

# Causal and compositional generative models in online perception

Michael Janner

with Ilker Yildirim, Mario Belledonne, Christian  
Wallraven, Winrich Freiwald, Joshua Tenenbaum  
MIT





## ResNet-18 predictions

zebra	<div style="width: 98%;"></div>	98%
dalmation	<div style="width: .2%;"></div>	.2%
park bench	<div style="width: .1%;"></div>	.1%
.		
.		
.		

# Outline

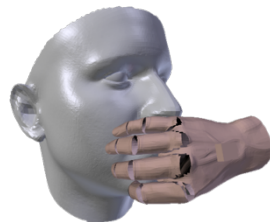
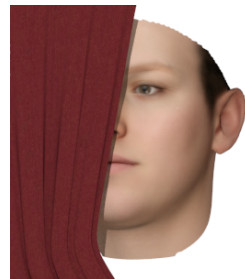
1. Occluded face perception with causal and compositional models





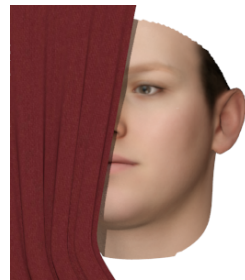
# Outline

1. Occluded face perception with causal and compositional models
2. Automatic training-free vision-to-touch crossmodal transfer

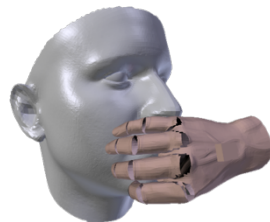


# Outline

1. Occluded face perception with causal and compositional models



2. Automatic training-free vision-to-touch crossmodal transfer



3. Learning visual causal models for generic object categories



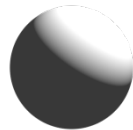
Composite



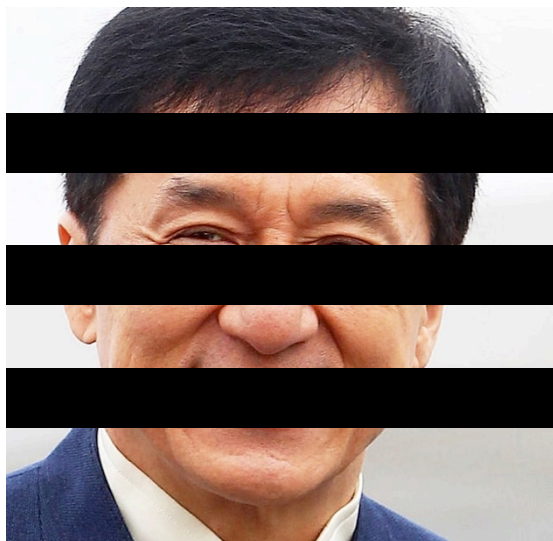
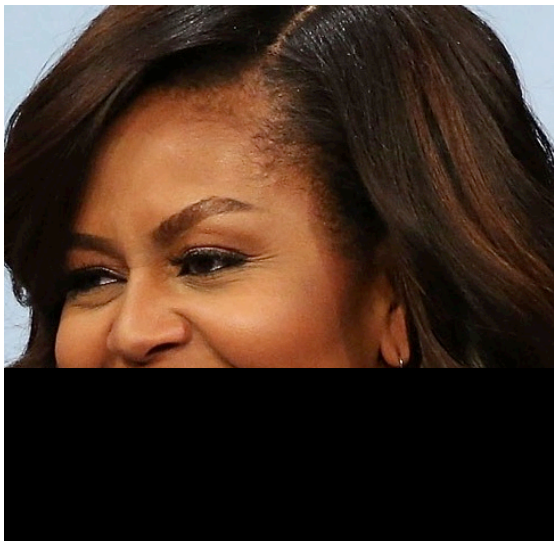
Texture



Shape



Lighting



# The generative model

## Intrinsic

Shape  $S$ ;  
Texture  $T$

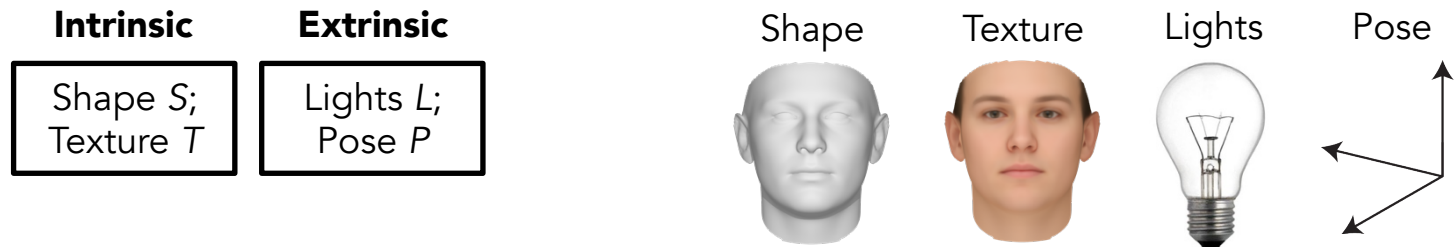
Shape



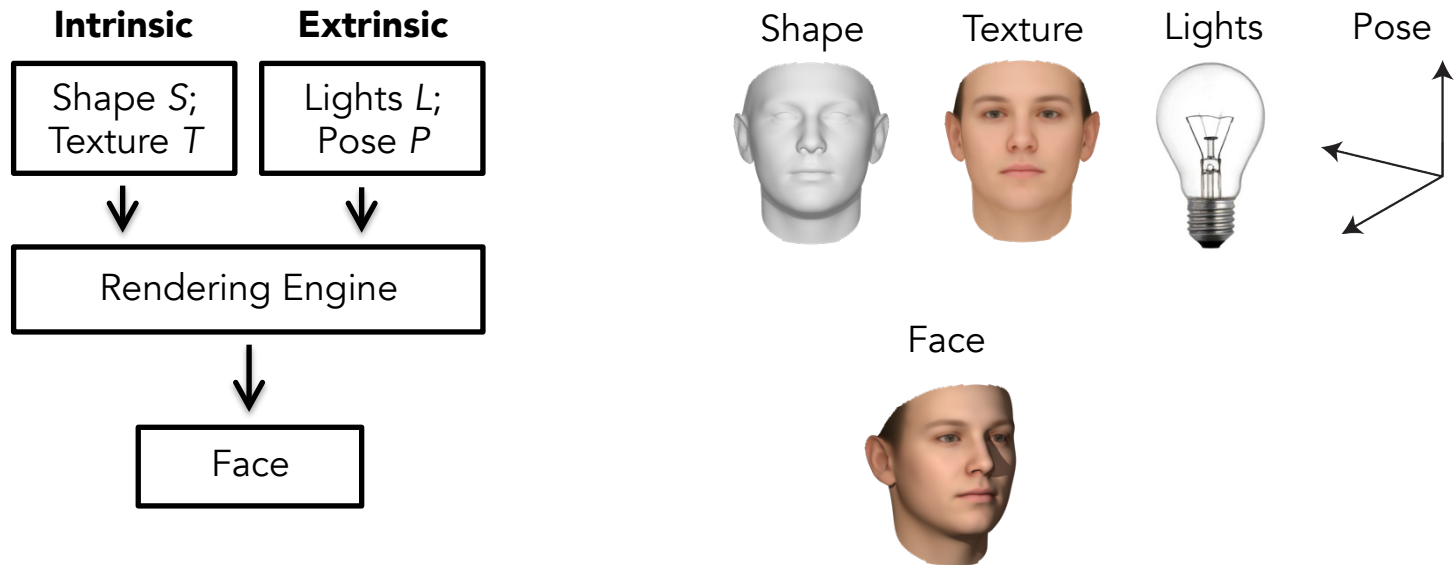
Texture



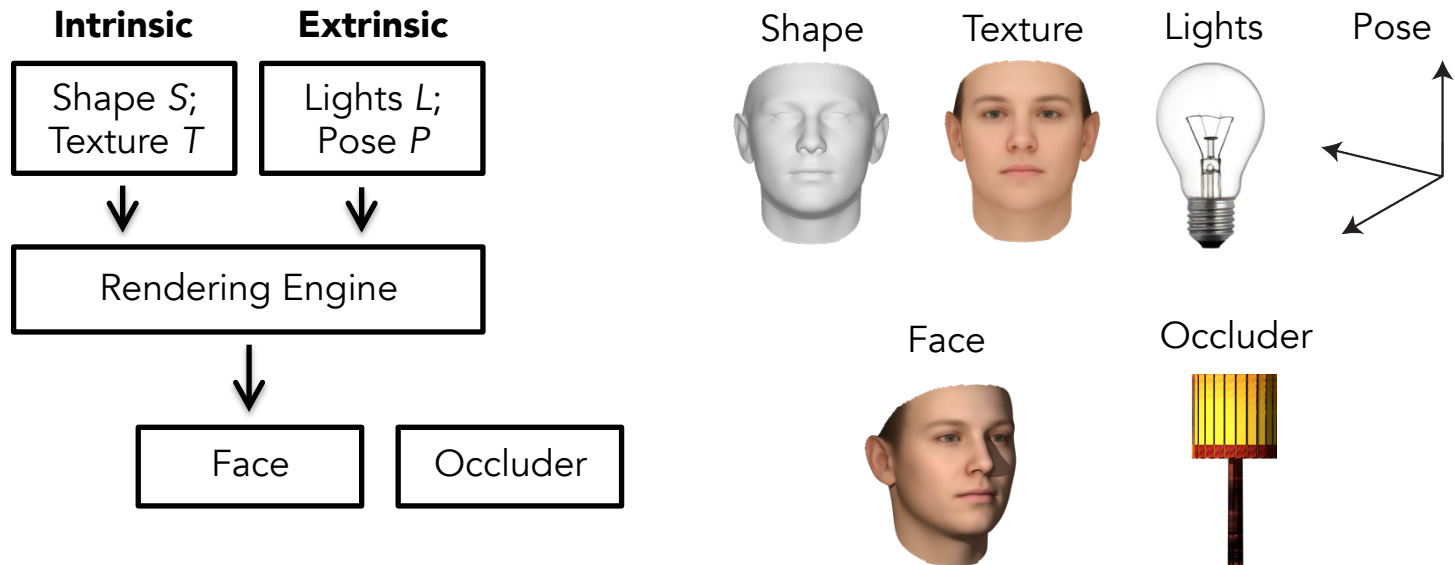
# The generative model



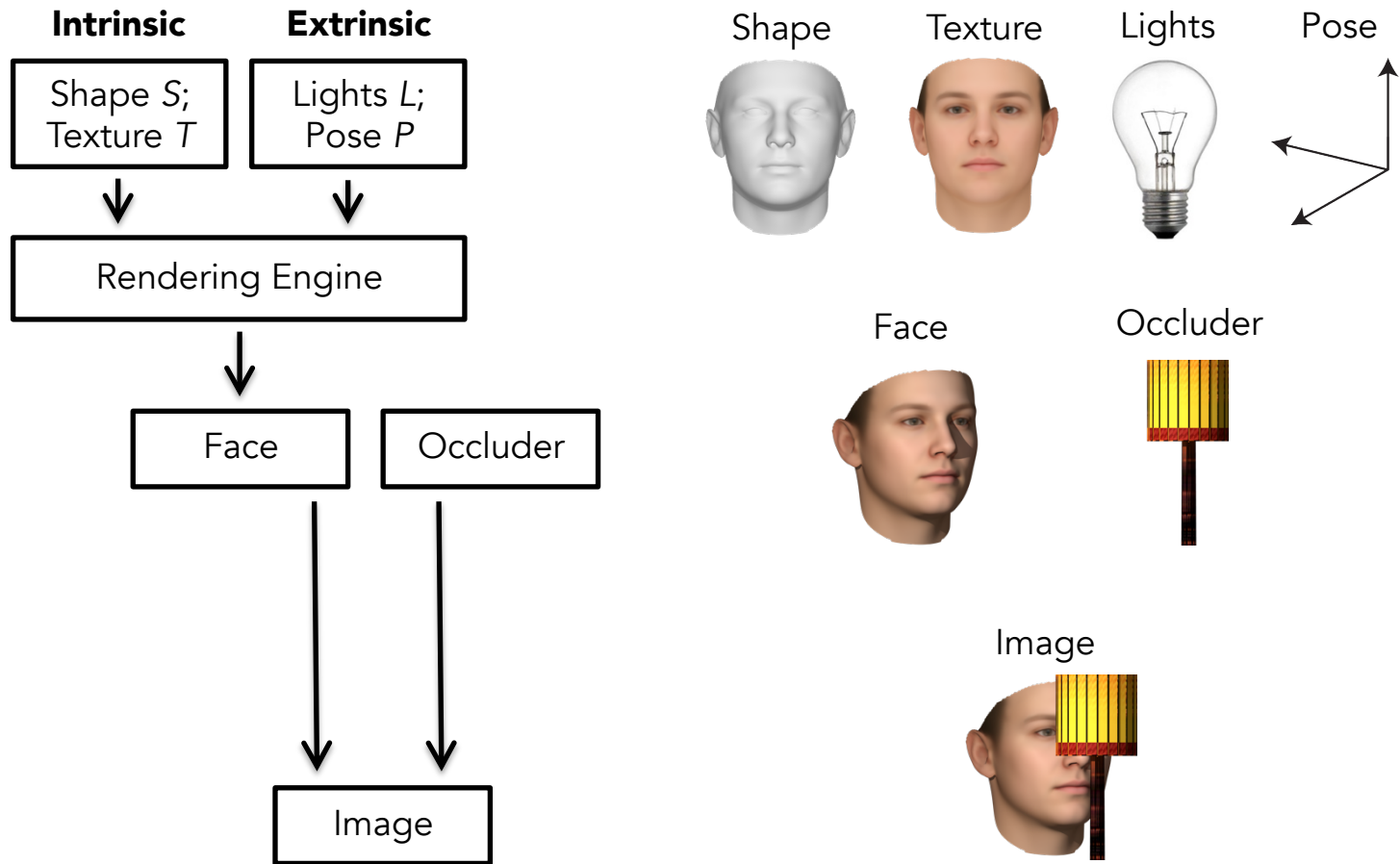
# The generative model



# The generative model

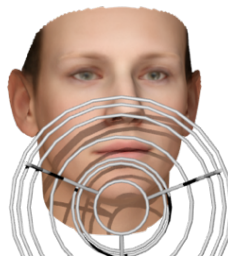
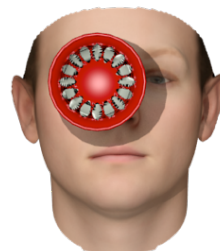
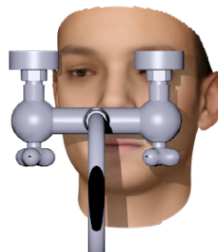
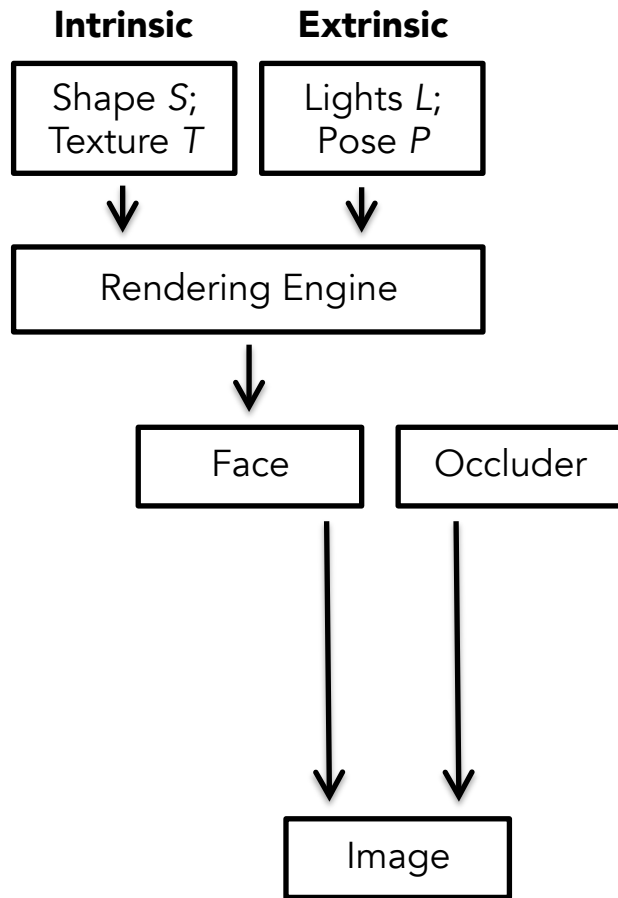


# The generative model

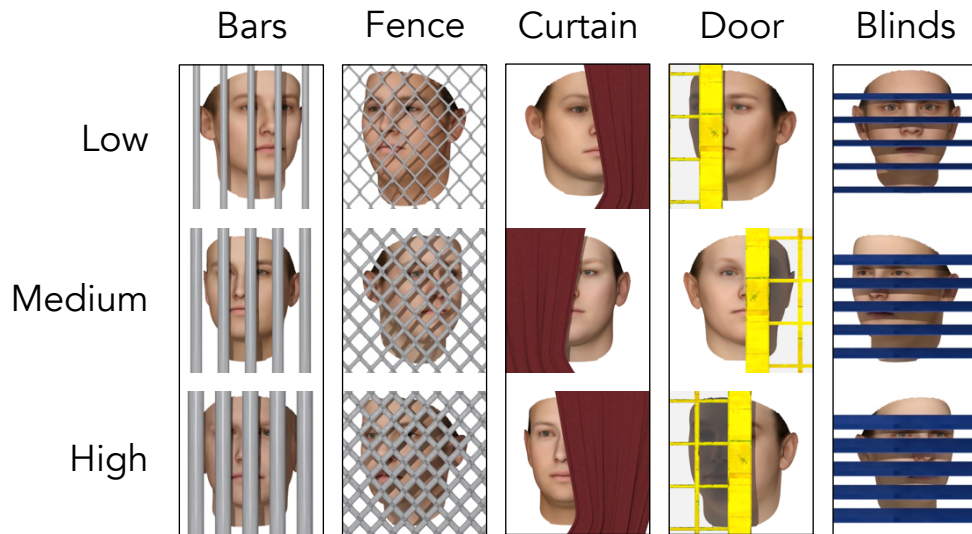




# Samples from the generative model



# Test occluders



- Occlude face in stripes, mesh patterns, and large patches
- Three levels of occlusion, ranging from 15-55% of face covered

# Occluded face perception task



# Occluded face perception task



# Occluded face perception task

**Occluded → Unoccluded**



**Same**

# Occluded face perception task



# Occluded face perception task



# Occluded face perception task

**Unoccluded → Occluded**



**Different**



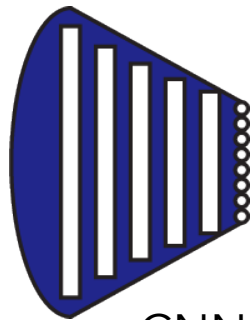
# Behavioral results

## Accuracy

	Occlusion Level		
	Low	Medium	High
Occluded → Unoccluded	.77	.73	.67
Unoccluded → Occluded	.78	.76	.70

Chance: .5

# Naïve model: Learn to ignore the occluder



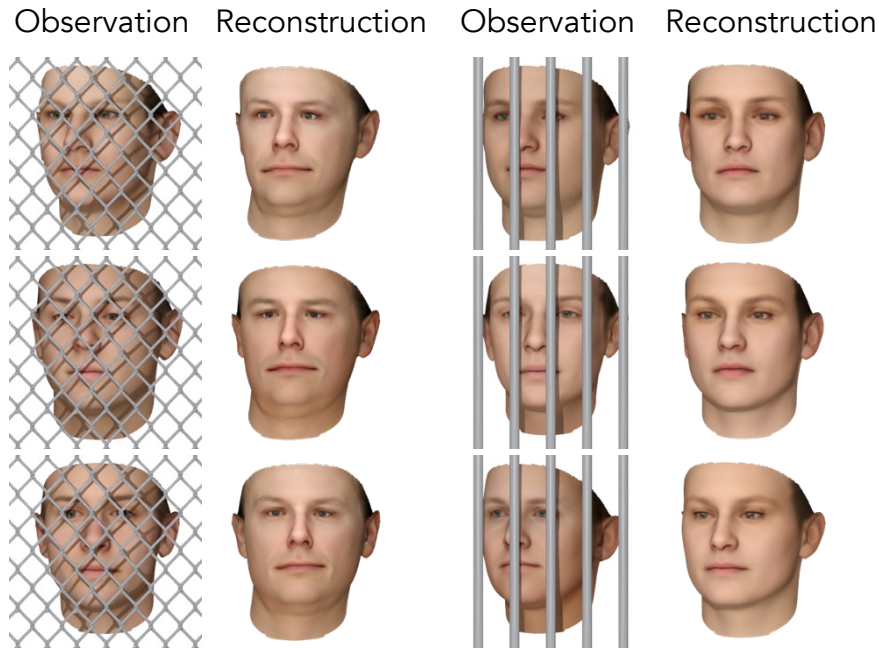
CNN

Shape  $S$ ;  
Texture  $T$

Predict intrinsic and extrinsic face parameters directly.

Can model become invariant to all types of occluders?

# Naïve model: Learn to ignore the occluder

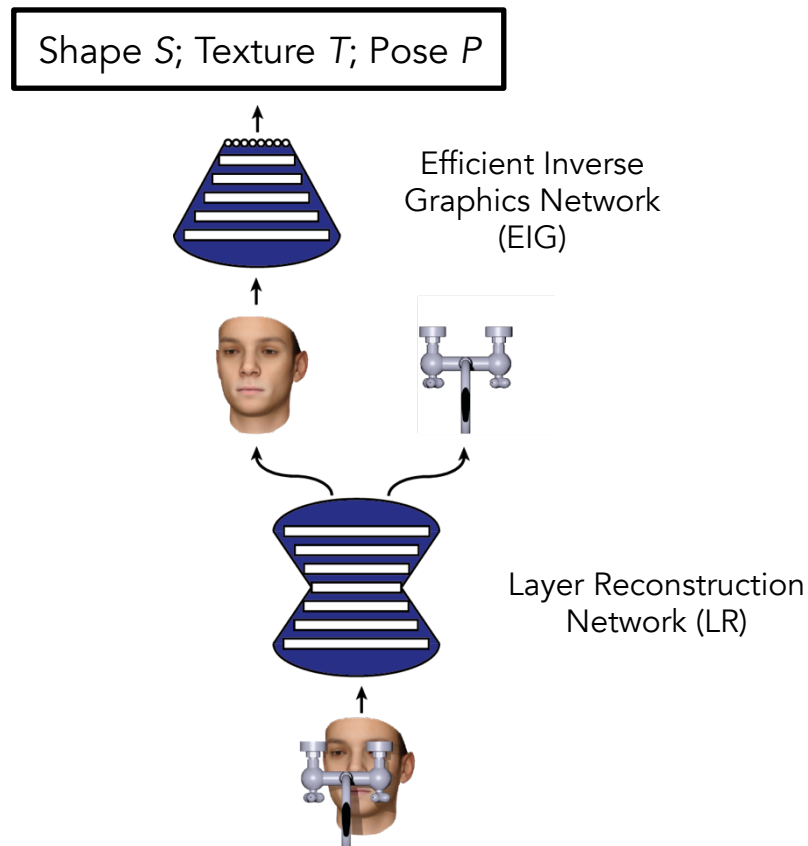
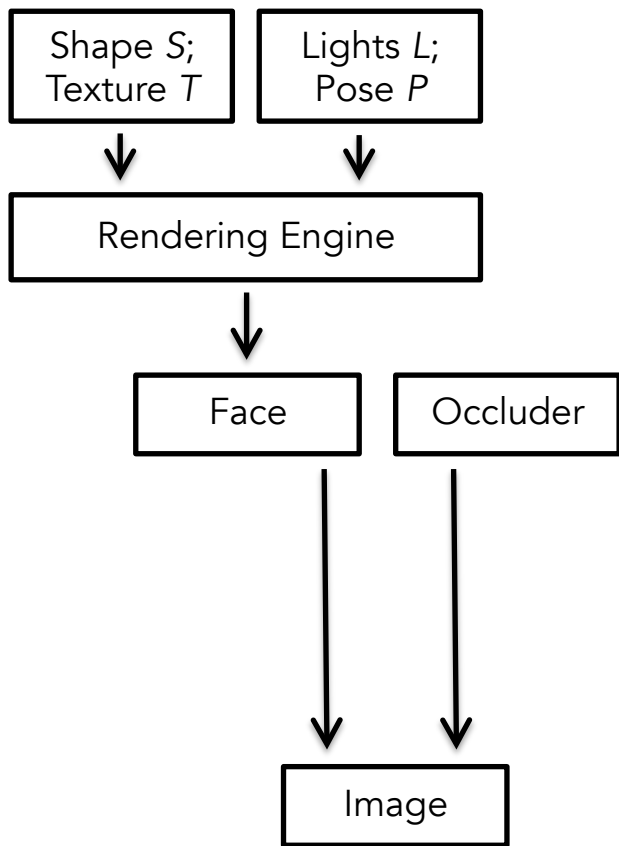


Predict intrinsic and extrinsic face parameters directly.

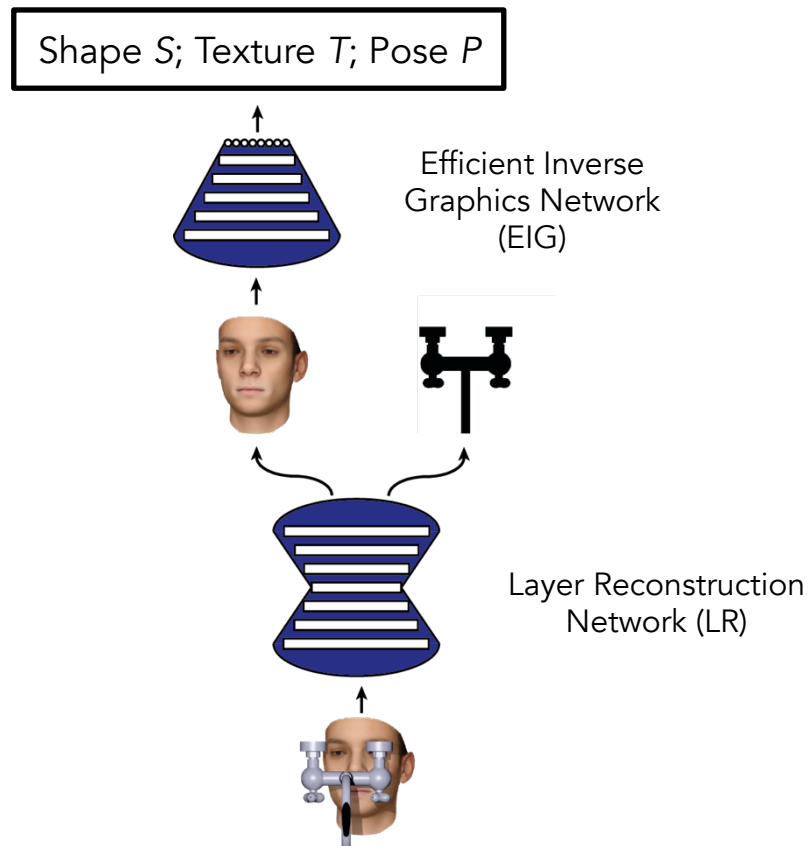
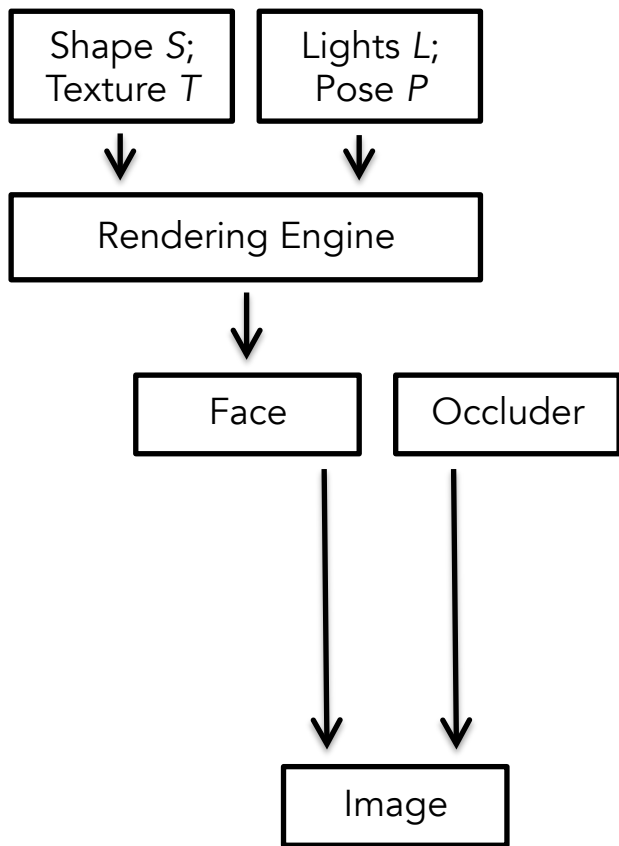
Can model become invariant to all types of occluders? **No.**

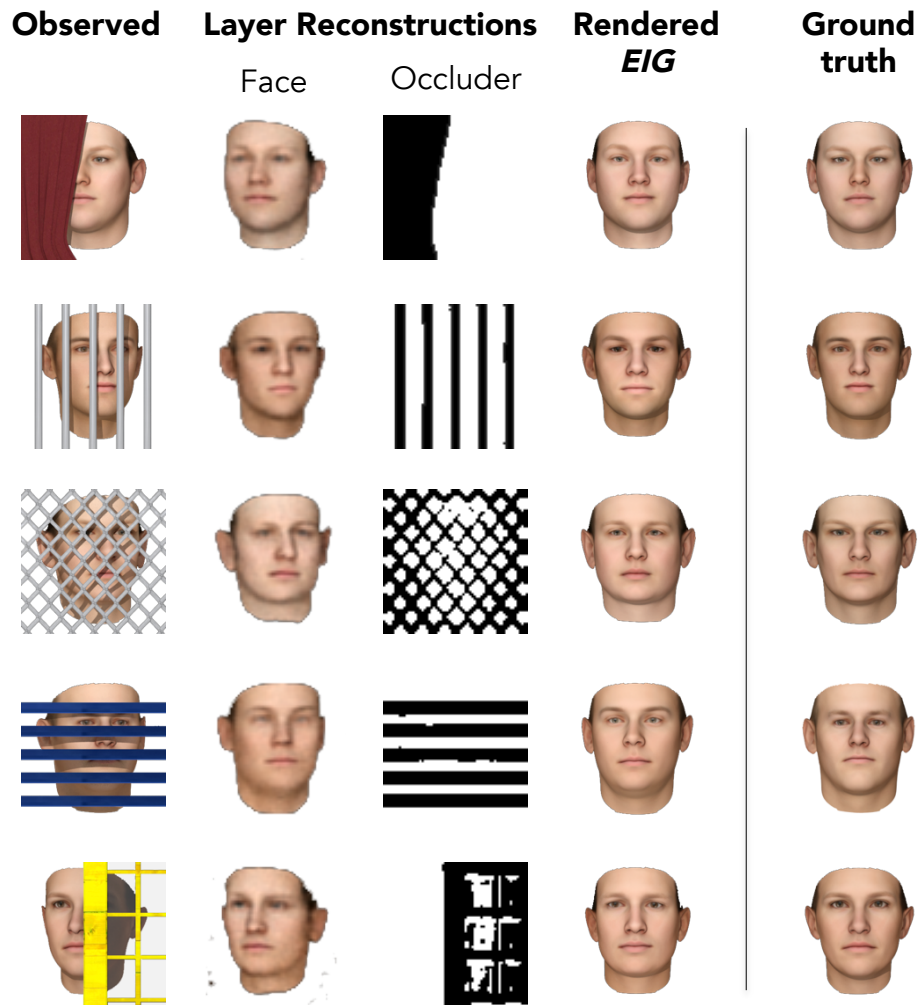
Test occluders heavily influence face predictions.

# Inverse graphics: Inverting the generative model



# Inverse graphics: Inverting the generative model





Modeling causes allows for:

- (1) face predictions that are invariant to occluders
- (2) better generalization to unseen occluders

# Prediction human judgments

**Occluded → Unoccluded**

**Unoccluded → Occluded**

Low

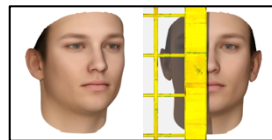
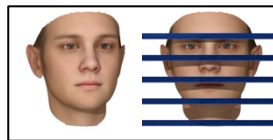
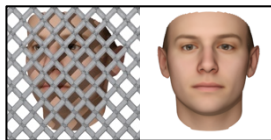
Medium

High

Low

Medium

High



Spearman rank correlation

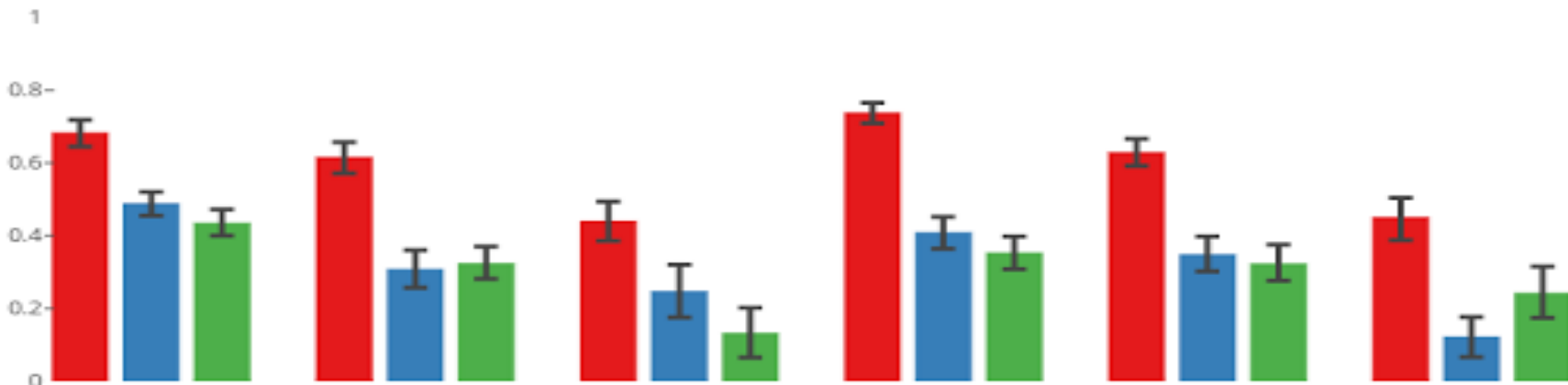


Image-space comparison

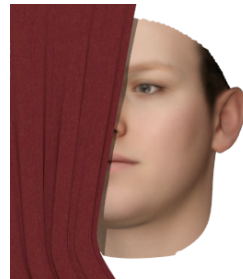


Latent-space baseline

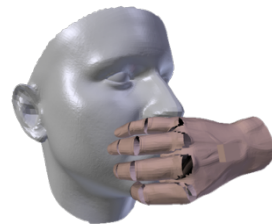


VGG

1. Occluded face perception with causal and compositional models



2. Automatic training-free vision-to-touch crossmodal transfer



3. Learning visual causal models for generic object categories



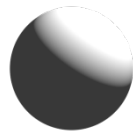
Composite



Texture



Shape



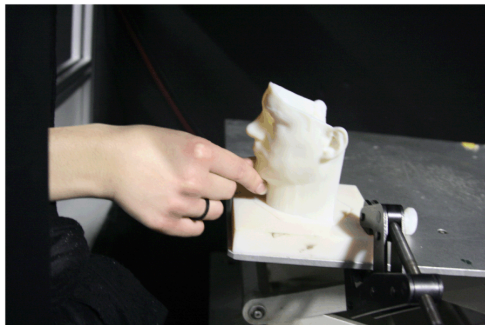
Lighting



# Vision-to-touch transfer



**Study**



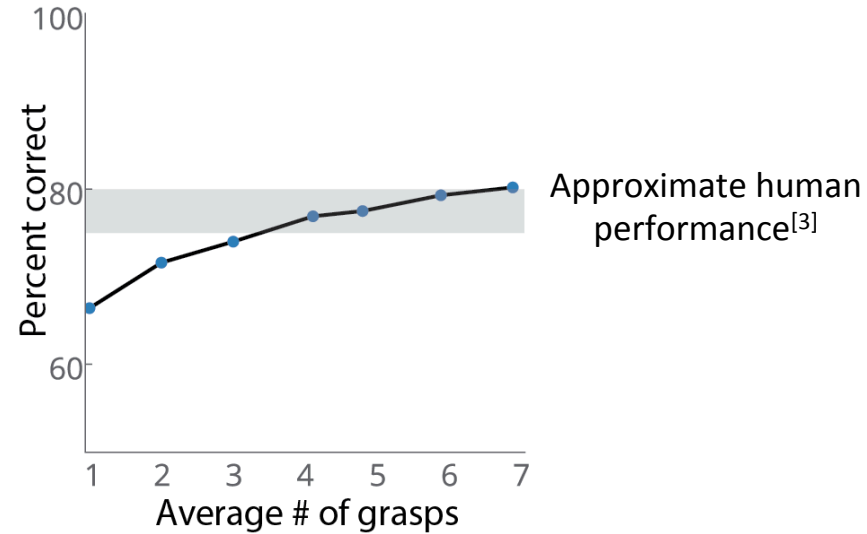
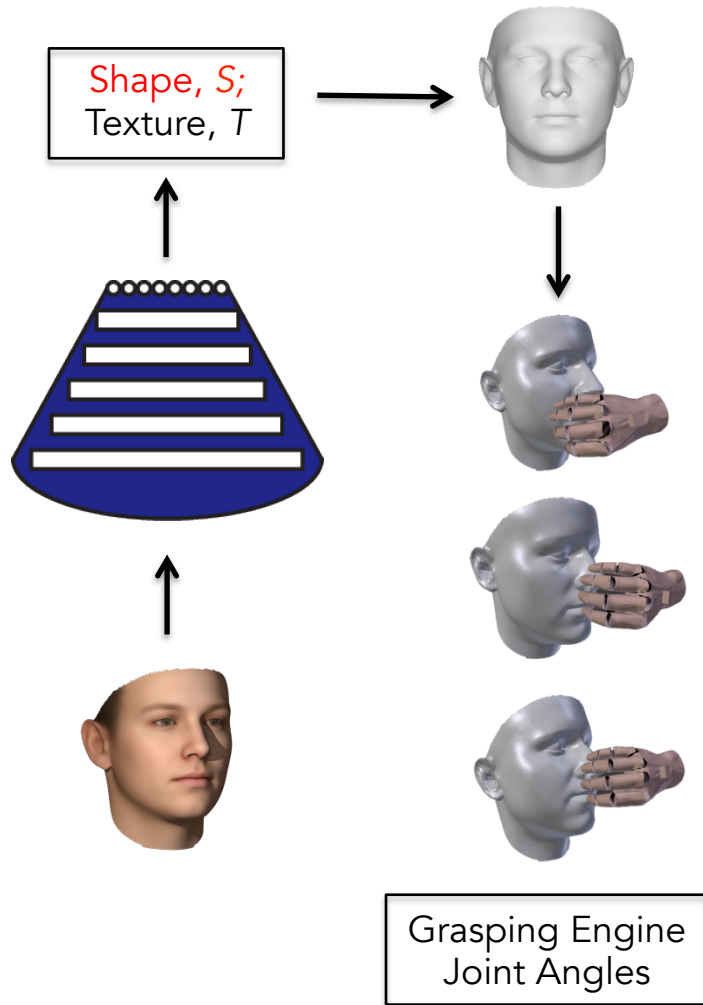
**Test**

Shape is explicitly represented

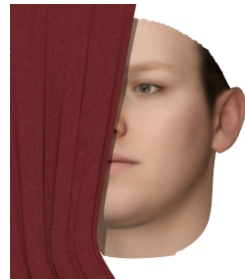
Immediate transfer to haptic domain

Can make judgments based on grasping joint angles

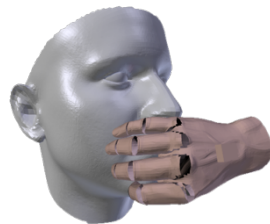
# Vision-to-touch transfer



1. Occluded face perception with causal and compositional models



2. Automatic training-free vision-to-touch crossmodal transfer



3. Learning visual causal models for generic object categories



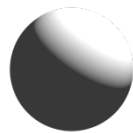
Composite



Texture



Shape



Lighting



# Decomposing generic objects

Inputs



Shader Outputs



Causal intermediate stage  
visual representations of  
reflectance, shape, and lighting

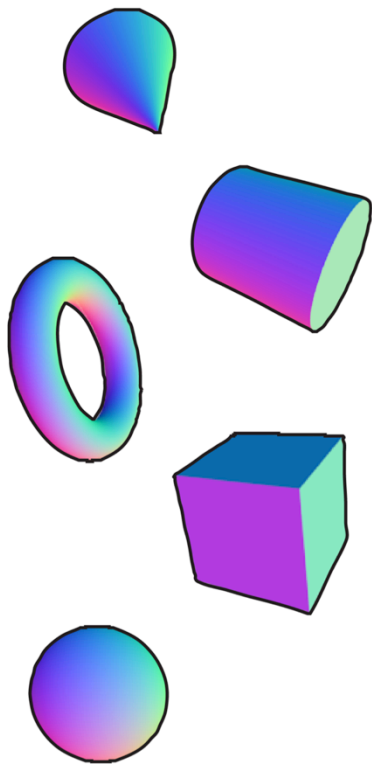
Learn both the recognition  
model and the rendering  
function

# Decomposing generic objects

Causal intermediate stage  
visual representations of  
reflectance, shape, and lighting

Learn both the recognition  
model and the rendering  
function

Representation learning driven  
by reconstruction error



**Train shapes**



**Test shapes**

# Decomposing generic objects

**Observation**



**Initial  
prediction**



**Unsupervised  
improvement**



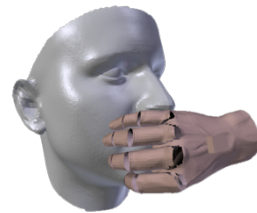
Causal intermediate stage  
visual representations of  
reflectance, shape, and lighting

Learn both the recognition  
model and the rendering  
function

Representation learning driven  
by reconstruction error

# Summary

- Causal and compositional models better predict human judgments
- Modeling causes allows for training-free crossmodal transfer
- Causal models can improve internal representations without supervision





# Thank you



Ilker  
Yildirim



Mario  
Belledonne



Christian  
Wallraven



Winrich  
Freiwald



Joshua  
Tenenbaum

# Comparison Pipeline

Predict latents of both faces

Render study face with pose of test face

Occlude rendering and test image with mask

Compare in **image** (or **latent**) space

