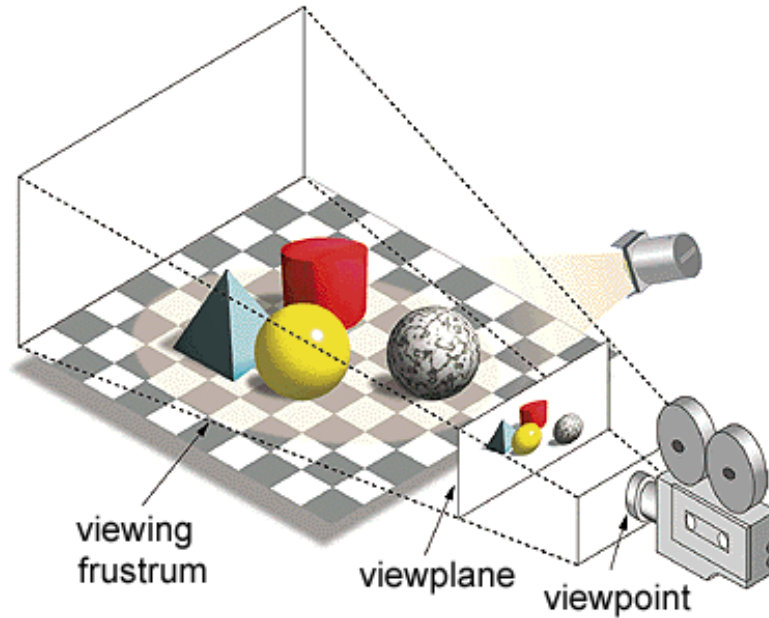


Neural Rendering

Rendering

From Computer Desktop Encyclopedia
Reprinted with permission.
© 1998 Intergraph Computer Systems

- 3D Scene:
- Material
 - Lighting
 - Geometry
(incl. animation)



Camera Def.

- Intrinsic
- Often:
 - focal length
 - principal point)

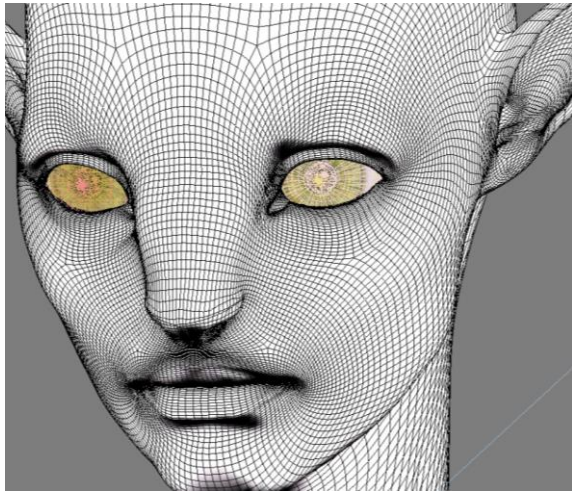
Camera View Point

- Extrinsic
- 6 DoF (3rot, 3trans)

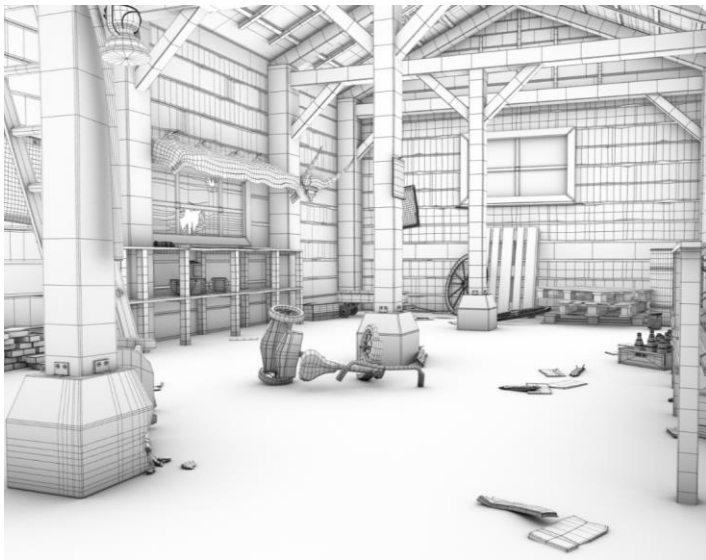
Photo-realistic Image Synthesis

The Rendering Equation [Kajiya 86]

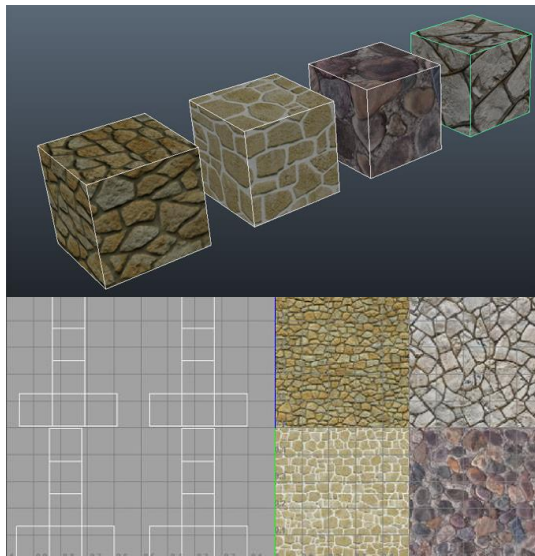
$$L_o(\mathbf{x}, \omega_o, \lambda, t) = L_e(\mathbf{x}, \omega_o, \lambda, t) + \int_{\Omega} f_r(\mathbf{x}, \omega_i, \omega_o, \lambda, t) L_i(\mathbf{x}, \omega_i, \lambda, t) (\omega_i \cdot \mathbf{n}) d\omega_i$$



Need 3D Content for Rendering



Geometry

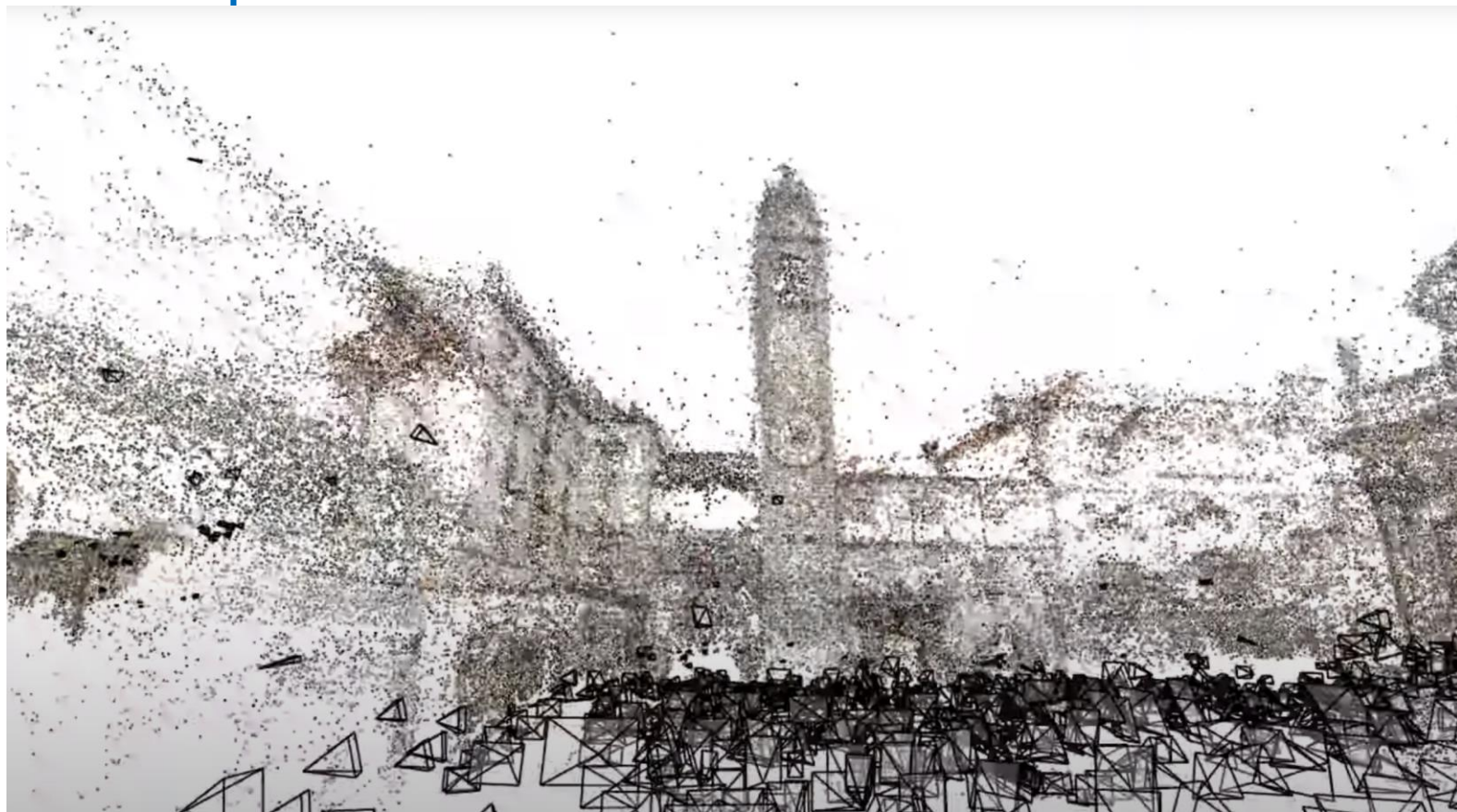


Textures



Material & Lighting

Computer Vision for Reconstruction



Prof. Leal-Taixé and Prof. Niessner

ICCV'09 [Agarwal et al.]: Building Rome in a Day

3D Digitization



Computer Graphics



Computer Vision

Traditional Graphics vs Deep Learning



3D Model + Textures + Shading -> Synthetic Image



Generative Adversarial Networks



Discriminator loss

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

Generator loss

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

Idea of Neural Rendering

Novel View point synthesis:

6 DoF Camera
Pose / View Point



Neural Network
-> Encodes entire
scene description,
lighting, materials,
etc.



Neural Rendering with Pix2Pix

Ground truth for training

- Pose + Target Image (e.g., observed from real world)
- Constrain with re-rendering loss

Testing

- Given unseen pose, generate image

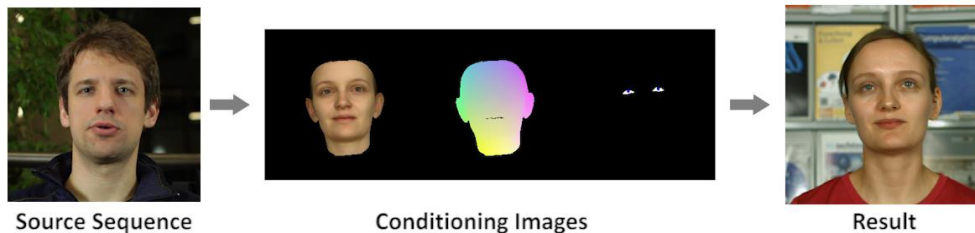
Neural Rendering with Pix2Pix

Pix2Pix [Isola et al. 2017]

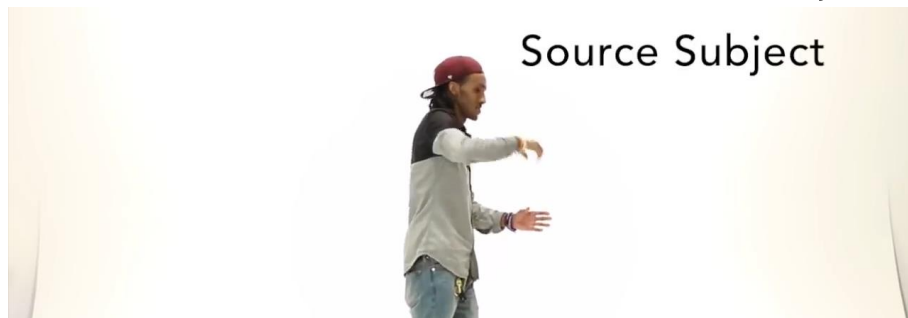


Other Neural Rendering

- Conditioned on Faces (Deep Video Portraits)



- Conditioned on Human Skeleton (Everybody Dance Now)



Neural Rendering with Pix2Pix

Pix2Pix [Isola et al. 2017]

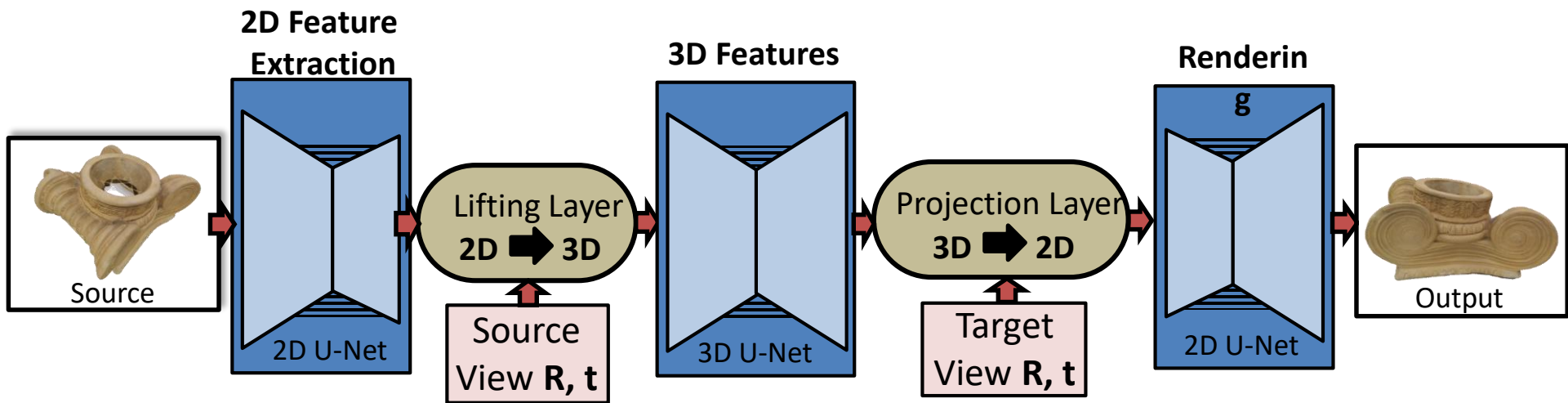


Deep Voxels

Deep Voxels

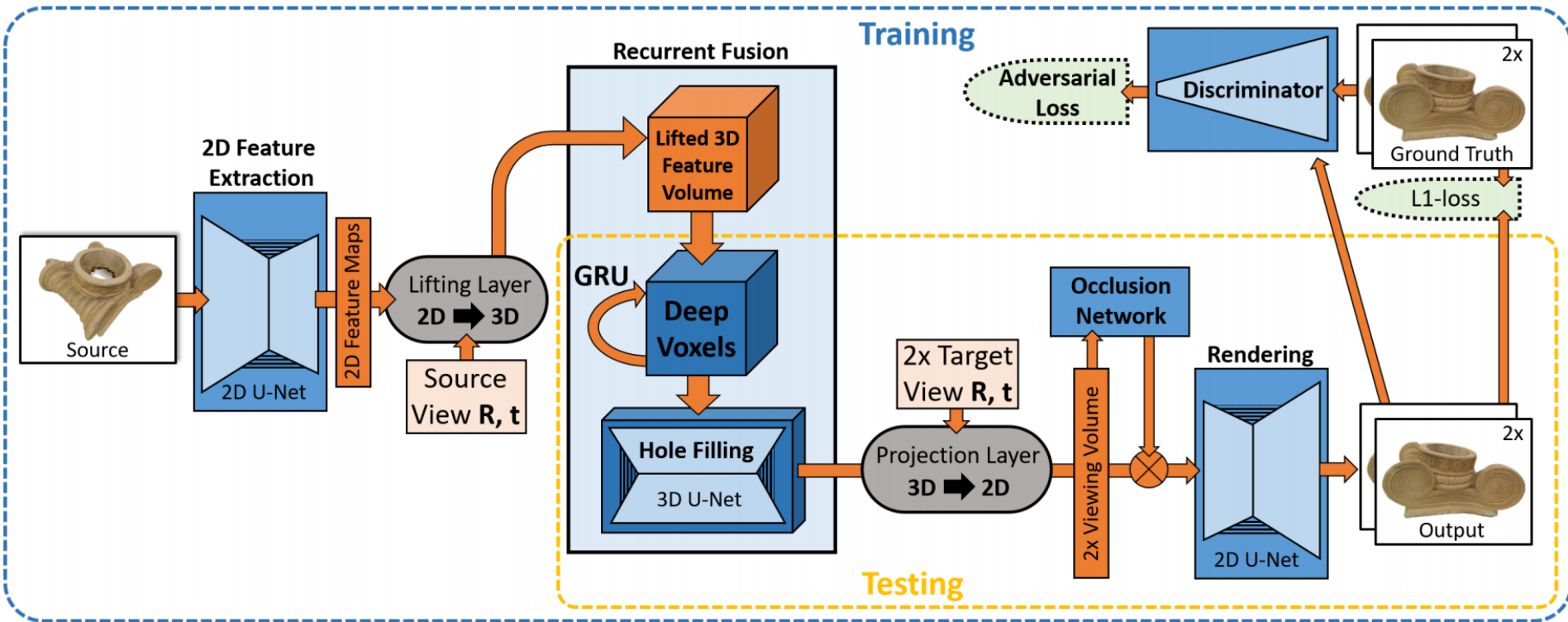
- Main idea for video generation:
 - Why learn 3D operations with 2D Convs !?!?
 - We know how 3D transformations work
 - E.g., 6 DoF rigid pose $[R | t]$
 - Incorporate these into the architectures
 - Need to be differentiable!
 - Example application: novel view point synthesis
 - Given rigid pose, generate image for that view

Deep Voxels



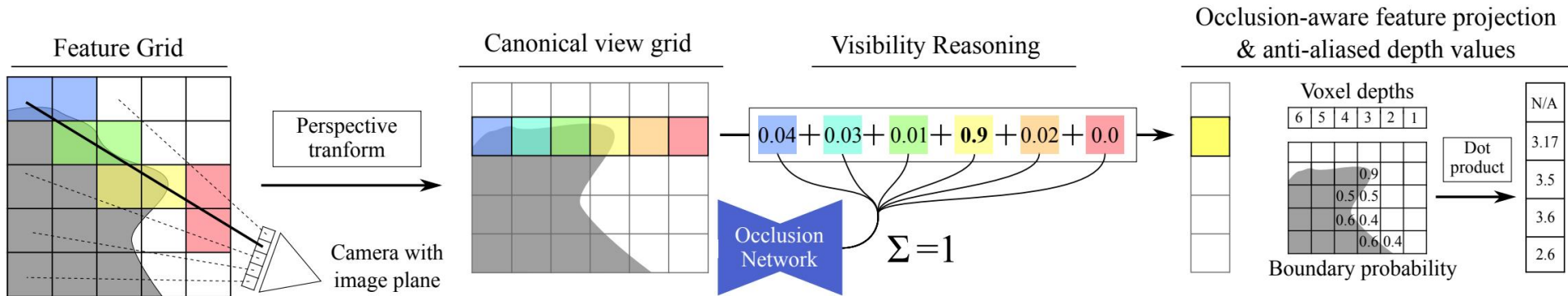
Simplified overview for novel view synthesis

Deep Voxels



Deep Voxels

Occlusion Network:



Issue: we don't know the depth for the target!

- > Per-pixel softmax along the ray
- > Network learns the depth

Deep Voxels

DeepVoxels

ABCDEFGHIJKLM
NOPQRSTUVWXYZ
ABCDEFGHIJKLM
NOPQRSTUVWXYZ
ABCDEFGHIJKLM
NOPQRSTUVWXYZ
ABCDEFGHIJKLM
NOPQRSTUVWXYZ
ABCDEFGHIJKLM



Best Baseline: Pix2Pix [Isola et al. 2017]

ABCDEFGHIJKLM
NOPQRSTUVWXYZ
ABCDEFGHIJKLM
NOPQRSTUVWXYZ
ABCDEFGHIJKLM
NOPQRSTUVWXYZ
ABCDEFGHIJKLM
NOPQRSTUVWXYZ
ABCDEFGHIJKLM

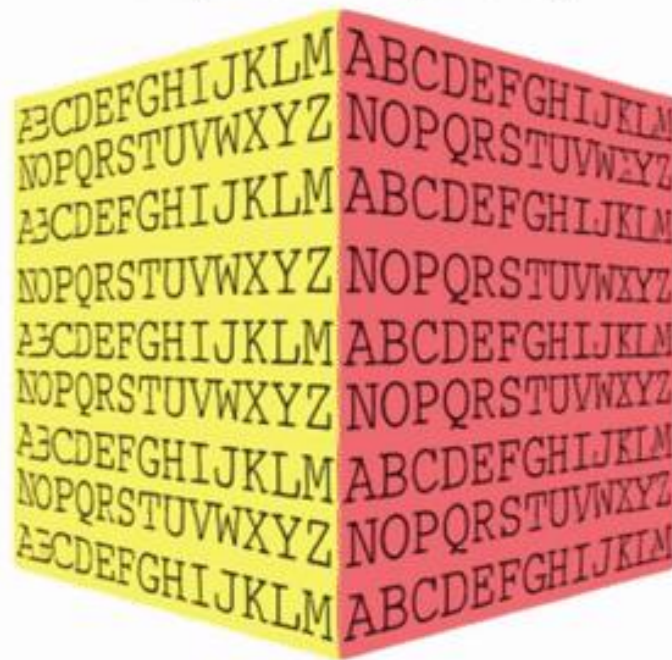


Deep Voxels

Pix2Pix [Isola et al. 2017]



DeepVoxels (Ours)

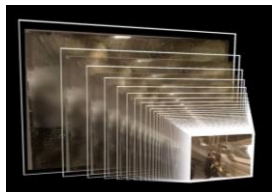


Deep Voxels: Insights

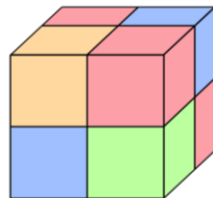
- Lifting from 2D to 3D works great
 - No need to take specific care for temp. coherency!
- All 3D operations are differentiable
- Currently, only for novel view-point synthesis
 - I.e., cGAN for new pose in a given scene
- But: limited resolution due to dense 3D voxel grid

Importing 3D structure from CG

Scene
Representa
tion



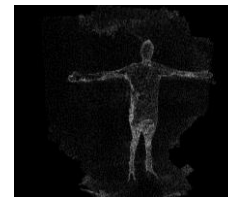
Multi-Plane Images



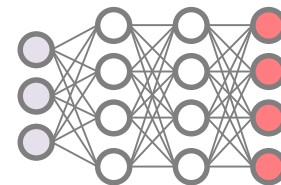
Voxelgrids



Image-based



Point Clouds



Implicit Function

Renderer

(Alpha) compositing

Volumetric
Ray-based

Rasterization

Splatting

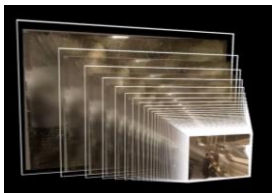
Sphere-Traced
Volumetric

Scene Representation

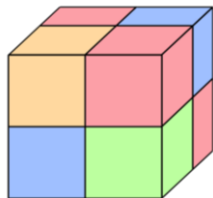
Differentiable Renderer

Importing 3D structure from CG

Scene
Representa
tion



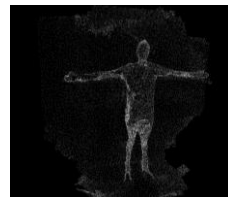
Multi-Plane Images



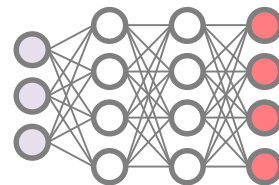
Voxelgrids



Image-based



Point Clouds



Implicit Function

Renderer

(Alpha) compositing

Volumetric
Ray-based

Rasterization

Splatting

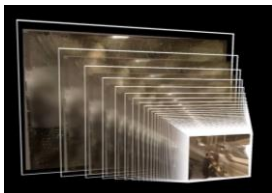
Sphere-Traced
Volumetric

Pros

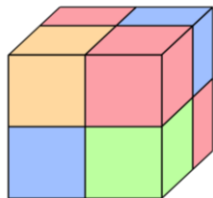
Cons

Importing 3D structure from CG

Scene
Representa
tion



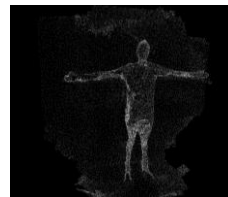
Multi-Plane Images



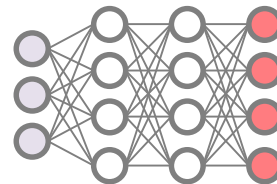
Voxelgrids



Image-based



Point Clouds



Implicit Function

Renderer

(Alpha) compositing

Volumetric
Ray-based

Rasterization

Splatting

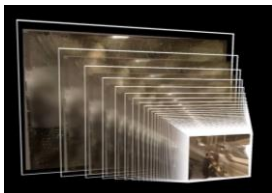
Sphere-Traced
Volumetric

Pros

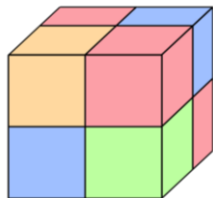
Cons

Importing 3D structure from CG

Scene
Representa
tion



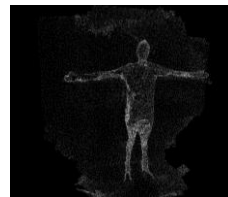
Multi-Plane Images



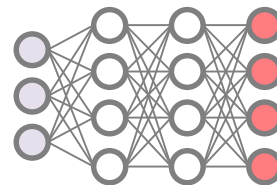
Voxelgrids



Image-based



Point Clouds



Implicit Function

Renderer

(Alpha) compositing

Volumetric
Ray-based

Rasterization

Splatting

Sphere-Traced
Volumetric

Pros

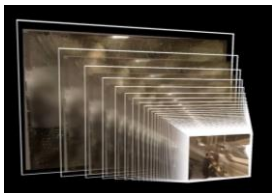
Fast rendering
High quality
Generalizes

Cons

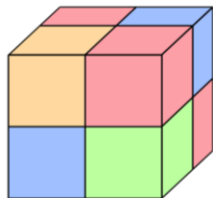
Only 2.5D
Size

Importing 3D structure from CG

Scene
Representa
tion



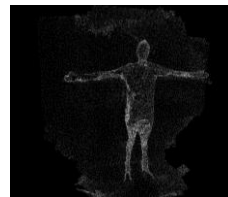
Multi-Plane Images



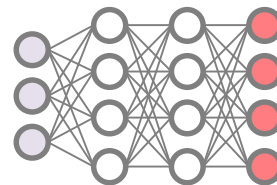
Voxelgrids



Image-based



Point Clouds



Implicit Function

Renderer

(Alpha) compositing

Volumetric
Ray-based

Rasterization

Splatting

Sphere-Traced
Volumetric

Pros

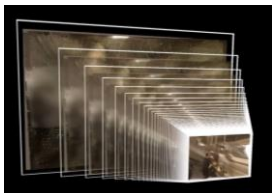
Fast rendering
High quality
Generalizes

Cons

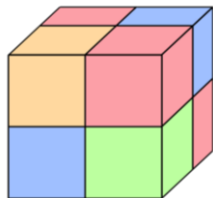
Only 2.5D
Size

Importing 3D structure from CG

Scene
Representa
tion



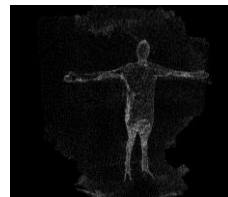
Multi-Plane Images



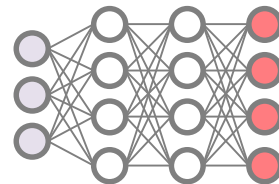
Voxelgrids



Image-based



Point Clouds



Implicit Function

Renderer

(Alpha) compositing

Volumetric
Ray-based

Rasterization

Splatting

Sphere-Traced
Volumetric

Pros

Fast rendering
High quality
Generalizes

“True 3D”
High quality

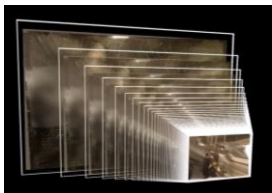
Cons

Only 2.5D
Size

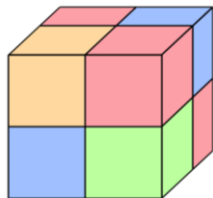
No reconstruction
priors
Memory $O(n^3)$

Importing 3D structure from CG

Scene Representation



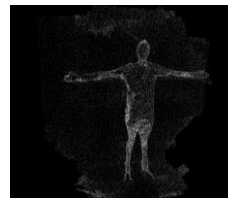
Multi-Plane Images



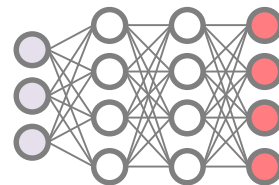
Voxelgrids



Image-based



Point Clouds



Implicit Function

Renderer

(Alpha) compositing

Volumetric Ray-based

Rasterization

Splatting

Sphere-Traced Volumetric

Pros

Fast rendering
High quality
Generalizes

"True 3D"
High quality

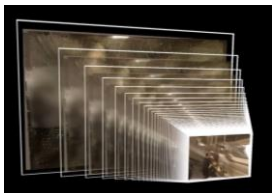
Cons

Only 2.5D
Size

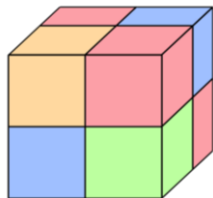
No reconstruction priors
Memory $O(n^3)$

Importing 3D structure from CG

Scene
Representa
tion



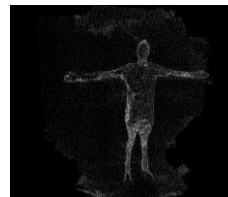
Multi-Plane Images



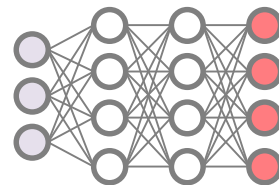
Voxelgrids



Image-based



Point Clouds



Implicit Function

Renderer

(Alpha) compositing

Volumetric
Ray-based

Rasterization

Splatting

Sphere-Traced
Volumetric

Pros

Fast rendering
High quality
Generalizes

“True 3D”
High quality

High quality

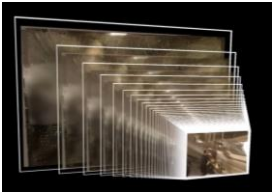
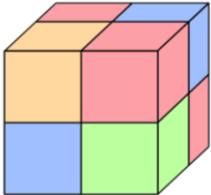

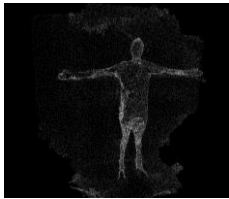
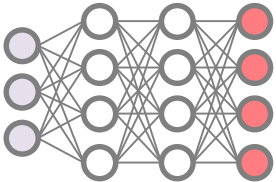
Cons

Only 2.5D
Size

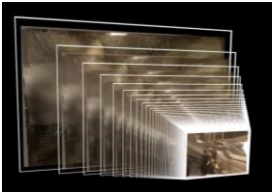
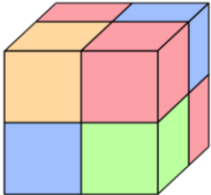

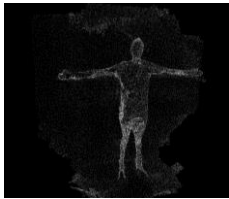
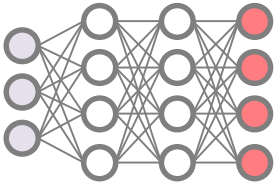
No reconstruction
priors
Memory $O(n^3)$

Requires good SFM
No compact
representation

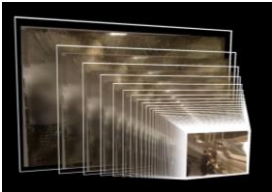
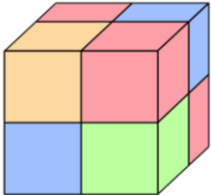

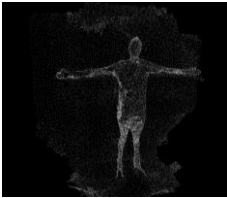
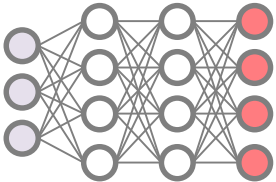
Importing 3D structure from CG

<p>Scene Representation</p>					
<p>Renderer</p>	<p>(Alpha) compositing</p>	<p>Volumetric Ray-based</p>	<p>Rasterization</p>	<p>Splatting</p>	<p>Sphere-Traced Volumetric</p>
<p>Pros</p>	<p>Fast rendering High quality Generalizes</p>	<p>“True 3D” High quality</p>	<p>High quality</p>	<p></p>	<p></p>
<p>Cons</p>	<p>Only 2.5D Size</p>	<p>No reconstruction priors Memory $O(n^3)$</p>	<p>Requires good SFM No compact representation</p>	<p></p>	<p></p>

Importing 3D structure from CG

<p>Scene Representation</p>					
	<p>Multi-Plane Images</p>	<p>Voxelgrids</p>	<p>Image-based</p>	<p>Point Clouds</p>	<p>Implicit Function</p>
<p>Renderer</p>	<p>(Alpha) compositing</p>	<p>Volumetric Ray-based</p>	<p>Rasterization</p>	<p>Splatting</p>	<p>Sphere-Traced Volumetric</p>
<p>Pros</p>	<p>Fast rendering High quality Generalizes</p>	<p>“True 3D” High quality</p>	<p>High quality</p>	<p>High quality</p>	
<p>Cons</p>	<p>Only 2.5D Size</p>	<p>No reconstruction priors Memory $O(n^3)$</p>	<p>Requires good SFM No compact representation</p>	<p>Requires good SFM</p>	

Importing 3D structure from CG

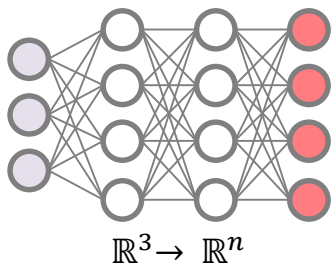
<p>Scene Representation</p>					
<p>Renderer</p>	<p>(Alpha) compositing</p>	<p>Volumetric Ray-based</p>	<p>Rasterization</p>	<p>Splatting</p>	<p>Sphere-Traced Volumetric</p>
<p>Pros</p>	<p>Fast rendering High quality Generalizes</p>	<p>“True 3D” High quality</p>	<p>High quality</p>	<p>High quality</p>	
<p>Cons</p>	<p>Only 2.5D Size</p>	<p>No reconstruction priors Memory $O(n^3)$</p>	<p>Requires good SFM No compact representation</p>	<p>Requires good SFM</p>	

Scene Representation Networks

Sitzmann et al., Neurips 2019

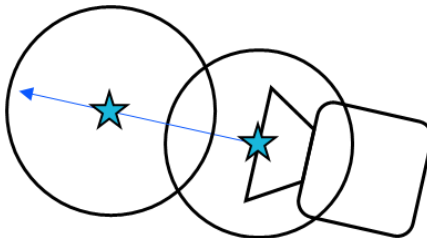
Scene
Representa
tion

ReLU MLP



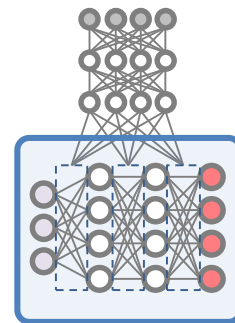
Renderer

Generalized (learned)
sphere-tracing



Generalizati
on

Hypernetwork



Scene Representation Networks

Sitzmann et al., Neurips 2019



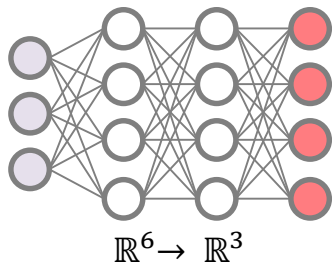
Full 3D Reconstruction from single image!

NERF: Neural Radiance Fields

Mildenhall et al., arXiv 2020

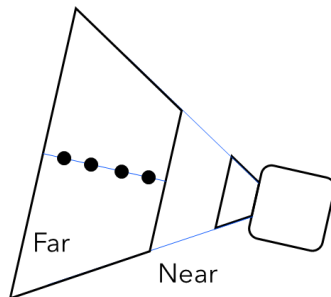
**Scene
Representa
tion**

ReLU MLP +
Positional Encoding
View Direction



Renderer

Volumetric,
stratified sampling

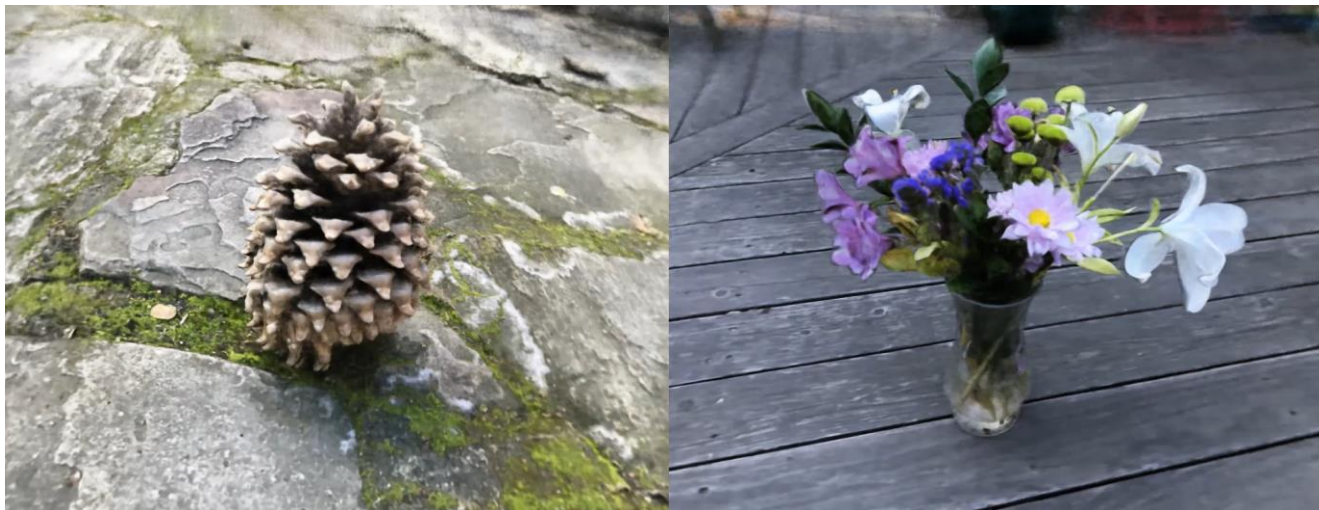


**Generalizati
on**

None.

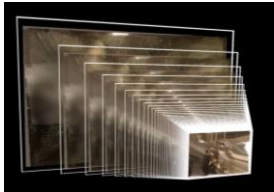
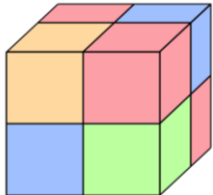

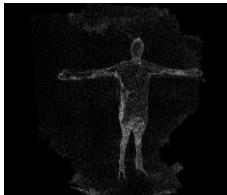
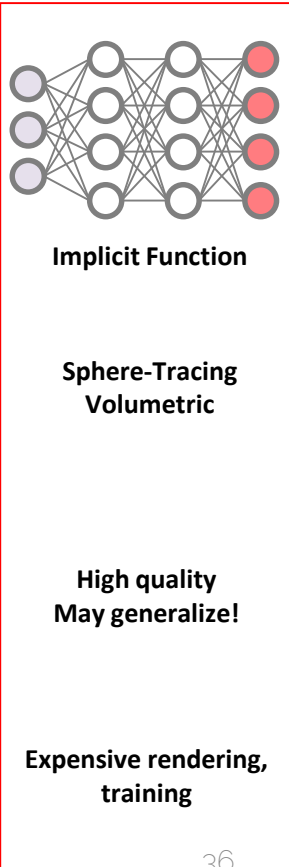
NERF: Neural Radiance Fields

Mildenhall et al., arXiv 2020



Photorealistic, including view-dependence!
(~100 images)

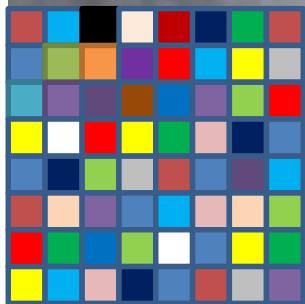
Requirements

Scene Representation					
Renderer	(Alpha) compositing	Volumetric Ray-based	Rasterization	Splatting	Sphere-Tracing Volumetric
Pros	Fast rendering High quality Generalizes	"True 3D" High quality	High quality	High quality	High quality May generalize!
Cons	Only 2.5D Size	No reconstruction priors Memory $O(n^3)$	Requires good SFM No compact representation	Requires good SFM	Expensive rendering, training

Neural Textures: Features on 3D Mesh

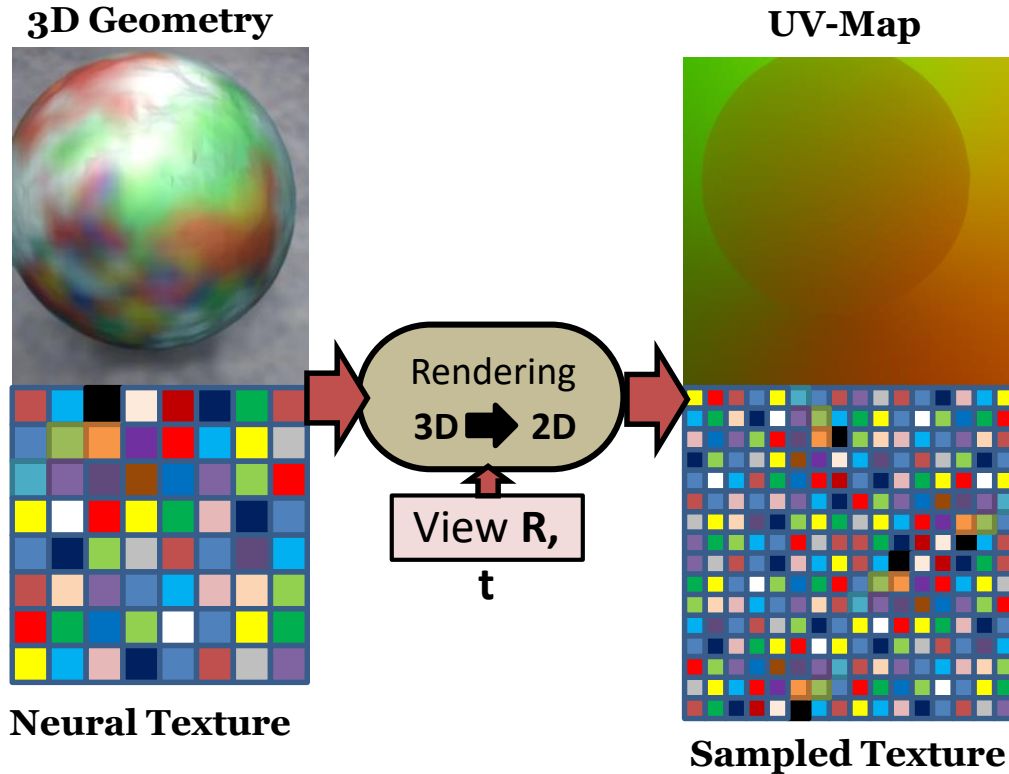
Neural Textures: Features on 3D Mesh

3D Geometry

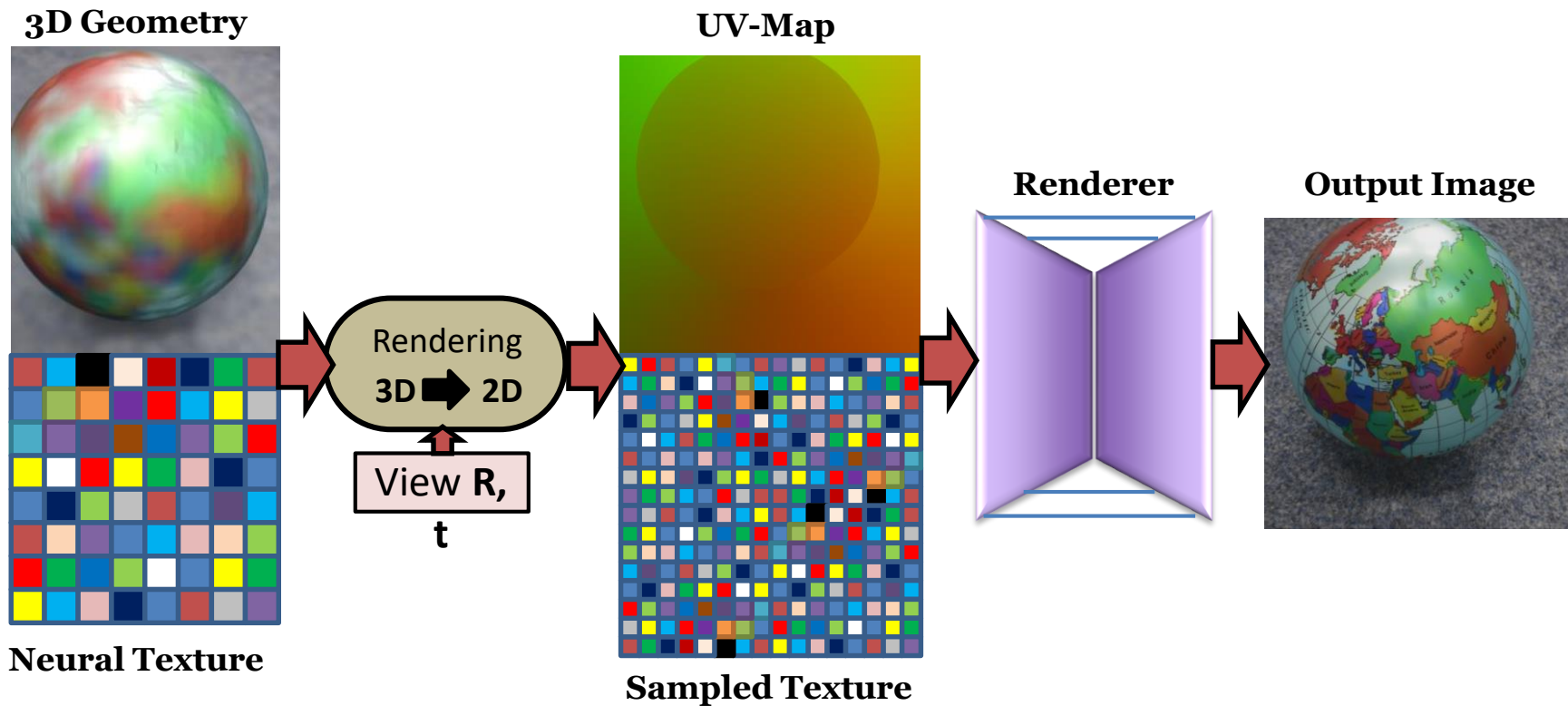


Neural Texture

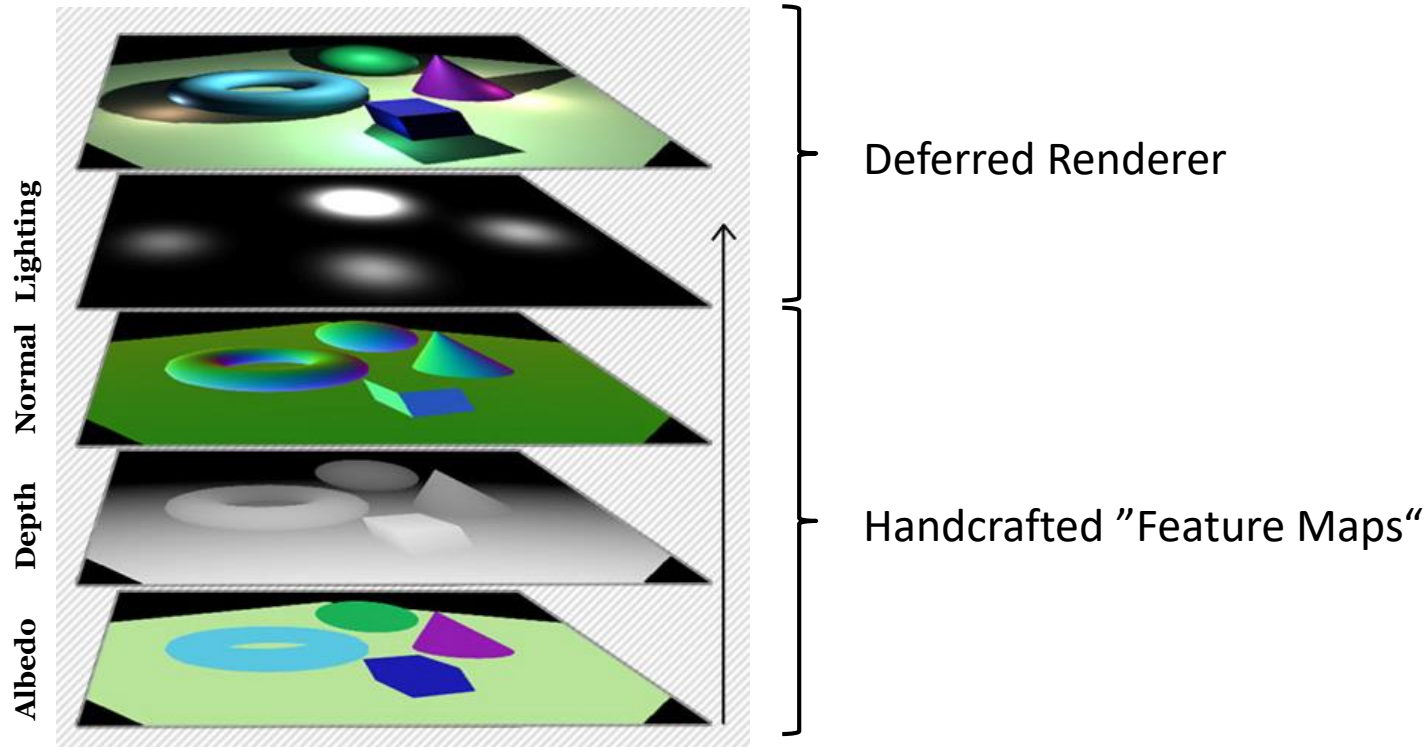
Neural Textures: Features on 3D Mesh



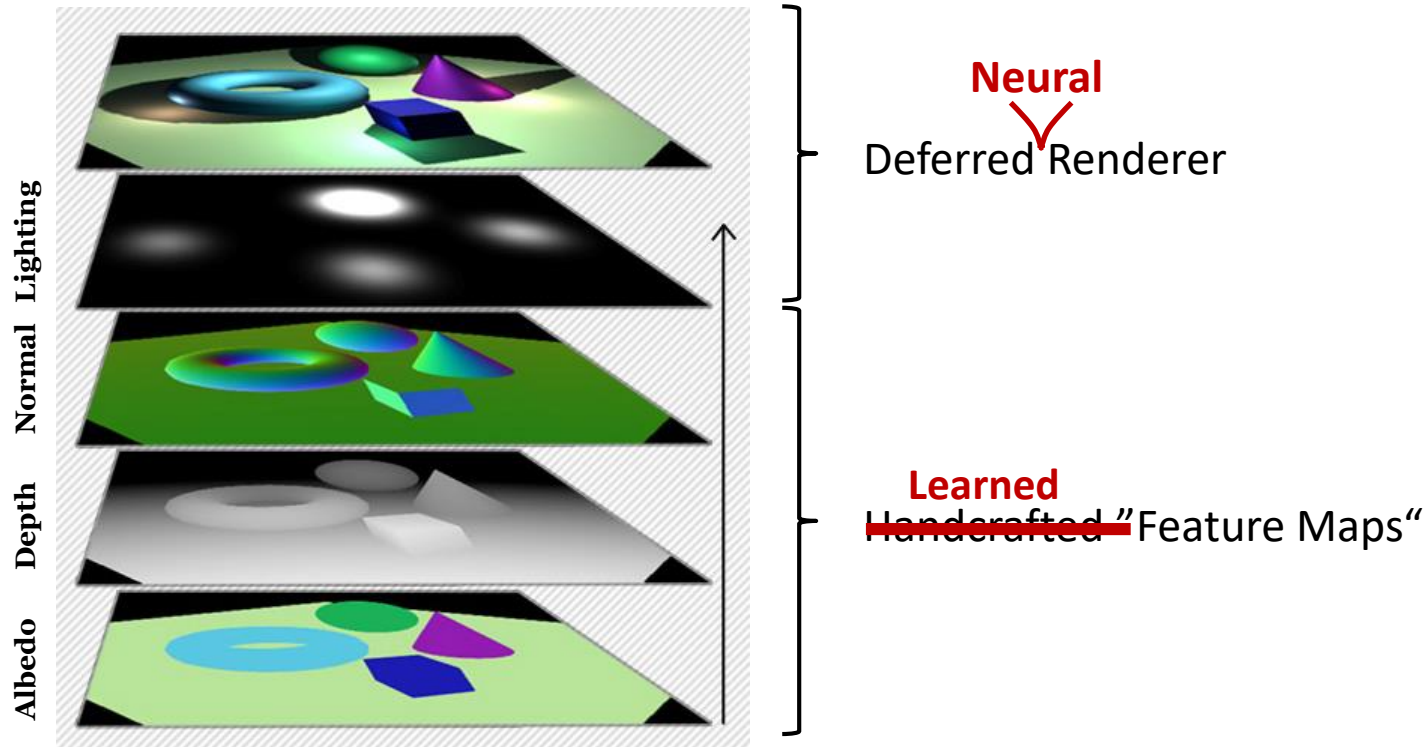
Neural Textures: Features on 3D Mesh



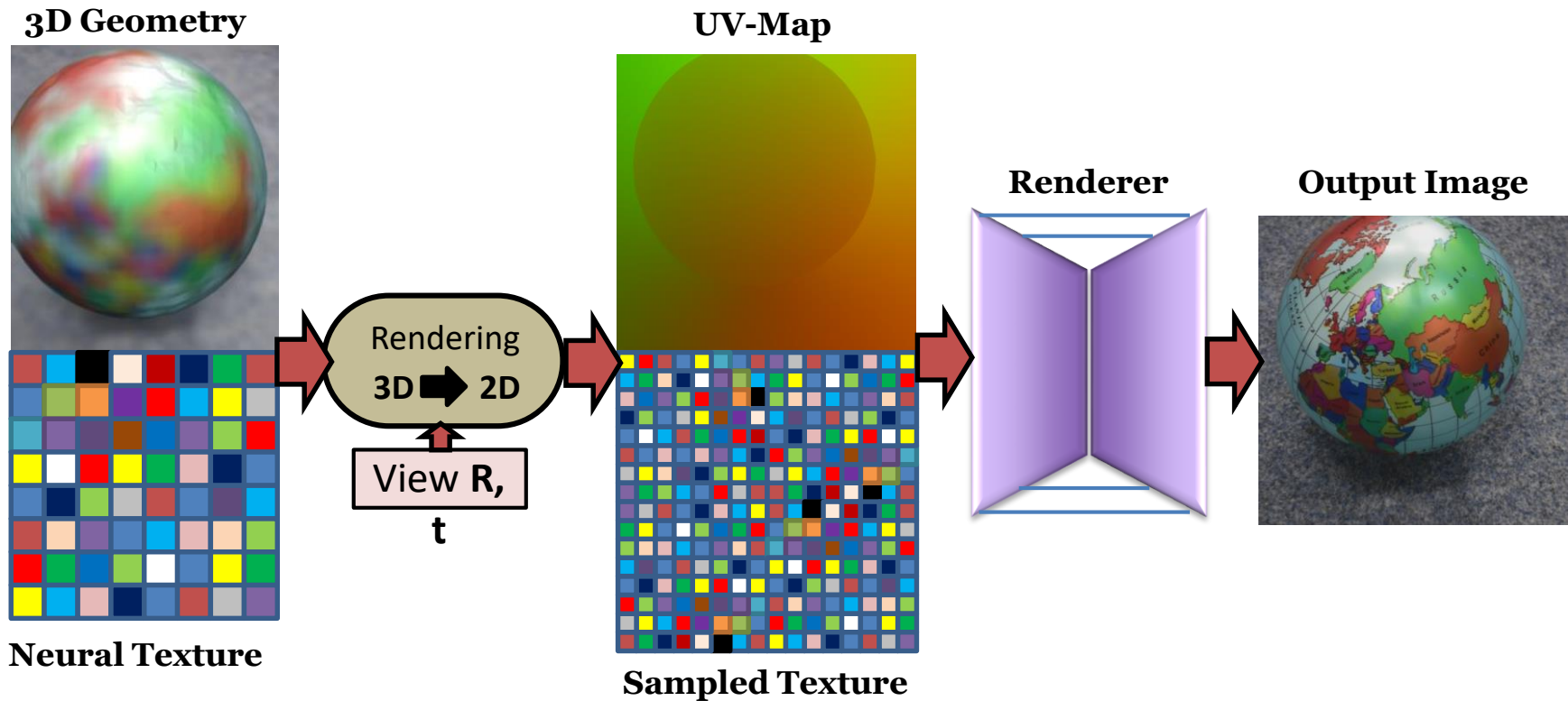
Deferred Neural Rendering



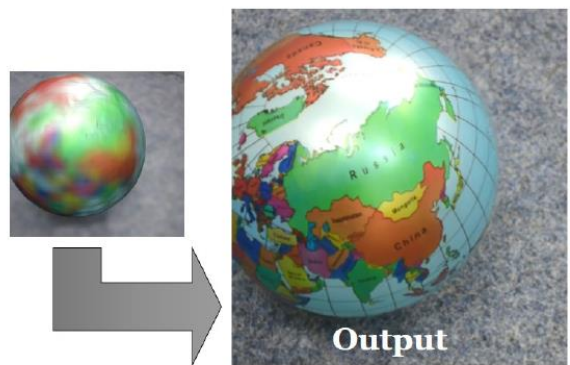
Deferred Neural Rendering



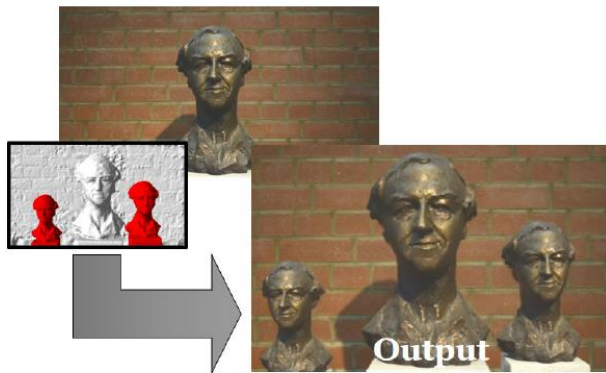
Deferred Neural Rendering



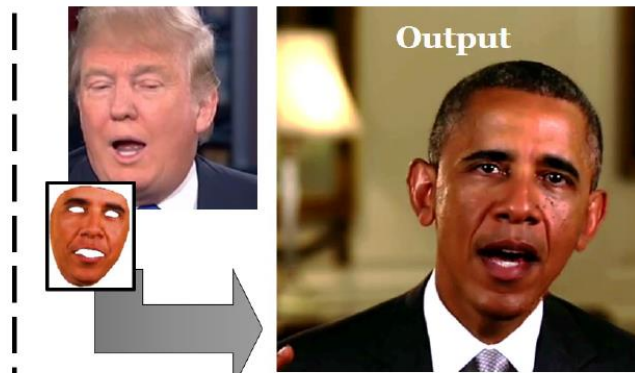
Neural Textures: Features on 3D Mesh



Novel View Synthesis

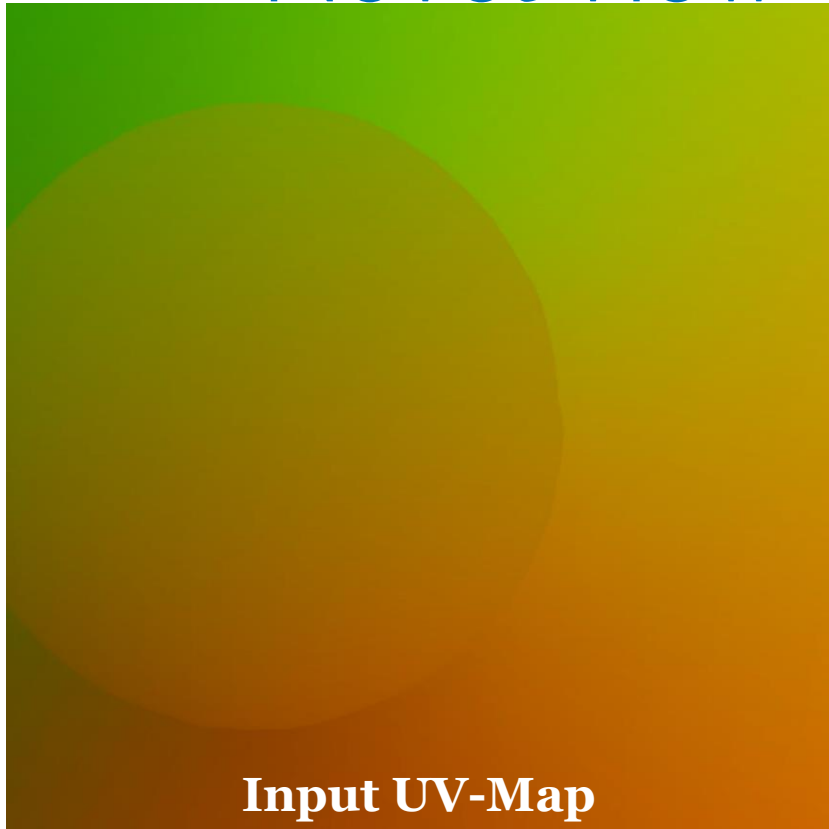


Scene Editing

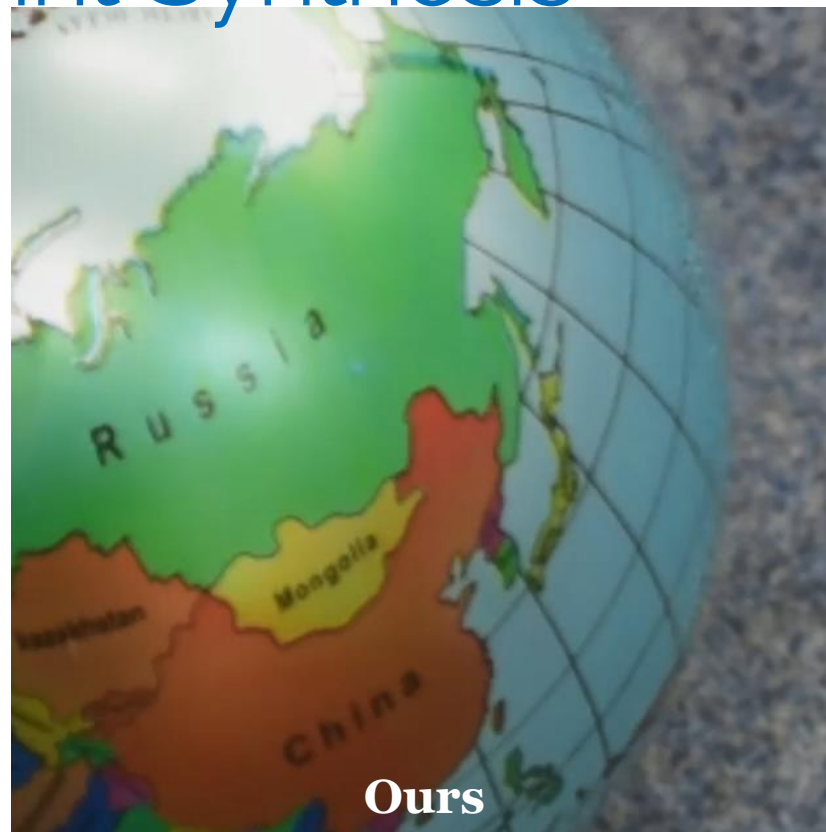
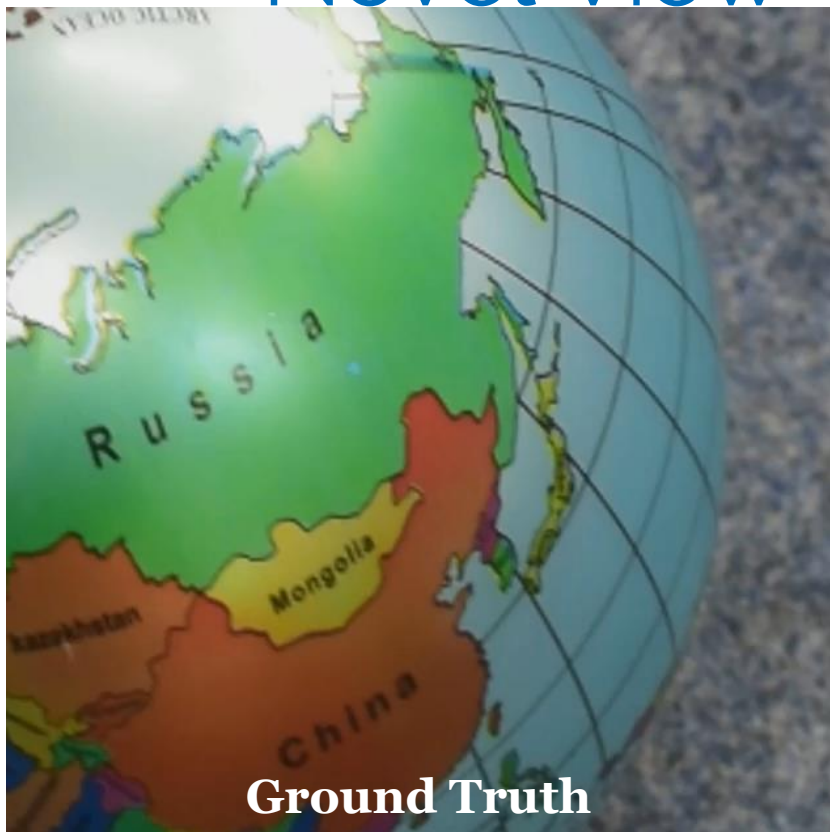


Animation Synthesis

Novel View-Point Synthesis



Novel View-Point Synthesis

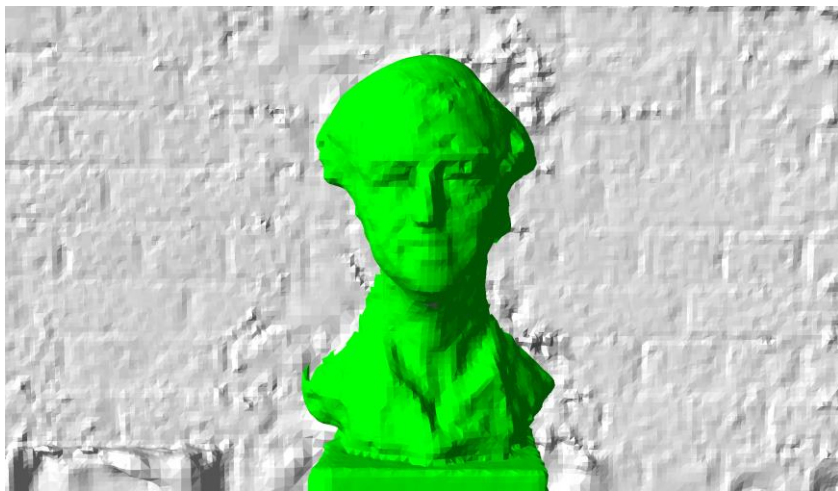


Scene Editing

**Input
Sequence**



**Geometry
Editing**



Scene Editing



Scene Editing



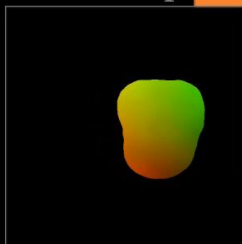
Facial Animation

Animation Synthesis

Source Actor



Target
UV-Map



Target
Background



Output



Facial Animation

Animation Synthesis

Source Actor



Output



Target
UV-Map

Target
Background



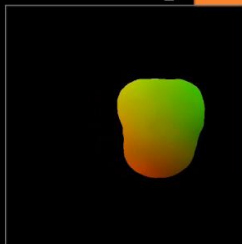
Facial Animation

Animation Synthesis

Source Actor



Target
UV-Map



Target
Background



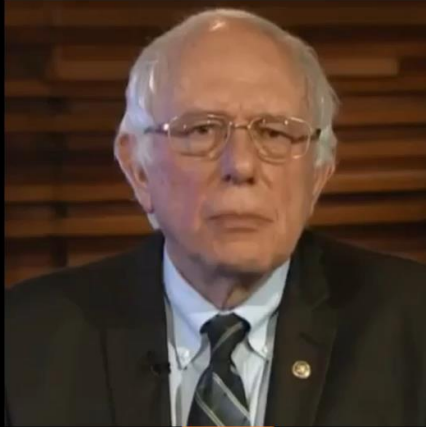
Output



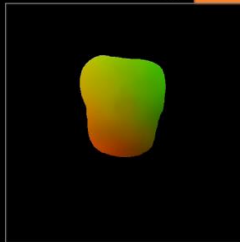
Facial Animation

Animation Synthesis

Source Actor



Target
UV-Map



Target
Background



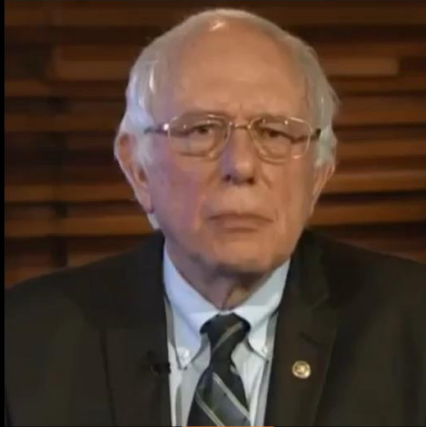
Output



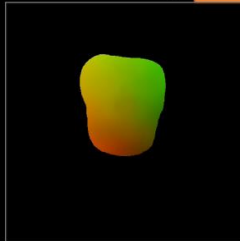
Facial Animation

Animation Synthesis

Source Actor



Target
UV-Map



Target
Background



Output



Deferred Neural Rendering

Animation Synthesis

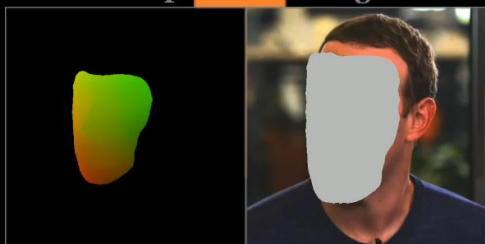
Source Actor



Target
UV-Map



Target
Background



Output



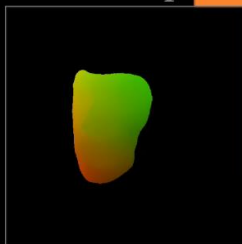
Deferred Neural Rendering

Animation Synthesis

Source Actor



Target
UV-Map



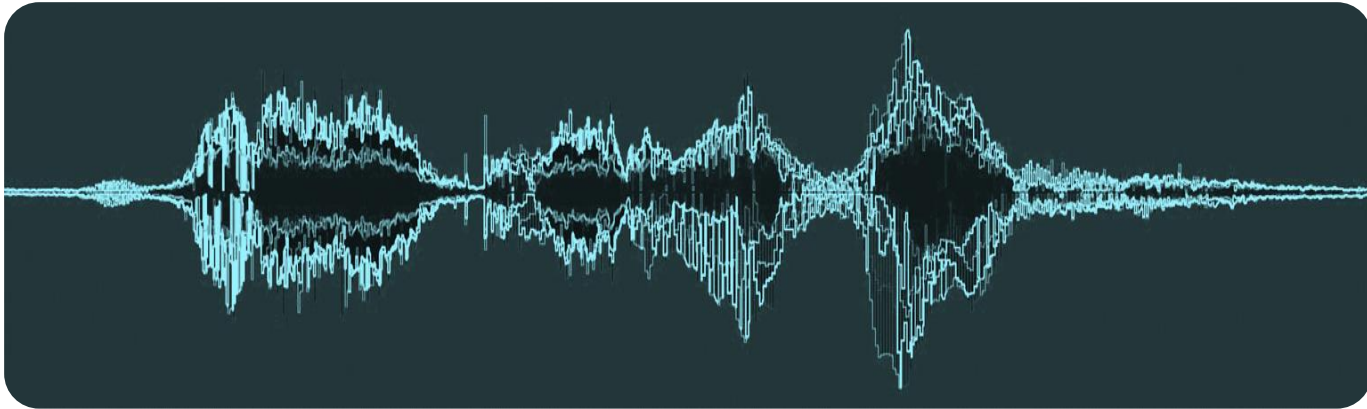
Target
Background



Output

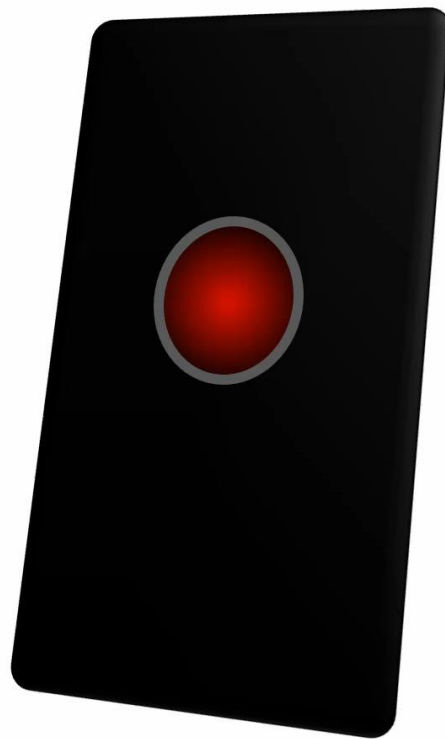


Neural Voice Puppetry



Neural Voice Puppetry

Hey Siri, can you
show me your face?

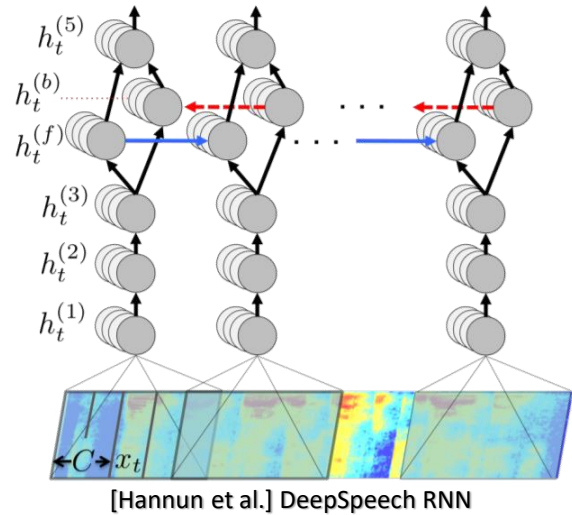


We use Siri as synonym for a digital assistant.

Neural Voice Puppetry

How does it work? **Pipeline Overview**

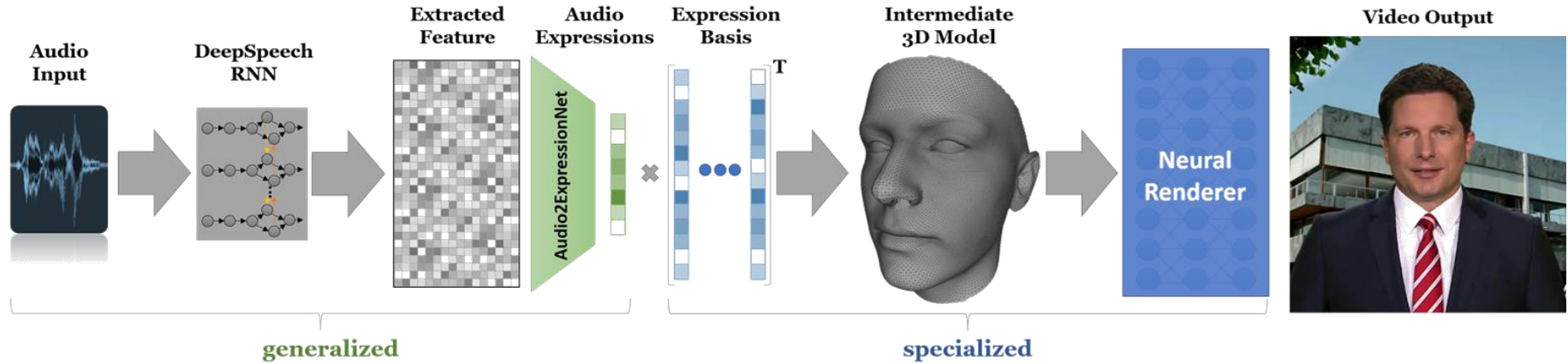
Neural Voice Puppetry



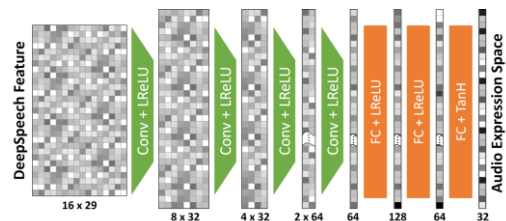
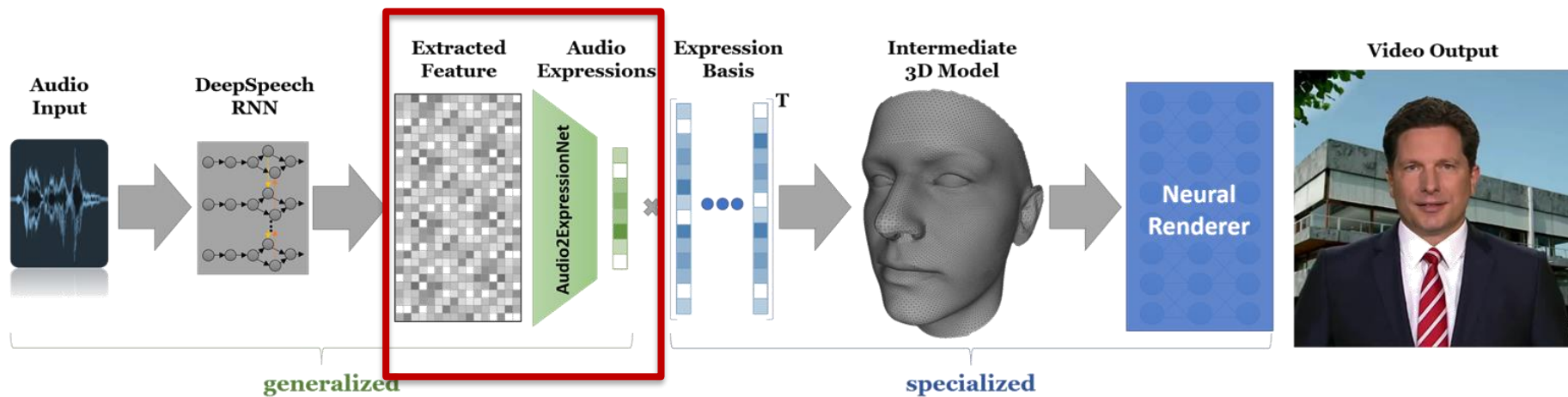
Output of the RNN of DeepSpeech:
- Logits of alphabet ($|\text{alphabet}|=29$)

We use a time window ($n=16$)

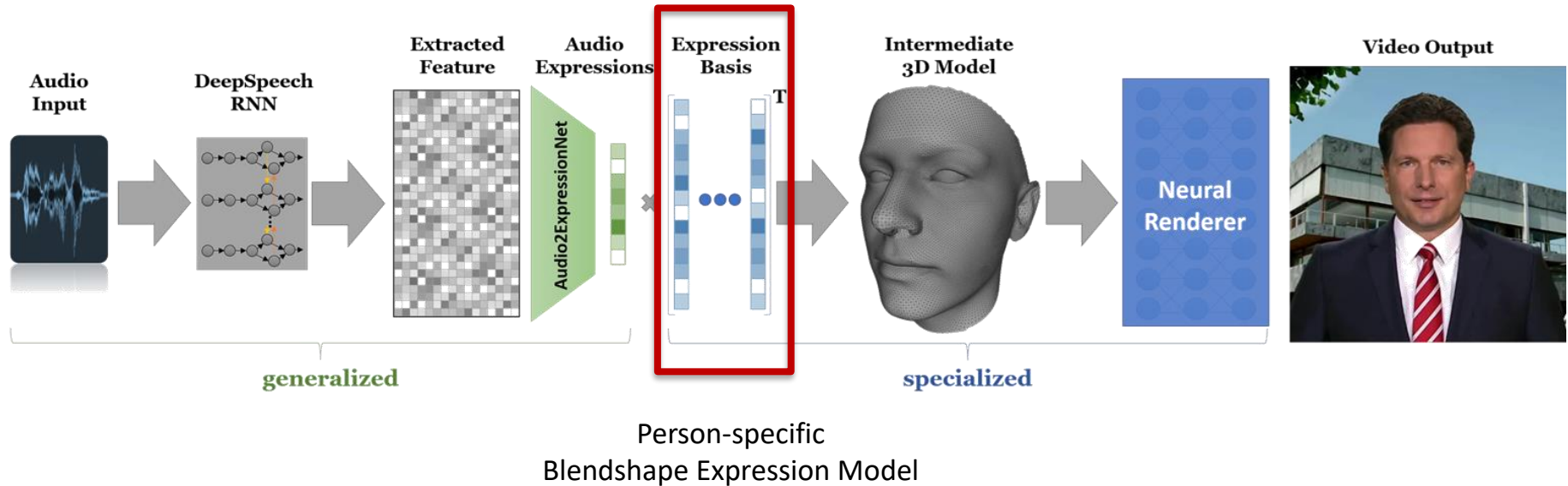
Neural Voice Puppetry



Neural Voice Puppetry

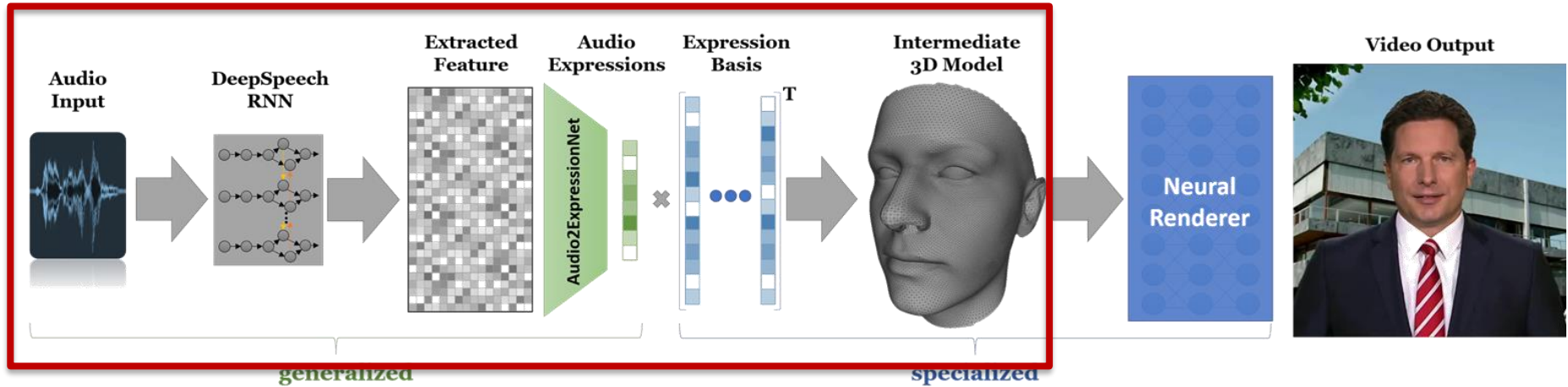


Neural Voice Puppetry

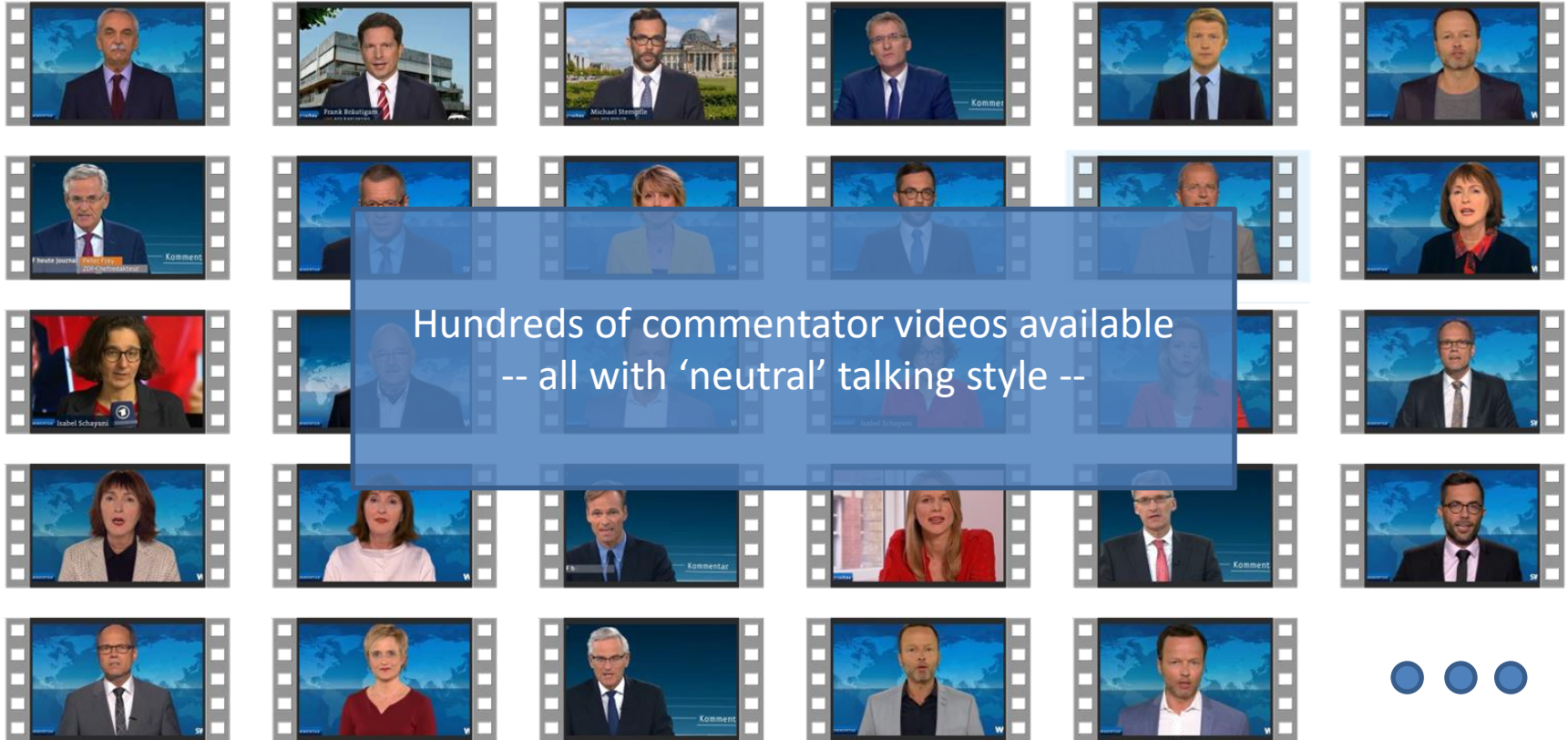


Neural Voice Puppetry

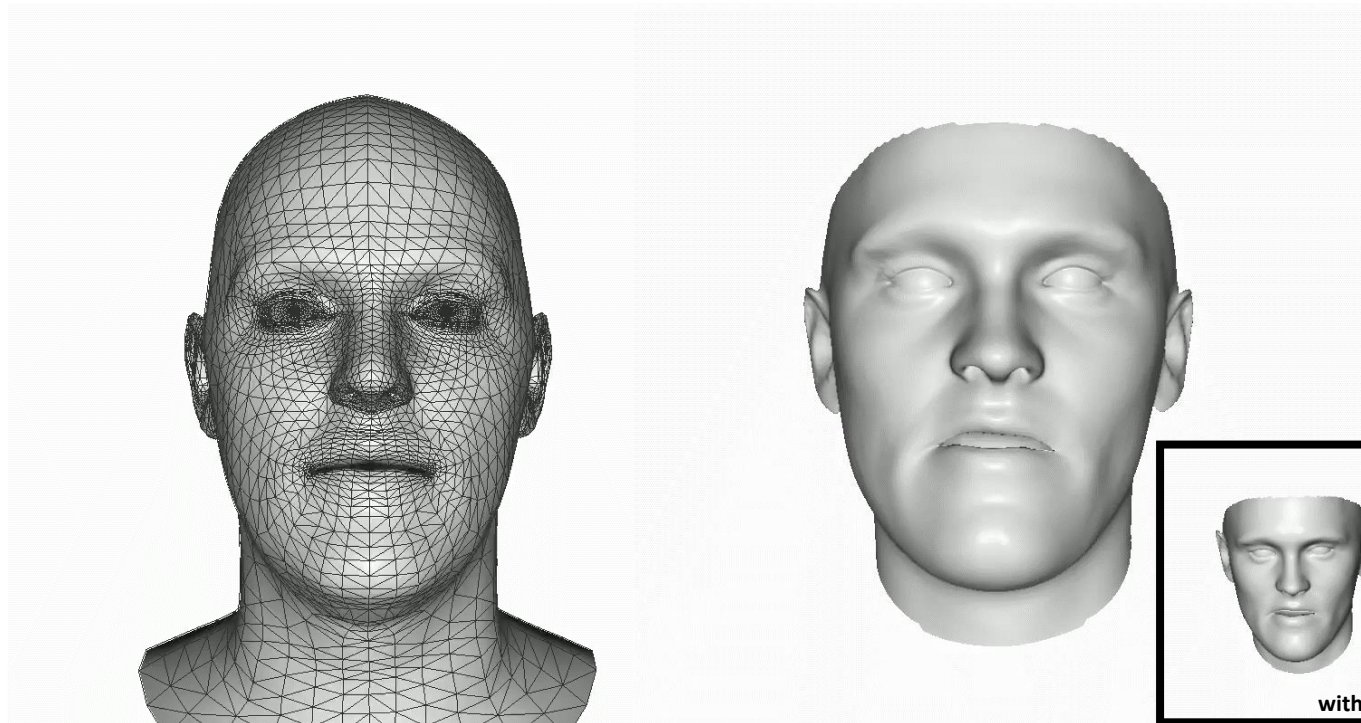
Audio2Expression Training



Neural Voice Puppetry



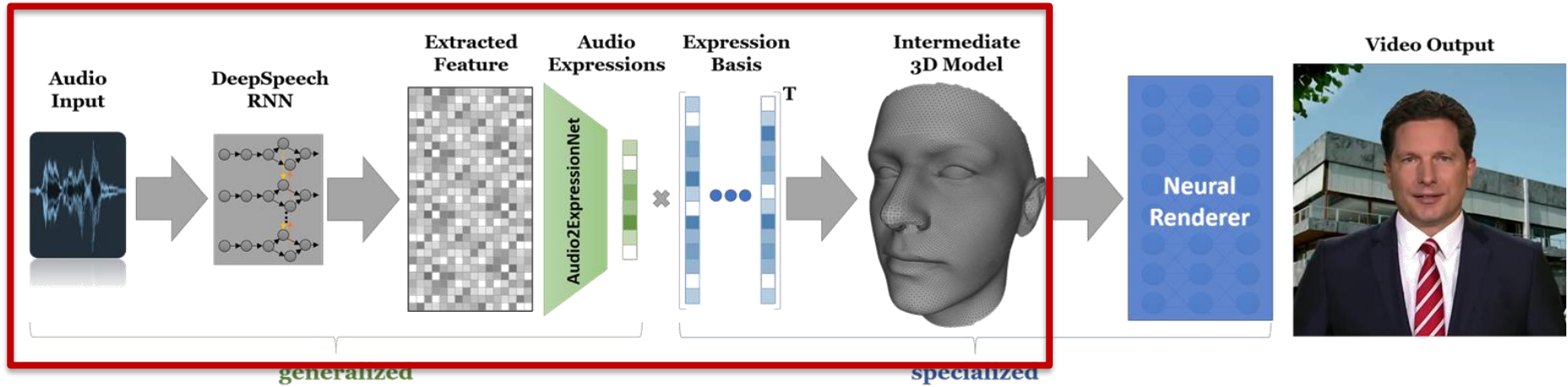
Neural Voice Puppetry



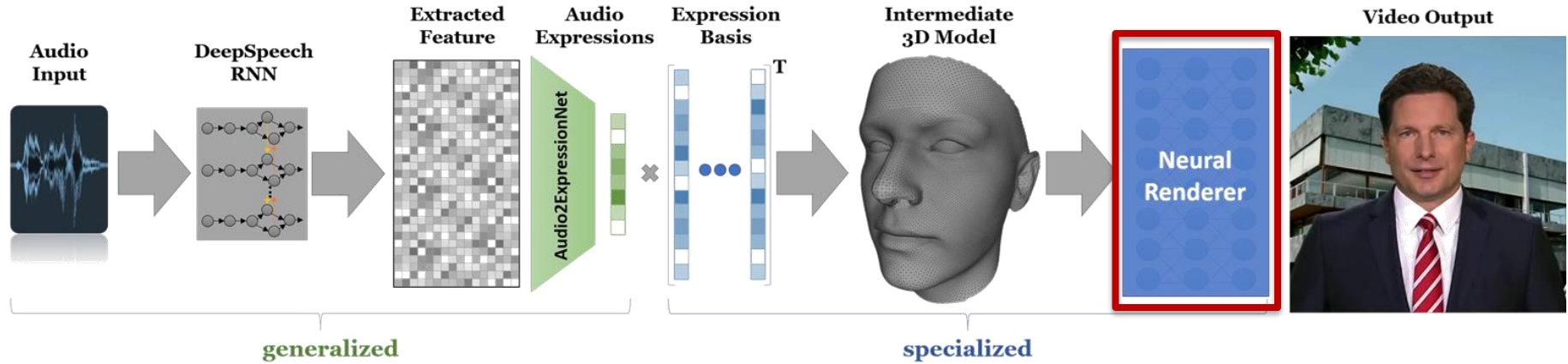
Flame Model

Basel Model

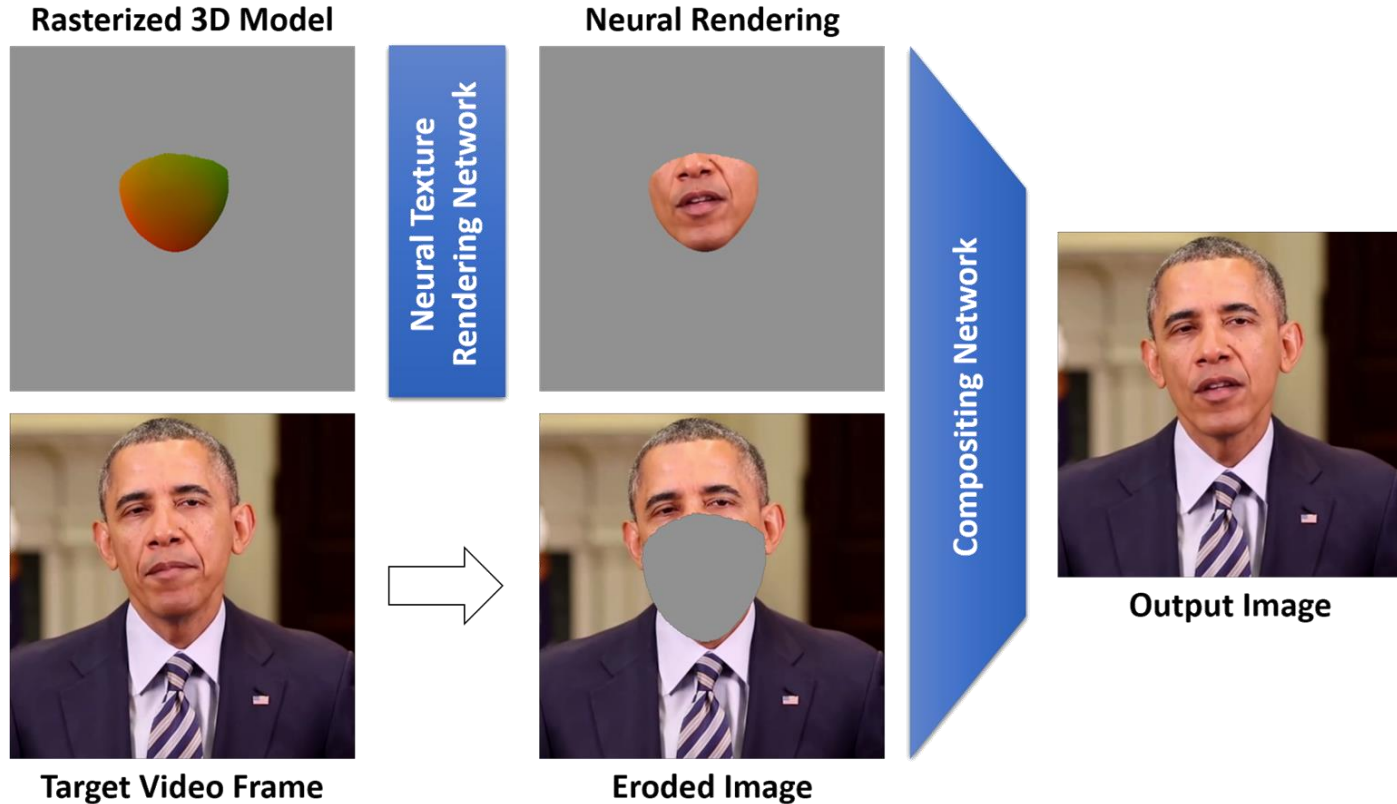
Neural Voice Puppetry



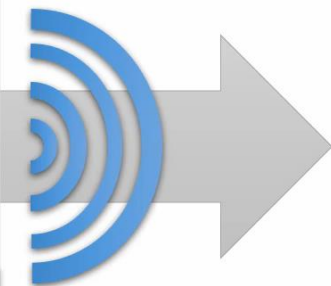
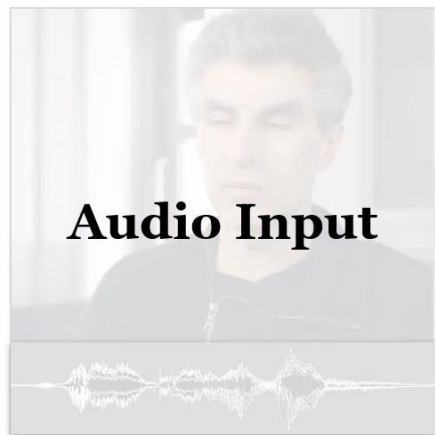
Neural Voice Puppetry



Neural Voice Puppetry



Neural Voice Puppetry



Big Open Challenges

Big Open Challenges



Photo-realistic Reconstruction

Big Open Challenges: How much can AI do?

Using a Bounding Box as Proxy



Input UV-Map



**Sampled
Texture**



Ours

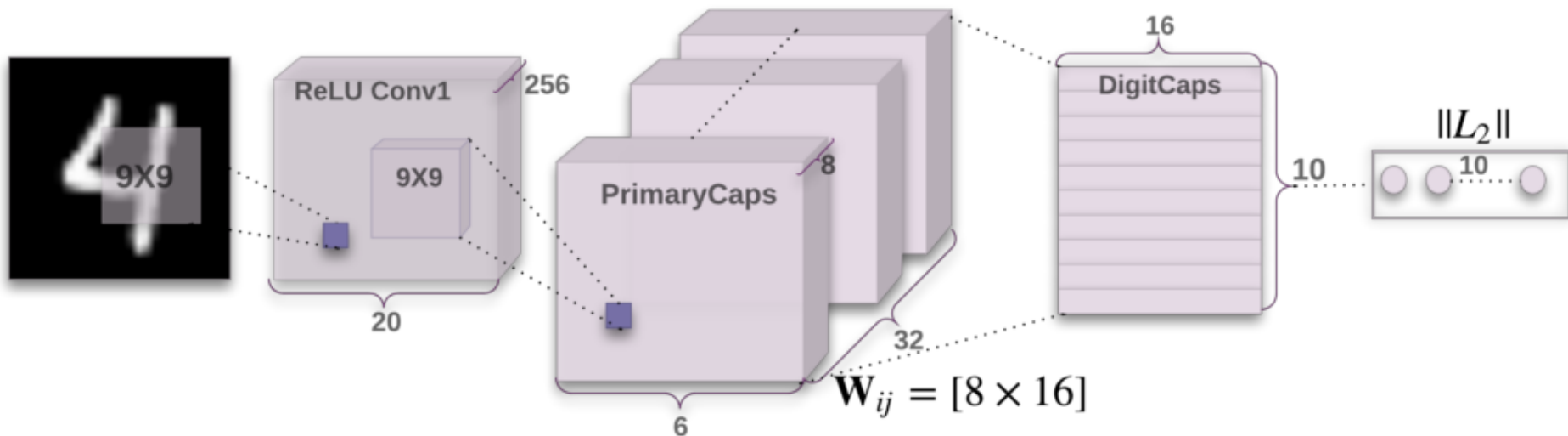


Ground Truth

Big Open Challenges: 3D in Networks

Why learn 3D operations, such as transformations?

-> *differentiate known operators*



Capsule networks are motivated by *inverse graphics* [Sabour et al. 17]

State of the Art on Neural Rendering

A. Tewari^{1*} O. Fried^{2*} J. Thies^{3*} V. Sitzmann^{2*} S. Lombardi⁴ K. Sunkavalli⁵ R. Martin-Brualla⁶ T. Simon⁴ J. Saragih⁴ M. Nießner³
R. Pandey⁶ S. Fanello⁶ G. Wetzstein² J.-Y. Zhu⁵ C. Theobalt¹ M. Agrawala² E. Shechtman⁵ D. B Goldman⁶ M. Zollhöfer⁴

¹MPI Informatics ²Stanford University ³Technical University of Munich ⁴Facebook Reality Labs ⁵Adobe Research ⁶Google Inc *Equal contribution.

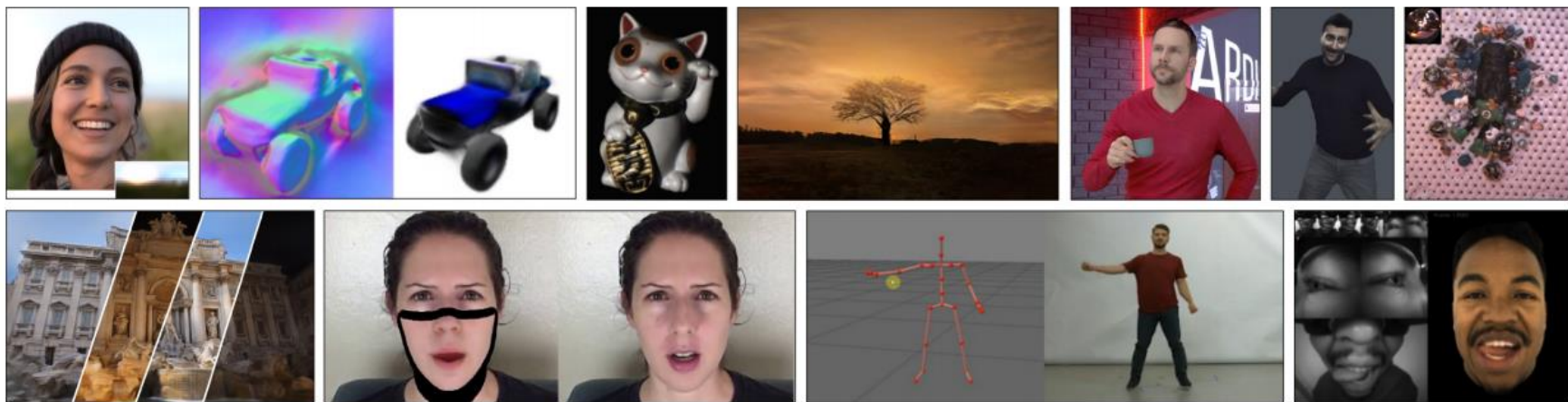


Figure 1: Neural renderings of a large variety of scenes. See Section 6 for more details on the various methods. Images from [SBT* 19, SZW19, XBS* 19, KHM17, GLD* 19, MBPY* 18, XSHR18, MGK* 19, FTZ* 19, LXZ* 19, WSS* 19].

See you next week 😊

Some Extra Slides:

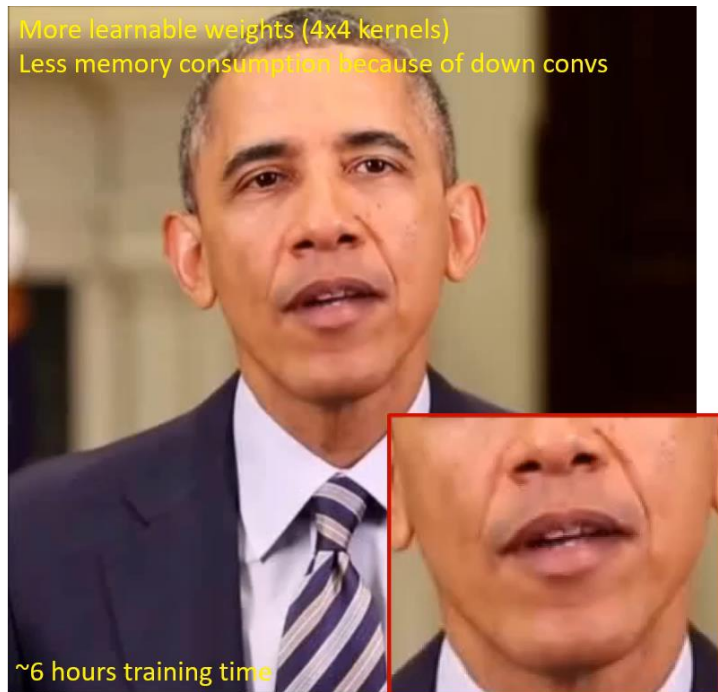
Neural Voice Puppetry

Comparisons **Model-based Methods**

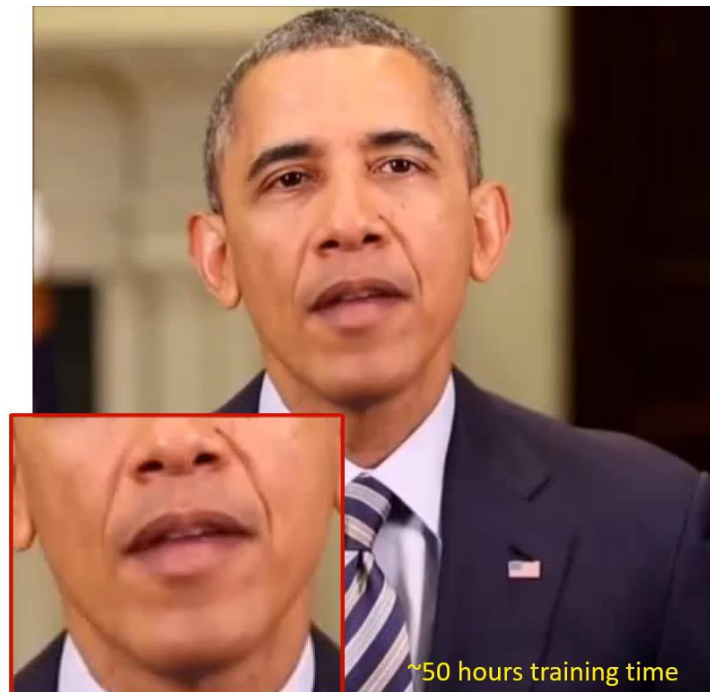
Neural Voice Puppetry

Comparisons **2D-based Methods**

Neural Voice Puppetry



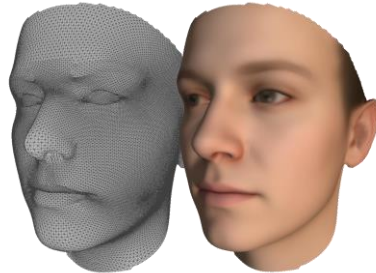
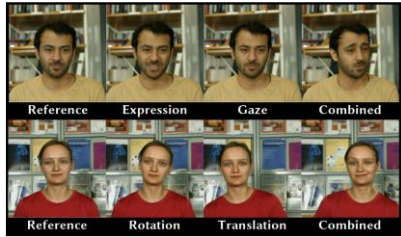
Strided Convolutions
(classical U-Net, 5 down & up convs,
kernel size 4)



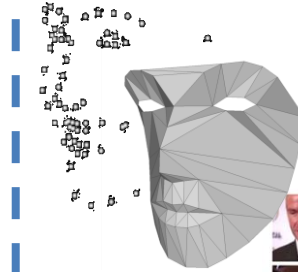
Dilated Convolutions
(U-Net, dilated instead of strided convs
increasing dilation per layer, kernel size 3)

Facial Reenactment

Dense

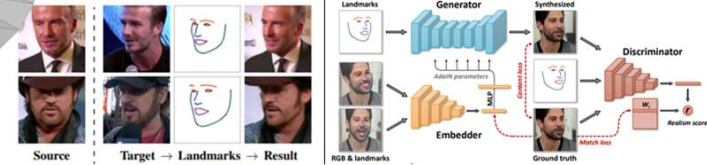


Sparse



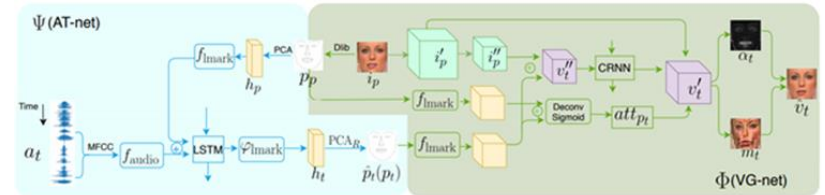
Few-Shot Adversarial Learning of Realistic Neural Talking Head Models

Egor Zakharov^{1,2} Aliaksandra Shysheya^{1,2} Egor Burkov^{1,2} Victor Lempitsky^{1,2}
¹Samsung AI Center, Moscow ²Skolkovo Institute of Science and Technology

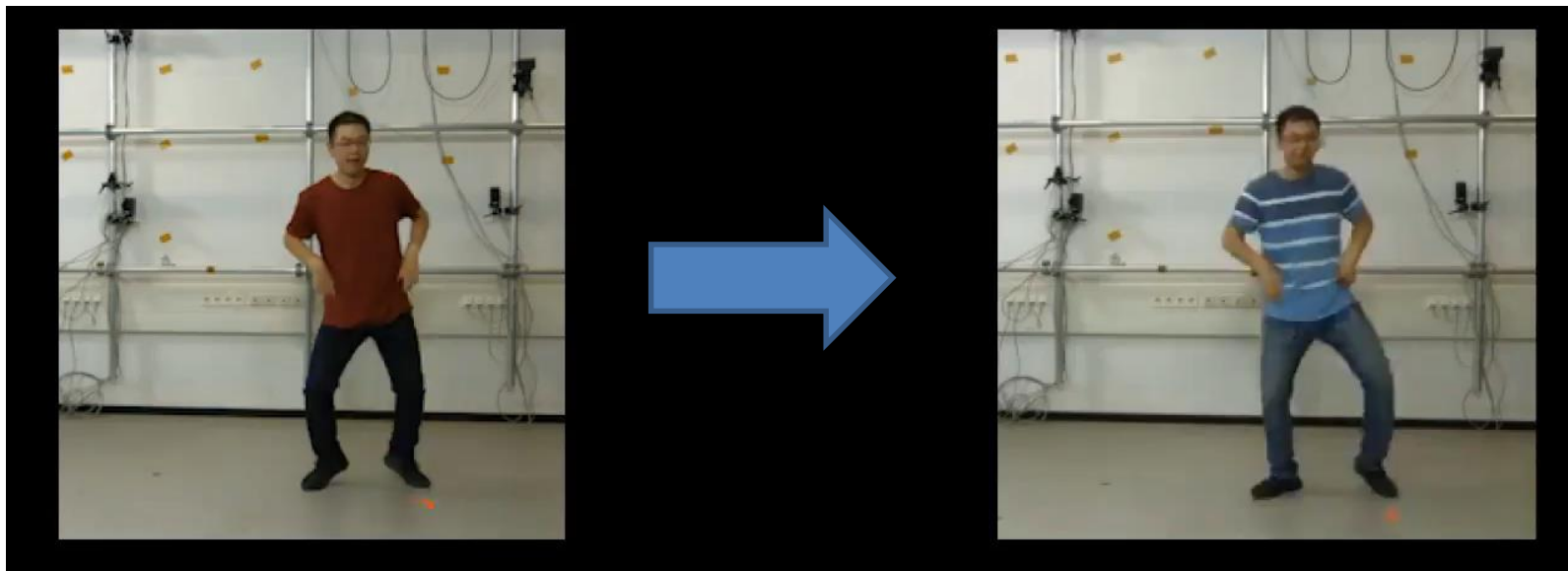


Hierarchical Cross-Modal Talking Face Generation with Dynamic Pixel-Wise Loss

Lele Chen Ross K. Maddox Zhiyao Duan Chenliang Xu
 University of Rochester, USA

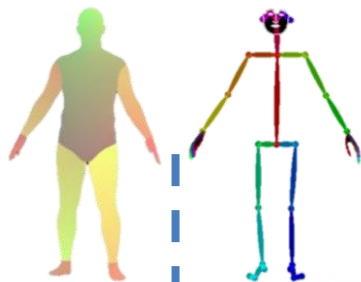


Neural Rendering and Reenactment of Human Actor Videos



Body Reenactment

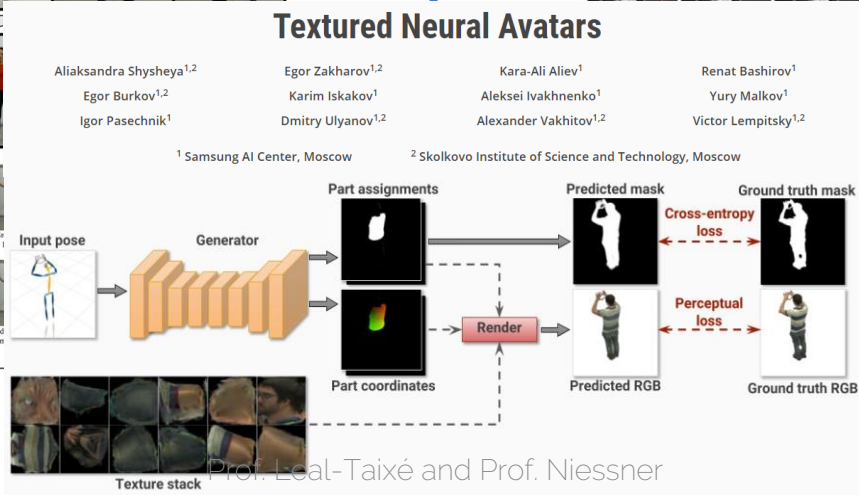
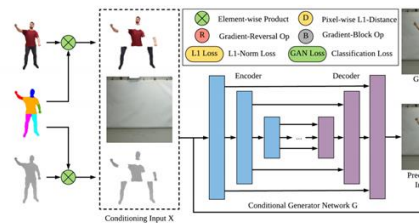
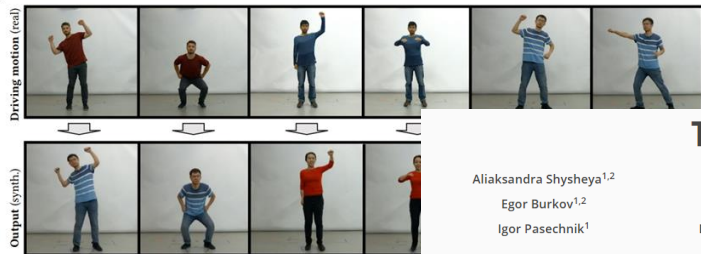
Dense



Sparse

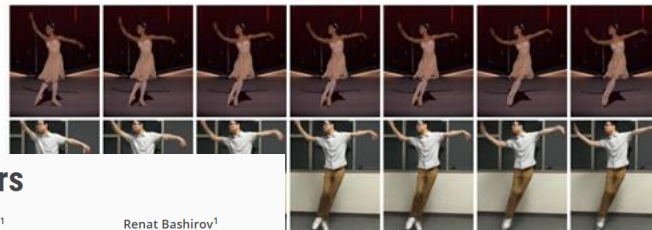
Neural Rendering and Reenactment of Human Actor Videos

LINGJIE LIU, University of Hong Kong, Max Planck Institute for Informatics
 WEIPENG XU, Max Planck Institute for Informatics
 MICHAEL ZOLLHÖFER, Stanford University, Max Planck Institute for Informatics
 HYEONGWOO KIM, FLORIAN BERNARD, and MARC HABERMANN, Max Planck Institute for Informatics
 WENPING WANG, University of Hong Kong
 CHRISTIAN THEOBALT, Max Planck Institute for Informatics

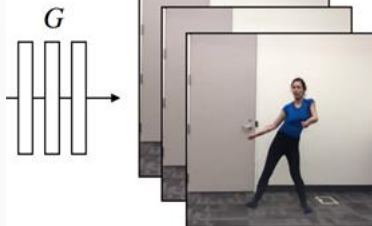


Everybody Dance Now

Caroline Chan* Shiry Ginosar Tinghui Zhou[†] Alexei A. Efros
 UC Berkeley



$G(x_1), \dots, G(x_t)$



Open Challenges

- Motion Capturing
- Person-specific Motions/Expressions
- Temporal Stability
- Image Quality

