetry

## Many Real-World Applications

#### Synthesia: Lip Sync



## synthesia

https://www.synthesia.io/







Victor Riparbelli CEO, Co-founder

Prof. Matthias Niessner Co-founder

Prof. Lourdes Agapito Co-founder



Steffen Tjerrild COO/CFO, Co-founder

Jason Lovell

VP Global Partnerships

Dr. Jonathan Starck CTO

Qi Liu Yin

**Research Engineer** 





Lead Al-Researcher

Dr. Karel Lebeda **Research Engineer** 



Research Engineer





Dr. Corneliu Ilisescu Research Engineer





#### Synthesia: Lip Sync



https://www.synthesia.io/

#### Synthesia: Lip Sync



https://www.synthesia.io/



Synthesia Dubbing

https://www.malariamustdie.com/



Synthesia Dubbing

https://www.malariamustdie.com/

### My Virtual Avatar



https://www.synthesia.io/

### My Virtual Avatar



https://www.synthesia.io/

#### Video Editing is Popular

#### Video Editing is Popular



## Video Editing is very Popular





IMU2Face



IMU2Face



Voice

NeuralVoicePuppetry



Voice

NeuralVoicePuppetry

Need to think about ethics and possible counter measures!

#### Study with over 200 participants



#### Al for Detection: FaceForensics



#### FaceForensics: Dataset

#### Source: 1,000 Videos (510,529 frames)

Methods	Train	Validation	Test
Pristine	366,847	68,511	73,770
DeepFakes	366,835	68,506	73,768
Face2Face	366,843	68,511	73,770
FaceSwap	291,434	54,618	59,640
NeuralTextures	291,834	54,630	59,672



- Publicly available!

- Over 2 million manipulated frames
- Three compression levels for each manipulated frame
- Over 1500 research groups

#### FaceForensics: Deep Detection Dataset

- Over 3000 manipulated videos
- From 28 actors
- Variety of scenes
- Provided by Google & JigSaw



#### FaceForensics: Supervised Detection











#### Unsupervised / Self-Supervised Forensics

Major challenges

- Self-supervised Learning
- Transfer Learning
- Unsupervised Learning



[Cozzolino et al. 19]: ForensicTransfer

#### Conclusion



#### Professor

#### PostDocs



Prof. Dr. Matthias Nießner **PhD Students** 



Dr. Justus Thies



Angel X. Chang Visiting Professor: Hans Fischer Fellow **Princeton University** 



Leonidas Guibas





Aljaž Božič



Norman Müller



Andreas Rössler



Yawar Siddiqui



Armen Avetisyan



Ji Hou



Dejan Azinović





Manuel Dahnert

Dave Zhenyu Chen





# Thank You!





# Thank You!



