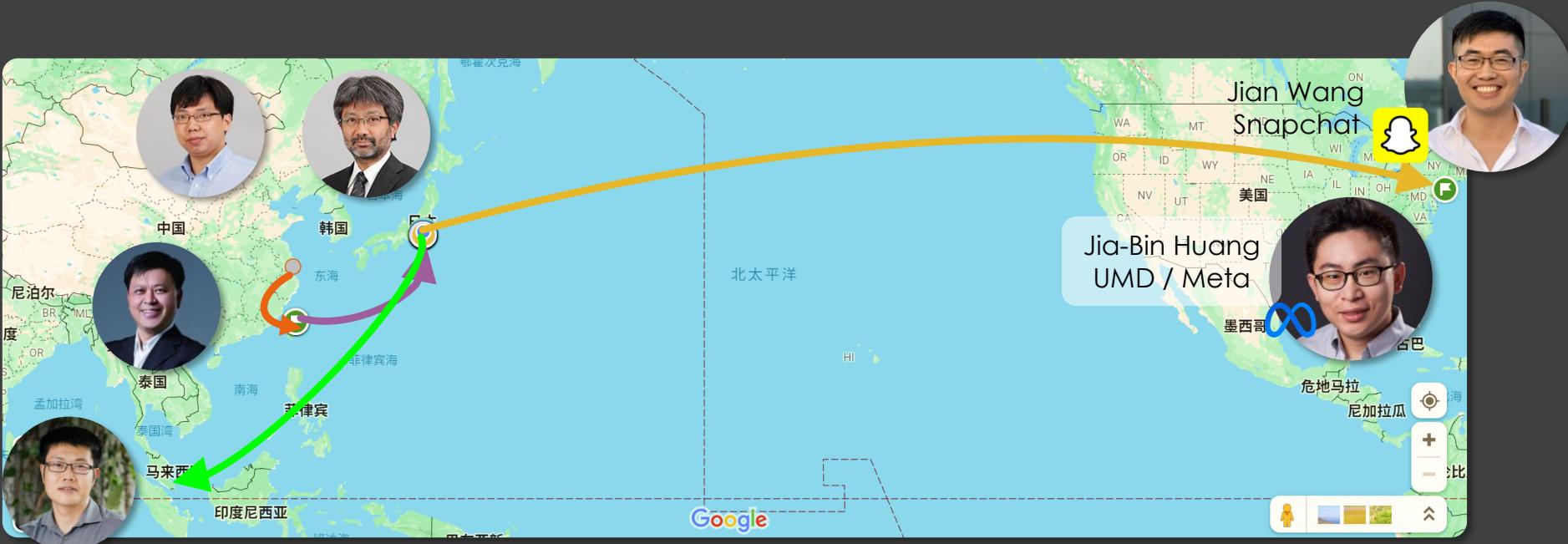


# **Image Factorization and Manipulation with Generative Regularizations**

Zhixiang Wang  
PhD candidate  
The University of Tokyo



# Academic Journey



南京師範大學  
Nanjing Normal University



國立臺灣大學  
National Taiwan University



東京大學  
THE UNIVERSITY OF TOKYO



NANYANG  
TECHNOLOGICAL  
UNIVERSITY  
SINGAPORE

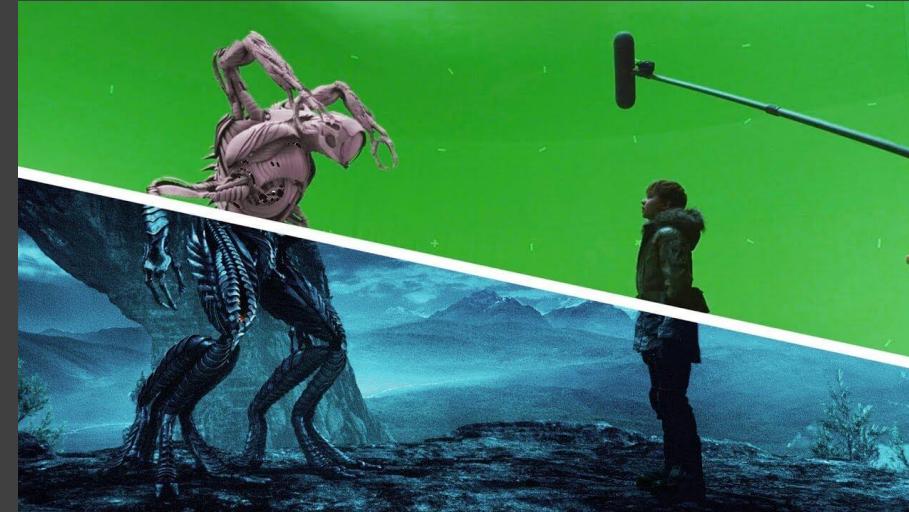


Snapchat

Sony Research

# Goal: GenAI + Advanced Cameras for VFX

Reduce actor, time, and money costs



# Research Works

## Special Hardwares

### Polarimetric Camera



Wang et al, CVPR 2019

### Infrared Camera



Wang et al, CVPR 2019  
Wei, Wang et al, AAAI'23

### Rolling Shutter Camera



Wang et al, CVPR 2022  
Ji, Wang et al, ICCV 2023

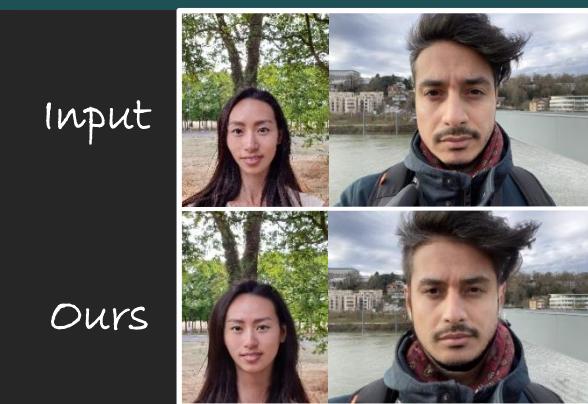
### Foggy Scene Understanding



Ma, Wang et al, CVPR 2022

## Generative Models

### Geometric Distortion Correction



Wang et al, IJCV 2024

### Background Replacement



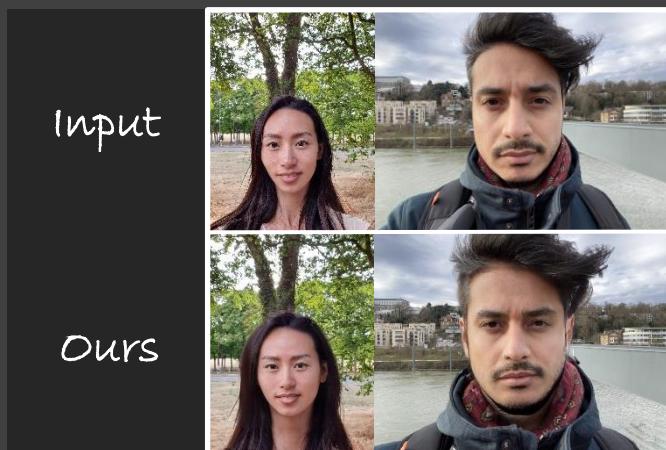
Wang et al, SIGGRAPH 2024

### Style Transfer



Chang, Wang et al, ECCV 2020

# Viewpoint + Lens



**Perspective Distortion Correction**

Wang et al, IJCV 2024

# Background Background



**Matting by Generation**

Wang et al, SIGGRAPH 2024

# Good Photos are Not Easy to Take

Examples of “bad/undesired” photos,  
caused by unwanted imaging factors



Device

Lighting

Viewpoint

Background

# Difficulty in Controlling Imaging Factors



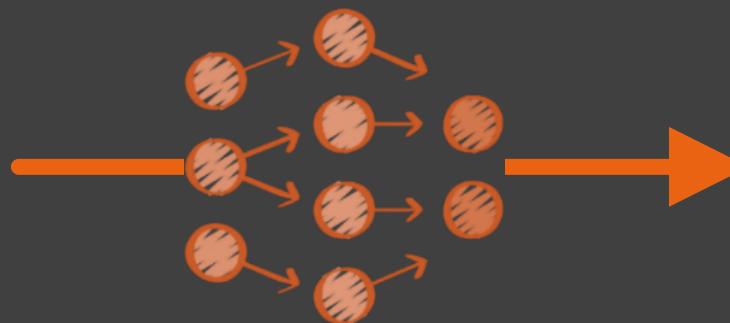
- ▶ Numerous factors
- ▶ Specialized equipment
  - ▶ Inflexible
  - ▶ Expensive
- ▶ Expertise
- ▶ Multiple trials

# Simple yet Popular DL-based Solution



Undesired  
Samplings

Image-to-Image Transform

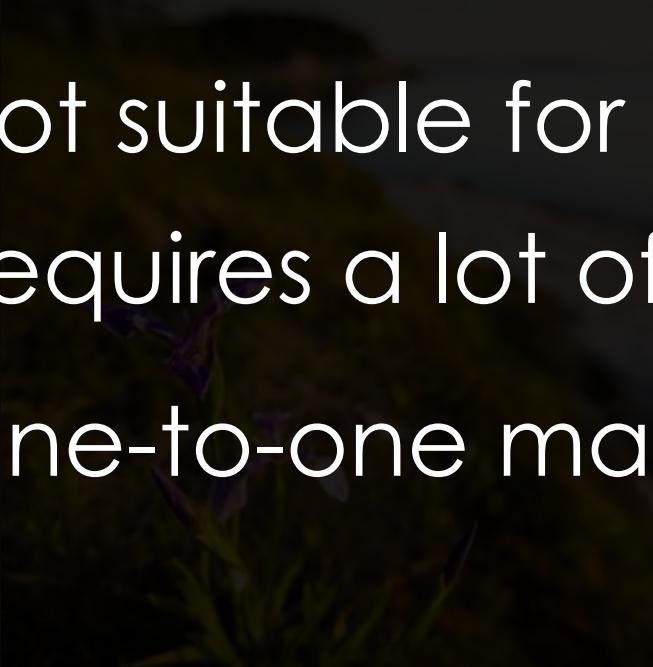


Desired  
Samplings

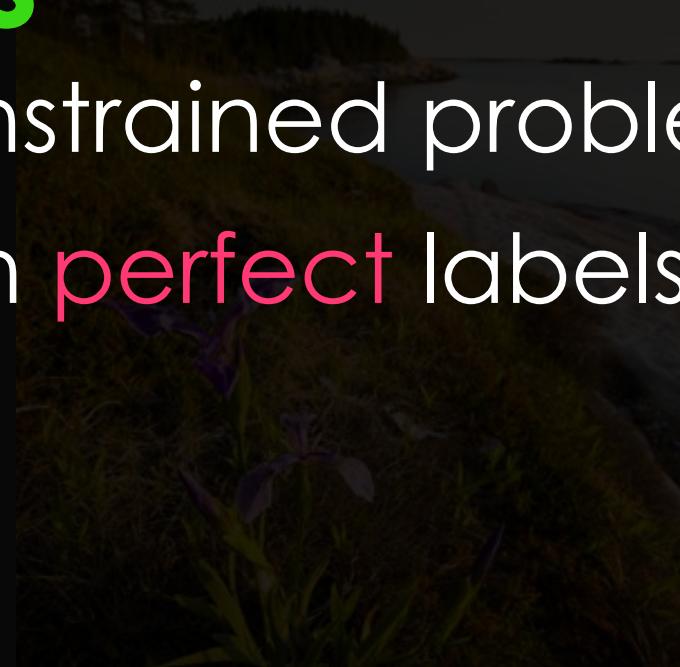
# Popular Approaches

## Challenges

- Not suitable for severe under-constrained problems
- Requires a lot of **paired** data with **perfect** labels
- One-to-one mapping

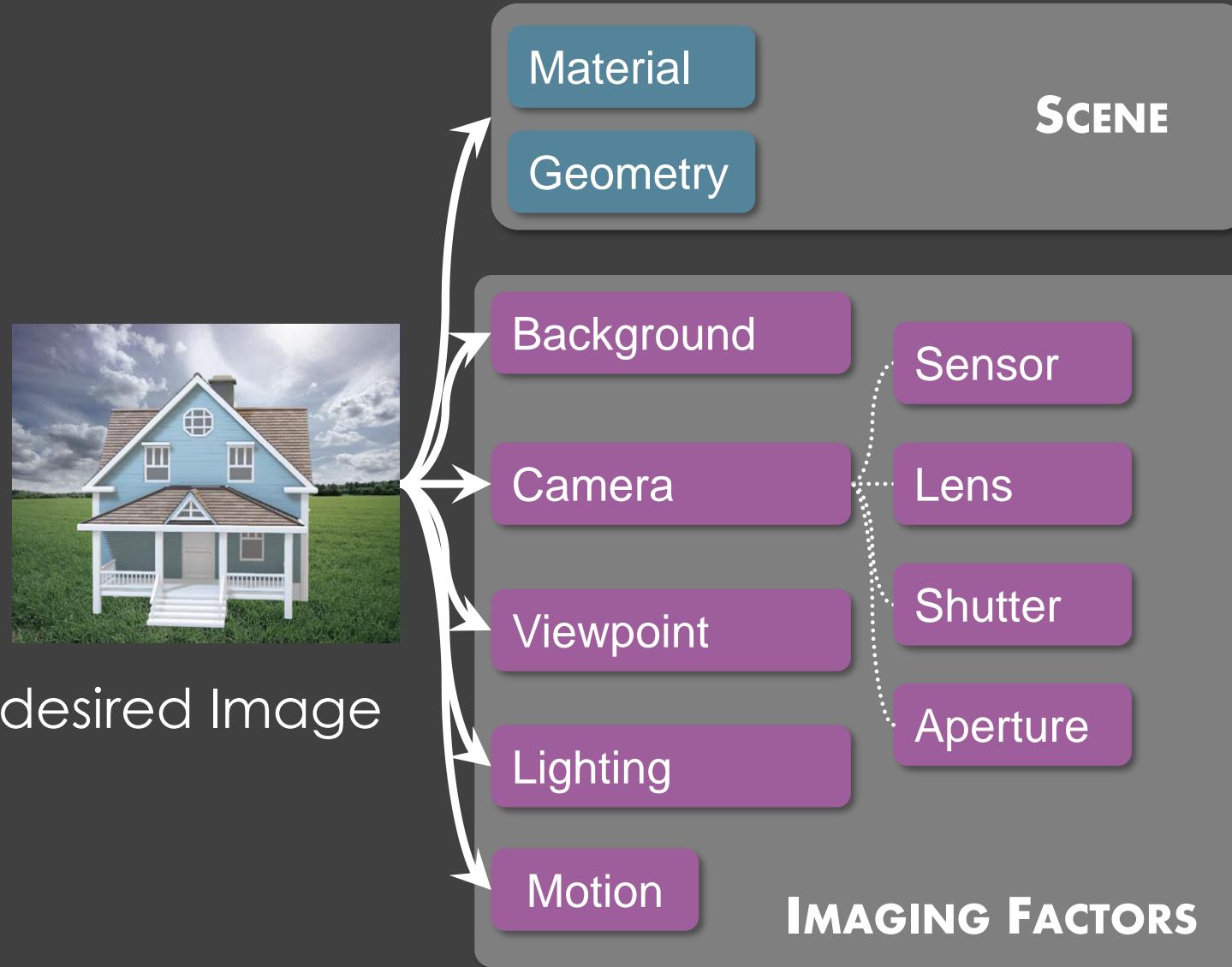


Undesired  
Samplings

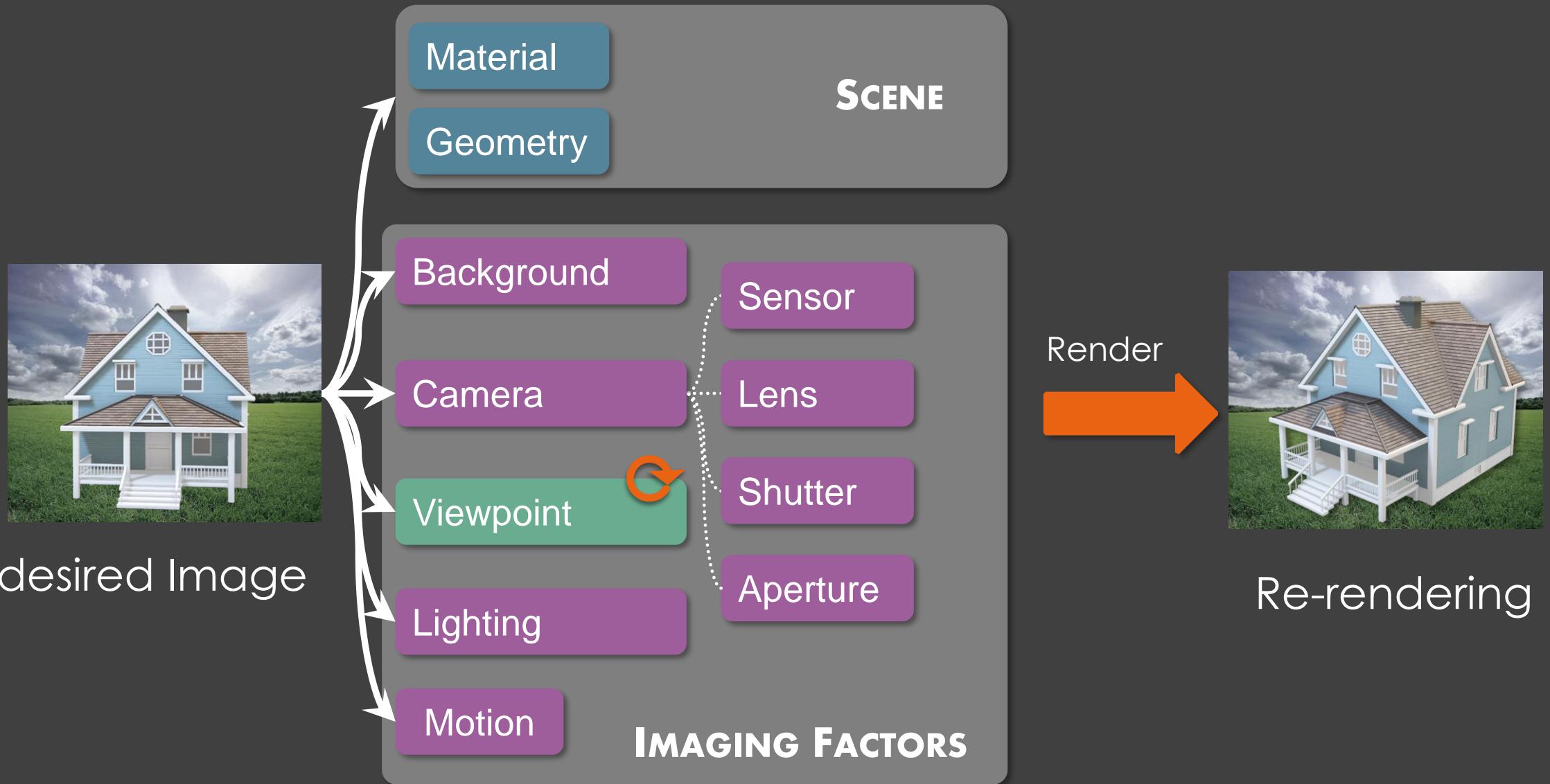


Desired  
Sampling

# Image Factors and Factorization



# Image Manipulation

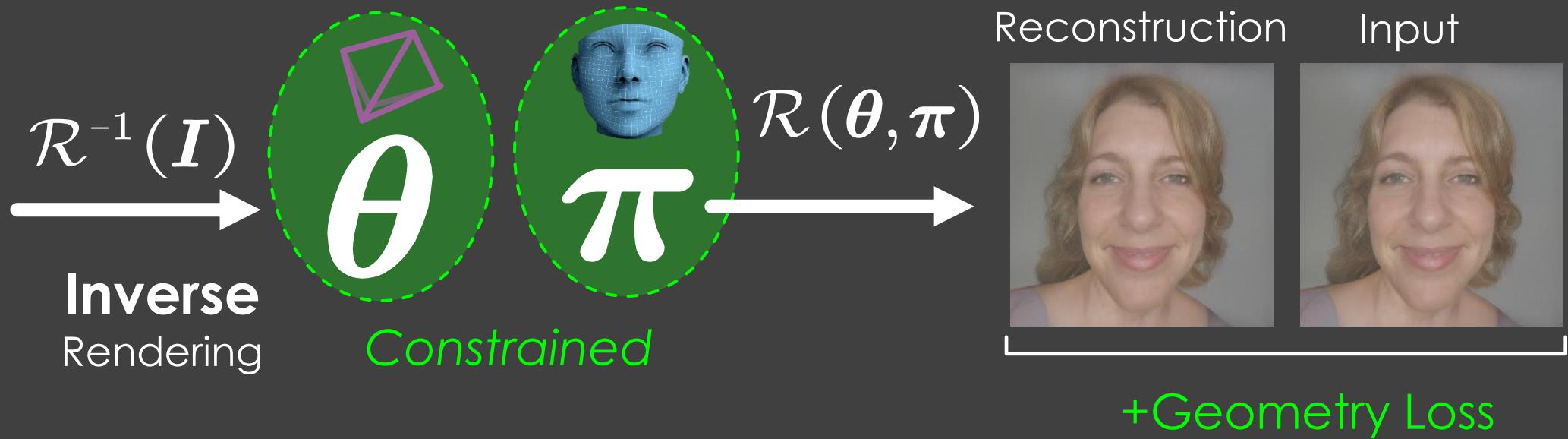


# Harness Pre-trained Generative Models

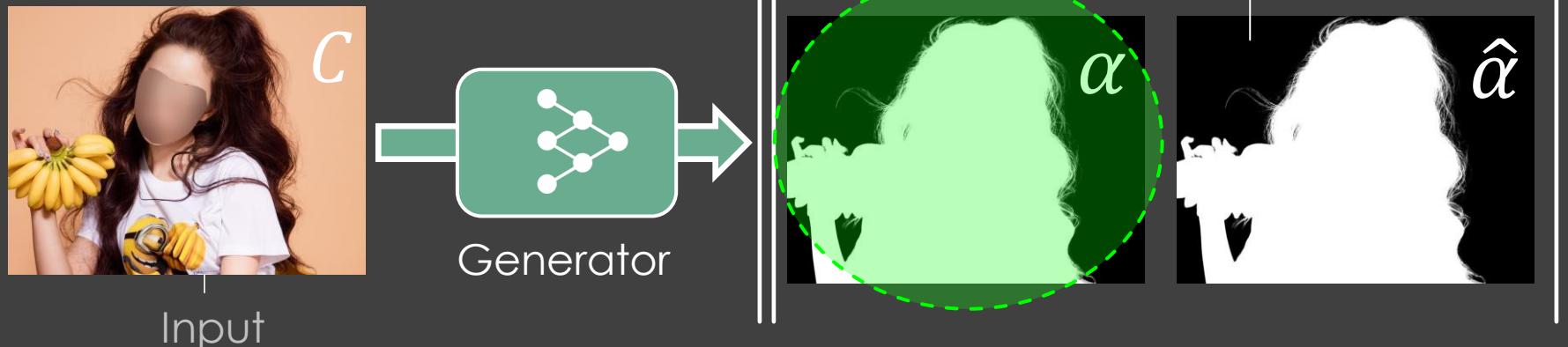
**Optimization-based:** no labels required



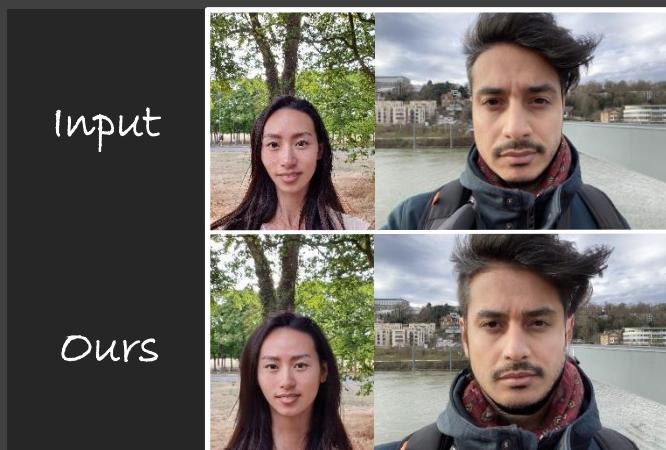
Input



**Learning with Labels:** imperfect labels



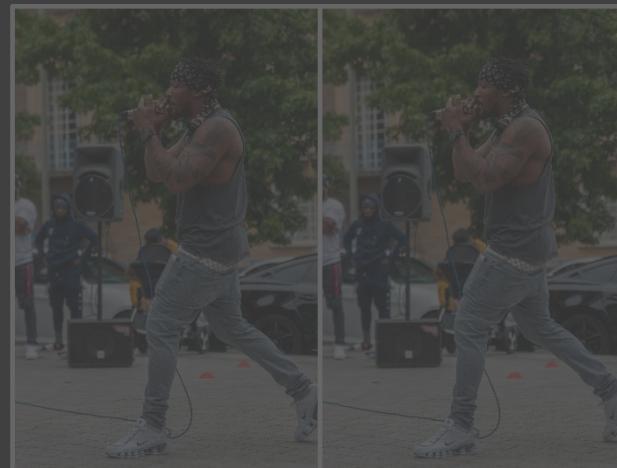
# Viewpoint + Lens



**Perspective Distortion Correction**

Wang et al, IJCV 2024

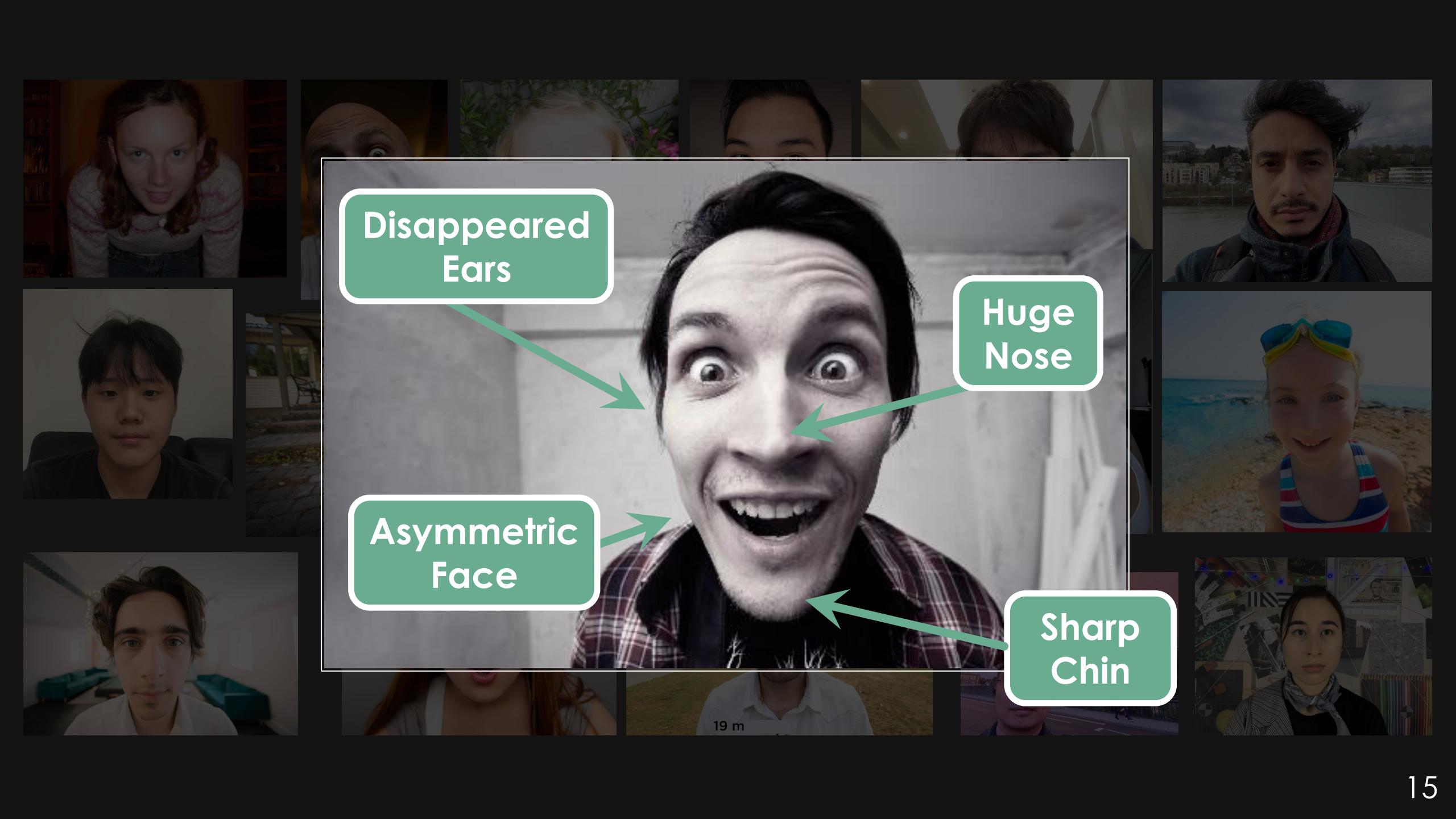
# Background



**Matting by Generation**

Wang et al, SIGGRAPH 2024





Disappeared  
Ears

Huge  
Nose

Asymmetric  
Face

Sharp  
Chin

# Short Camera-to-Subject Distance



# Perspective Projection

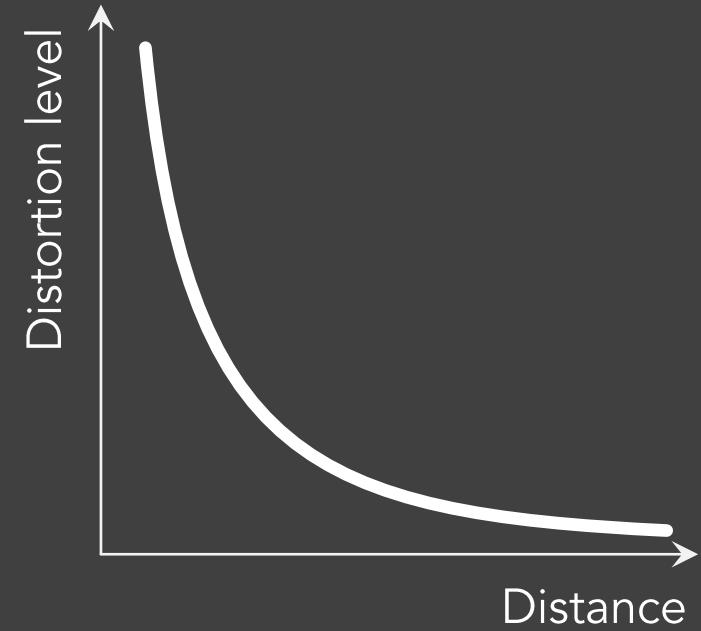


Depth Variation

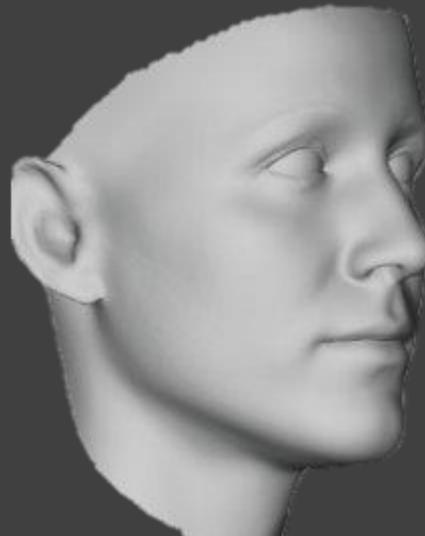
$$\Delta d$$



Perspective



# Weak-perspective Projection



Perspective



Weak-perspective



Depth Variation



$D_{\text{close}}$



$D_{\text{far}}$



Camera-to-Subject Distance

$$D_{\text{far}} \gg \Delta d$$



# Manipulate Viewpoint and Lens

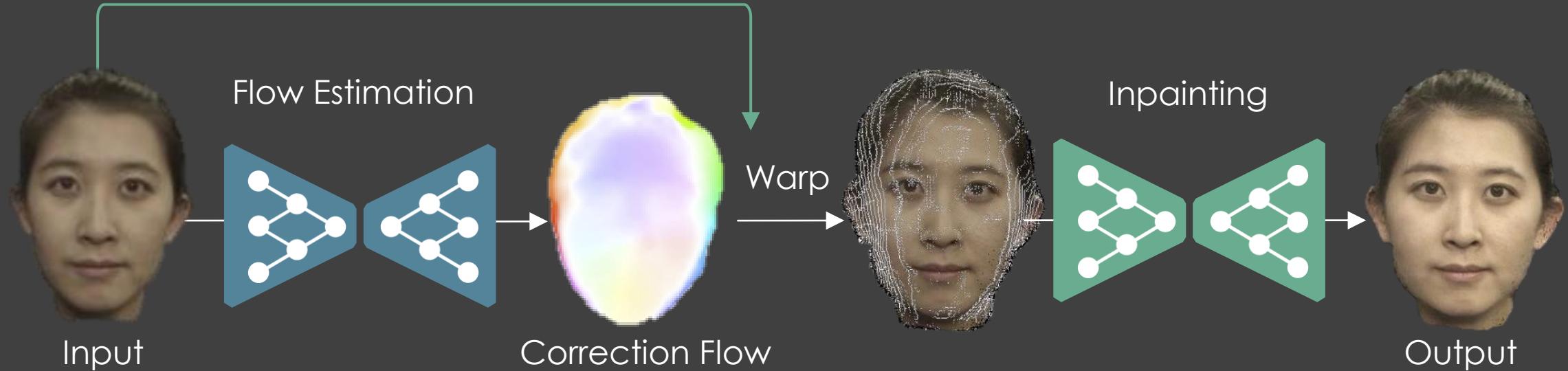
Perspective



Weak-perspective



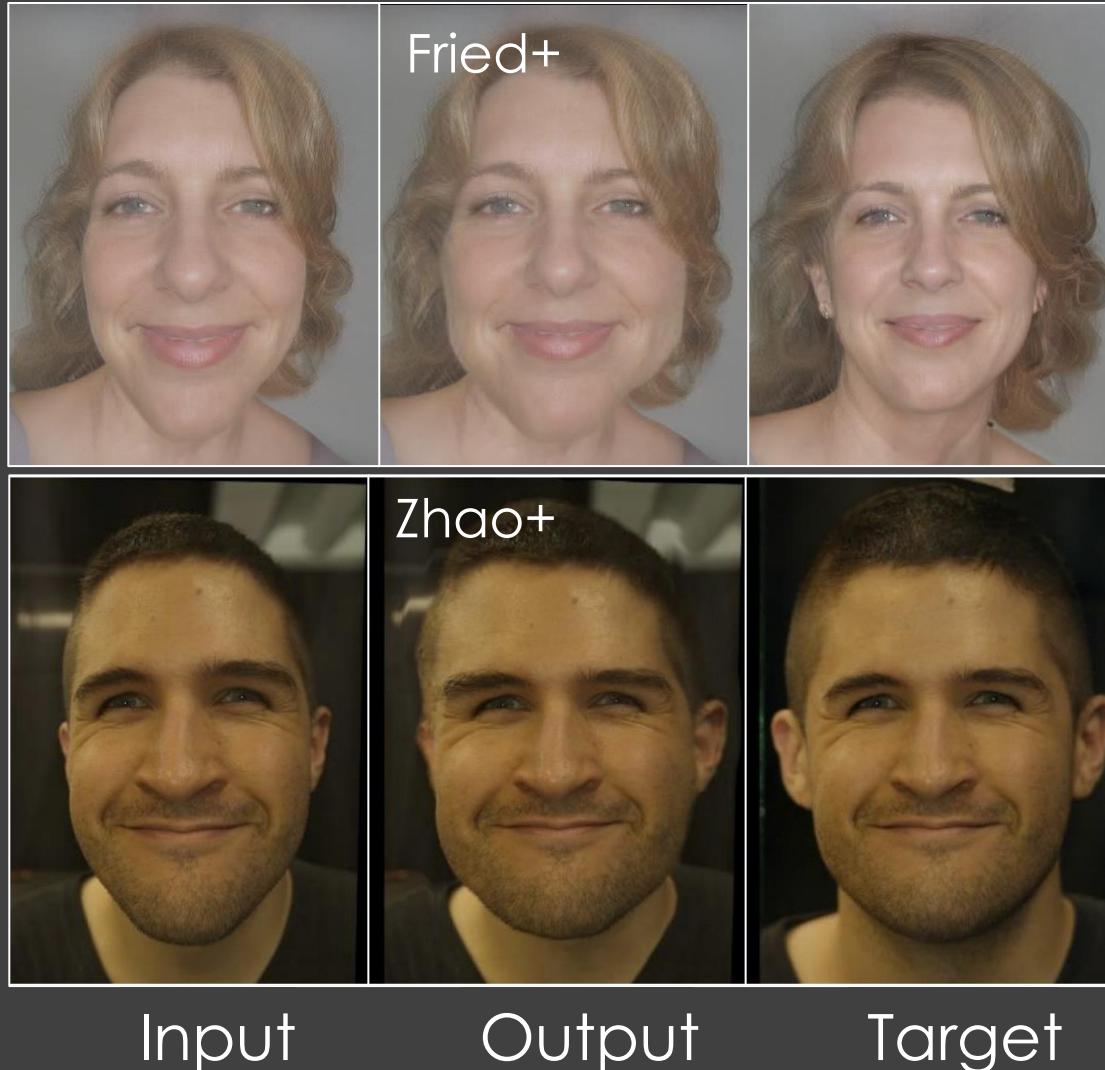
# Existing Methods – Warping-based



Fried et al, SIGGRAPH'16

Zhao et al, ICCV'19

# Limitations of Existing Methods



- ▶ **Flow warping only repeats existing pixels**
  - ▶ CANNOT reveal occluded regions
    - ▶ Invisible ear, cheek, neck ...
  - ▶ CANNOT deal with serious distortion
    - ▶ When camera-to-face distance is 20–40cm
  - ▶ Not 3D-aware
    - ▶ Face shape is flawed
- ▶ **Learning-based method (Zhao+) is worse**
  - ▶ Require a lot of training data
  - ▶ Hard to generalize
  - ▶ CANNOT continuously change

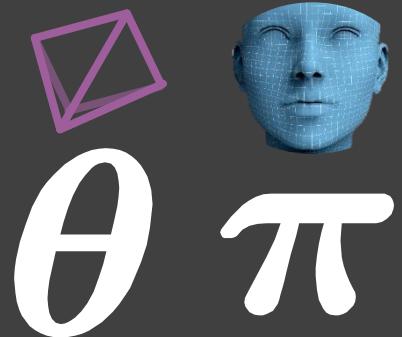
# Optimization-based Factorization

Input:  $I$   
**Single** Image



$$\mathcal{R}^{-1}(I)$$

.....  
**Inverse**  
Rendering



$$\mathcal{R}(\theta, \pi)$$

**Forward**  
Rendering

Reconstruction      Input



# Optimization-based Factorization

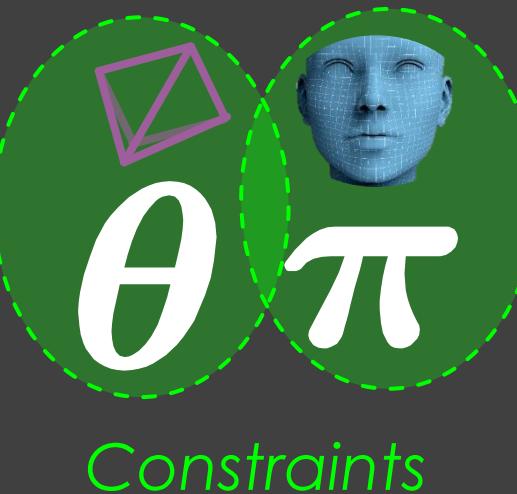
**Challenge:** ill-posed/unconstrained

Input:  $I$   
Single Image



$$\mathcal{R}^{-1}(I)$$

Inverse  
Rendering



$$\mathcal{R}(\theta, \pi)$$

Forward  
Rendering

Reconstruction      Input



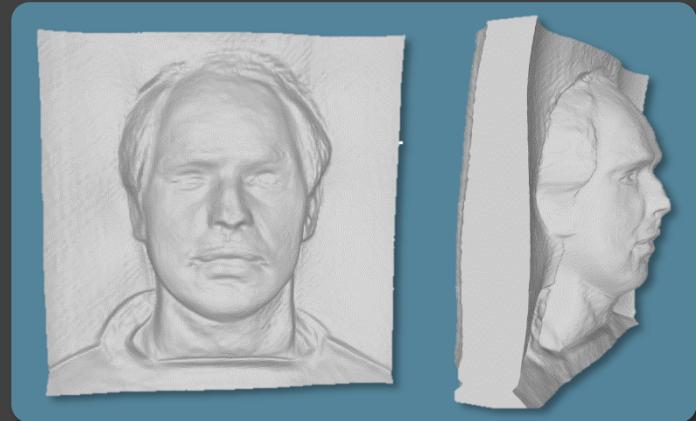
+Geometry Loss

# Ambiguity of Parameters

Many **combinations**  
resemble input image



...



Face Shape

Flat



Small

Camera-to-Subject Distance

Large



Small

Focal Length

Large



# 3D GAN Prior as Face Constraint

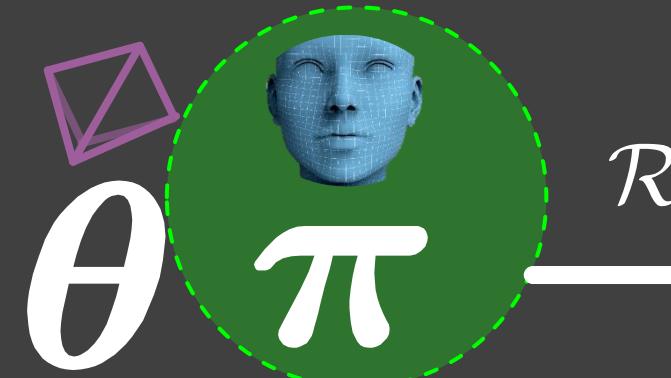
Single Image



$$\mathcal{R}^{-1}(I)$$

**$I$**

Inverse  
Rendering



Unconstrained

$$\mathcal{R}(\theta, \pi)$$

Reconstruction



3D GAN

Random  
noise

Generator

Implicit  
Representation



**Probabilistic Representation**

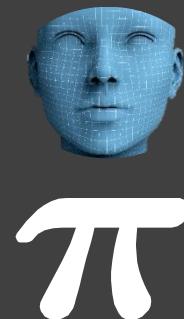
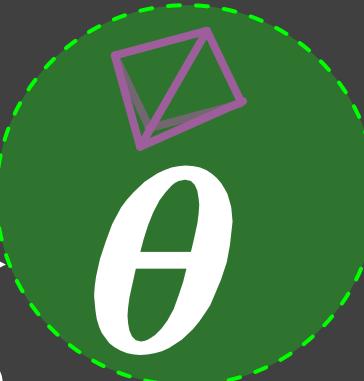
# Camera Regularization (CR)

Input:  $I$   
Single Image



$$\mathcal{R}^{-1}(I)$$

Inverse  
Rendering



*Unconstrained*

$$\mathcal{R}(\theta, \pi)$$

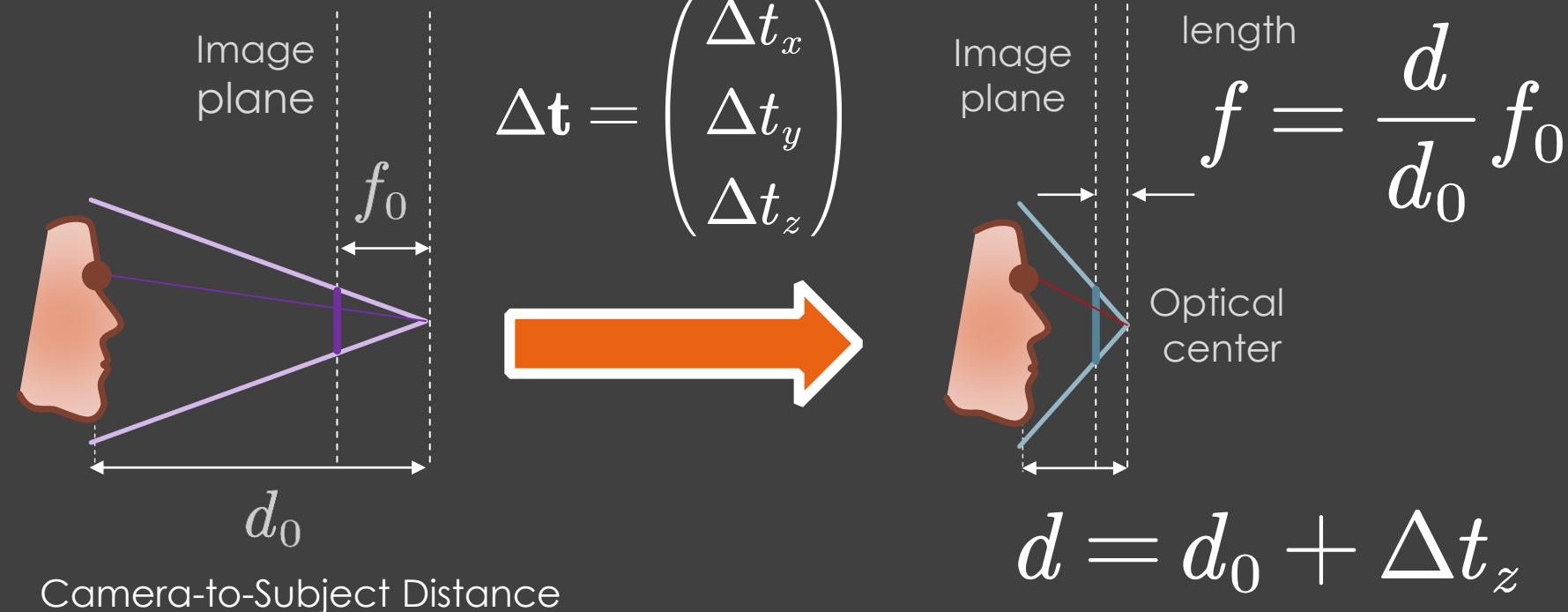
Reconstruction



3 Strategies

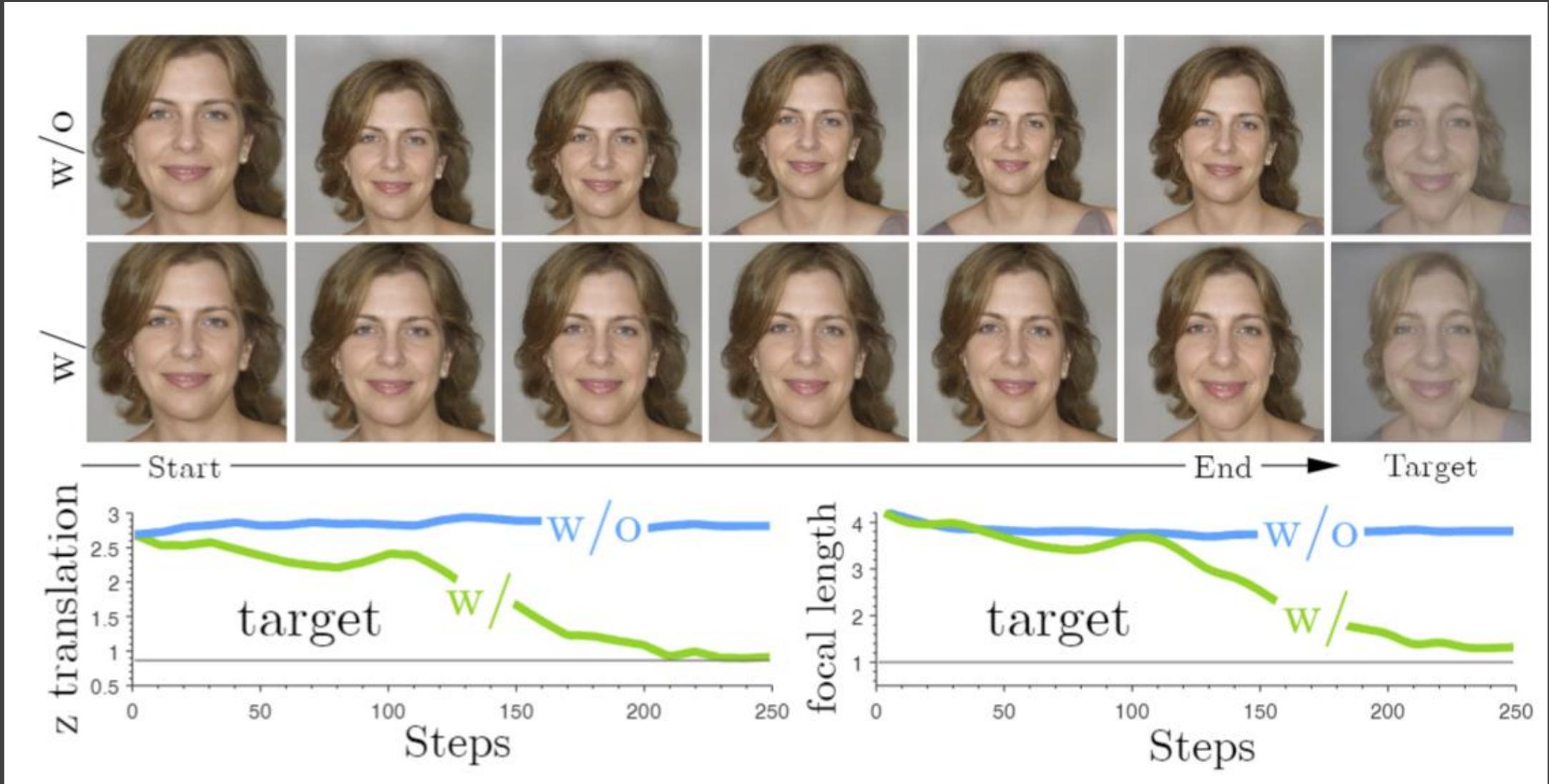
# CR 1: Focal Length Re-parameterization

- Focal Length  
(simplified approximation)



**Motivation:** Reduce unknown parameters and decouple

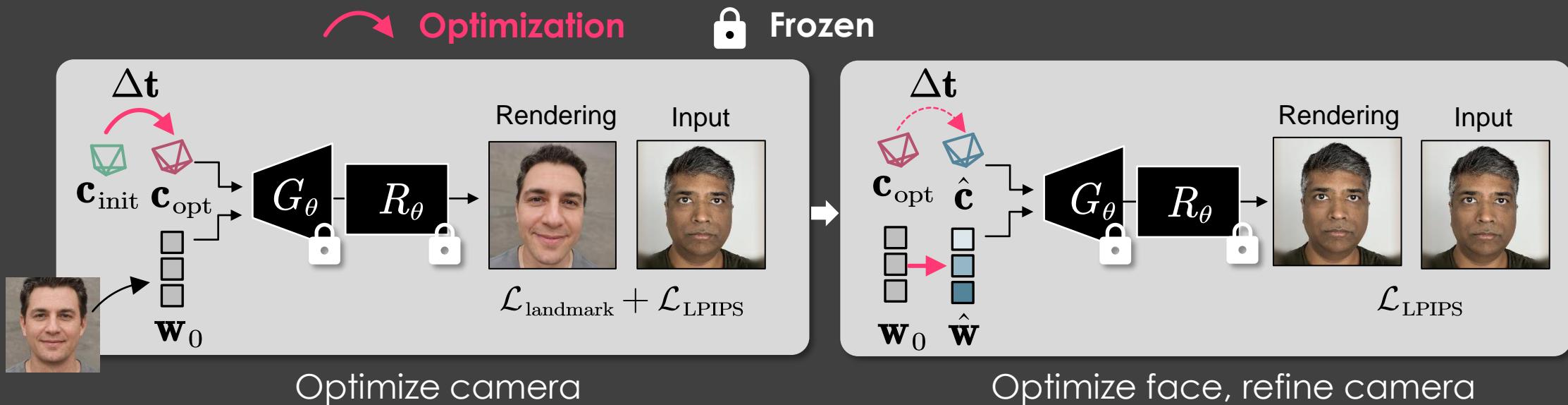
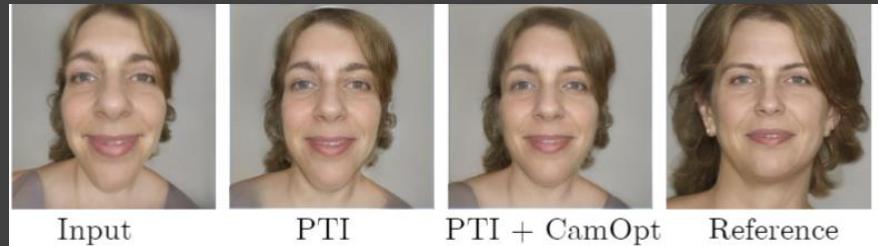
# CR 1: Focal Length Re-parameterization



# CR 2: Optimization Scheduling

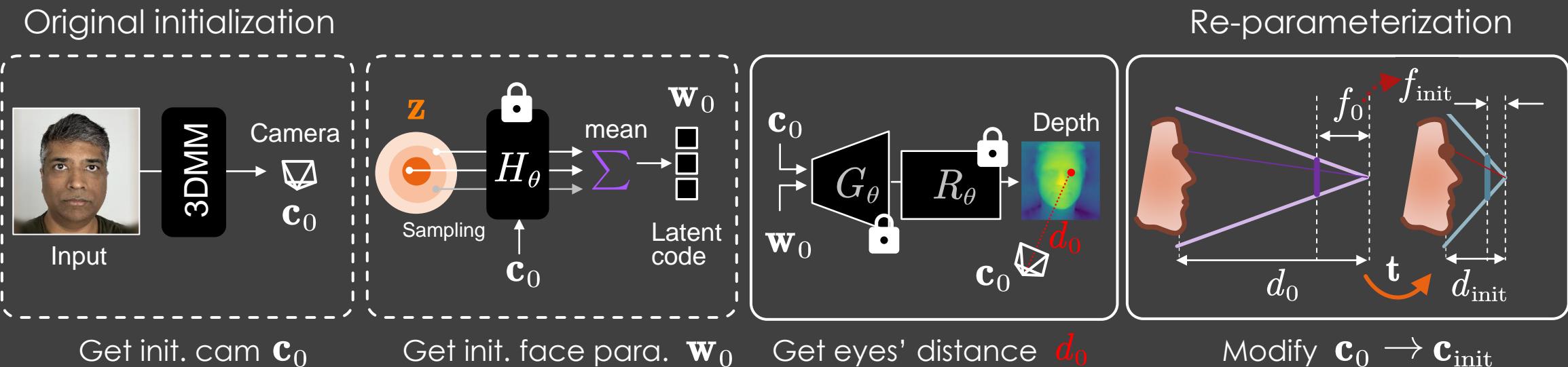


**Motivation:** Face is **easier** to fall into **sub-optimum** than camera



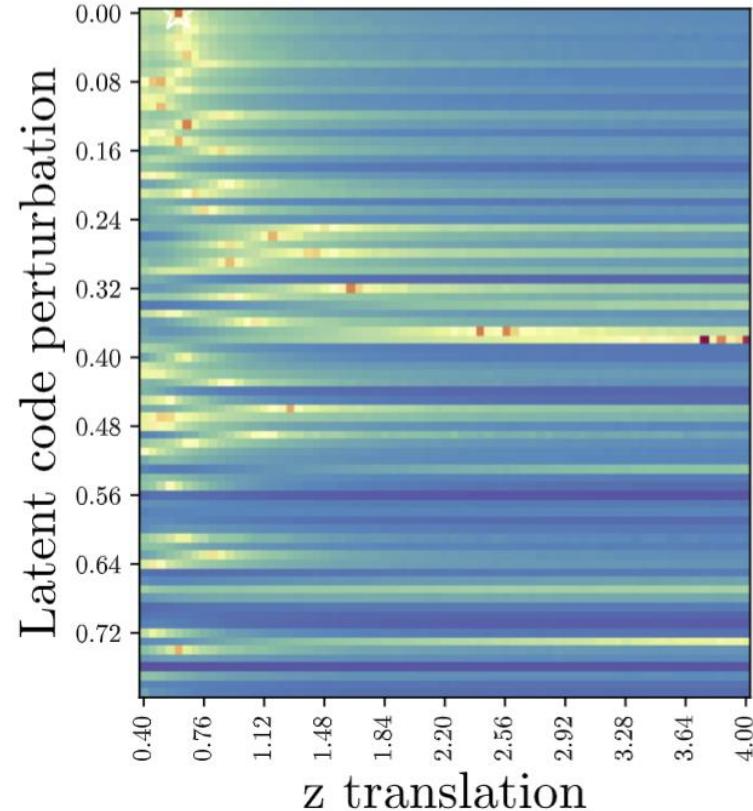
# CR 3: Better Initialization

Start from a close-up camera position

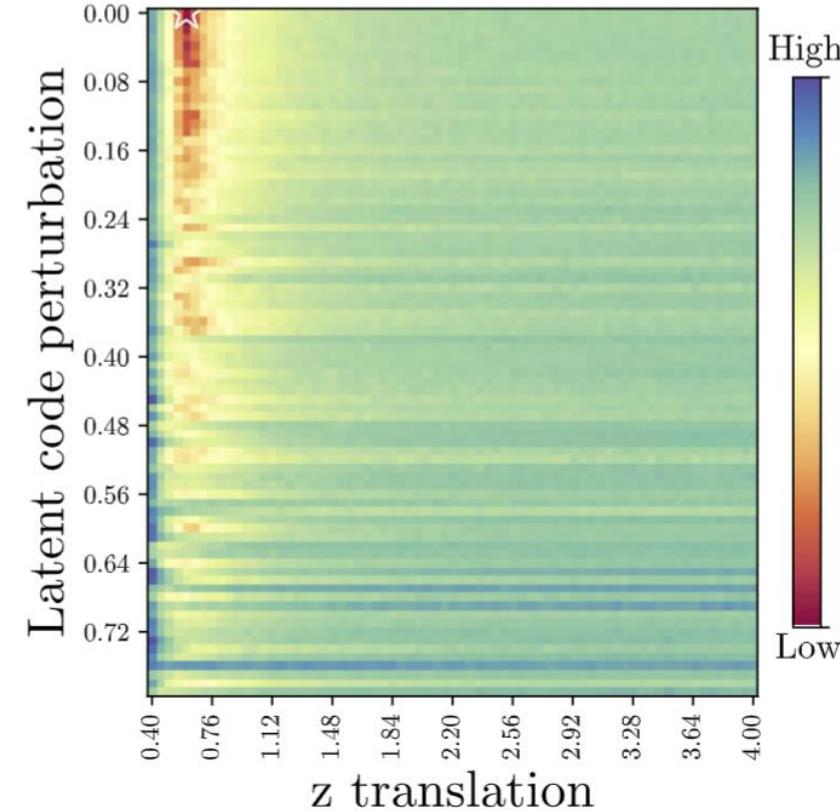


# Ambiguity Caused by Loss

Pixel loss is **very sensitive** to pixel change



(a)  $L_2$  loss



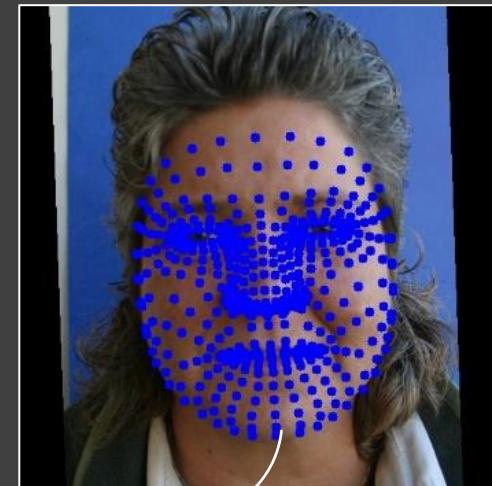
(b) Landmark loss

# Geometric Regularization

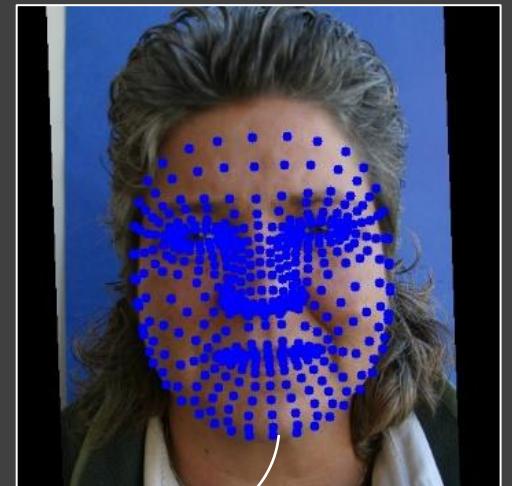
Uncertainty-based Loss

$$\sum_{i=1}^{\|\mathcal{M}\|} \left( \underbrace{\log(\sigma_i^2)}_{\text{Uncertainty term}} + \frac{\|m_i - m'_i\|_2^2}{2\sigma_i^2} \right)$$

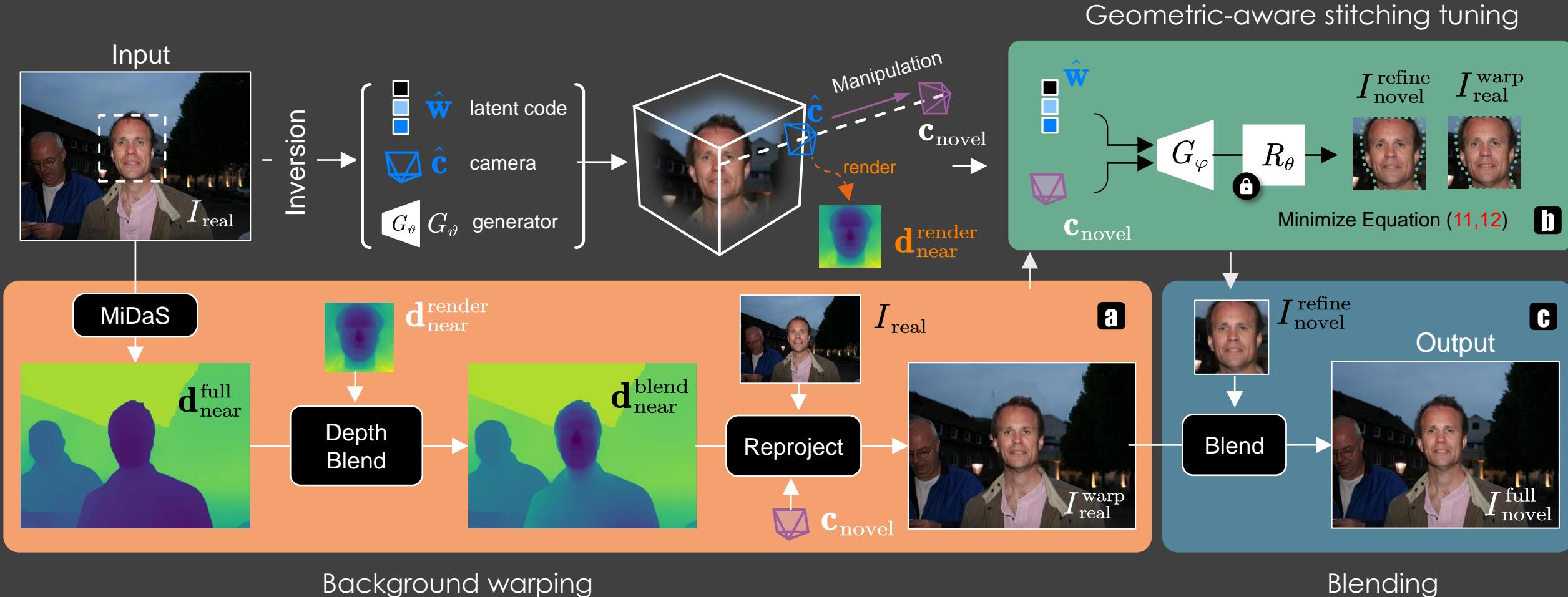
prediction



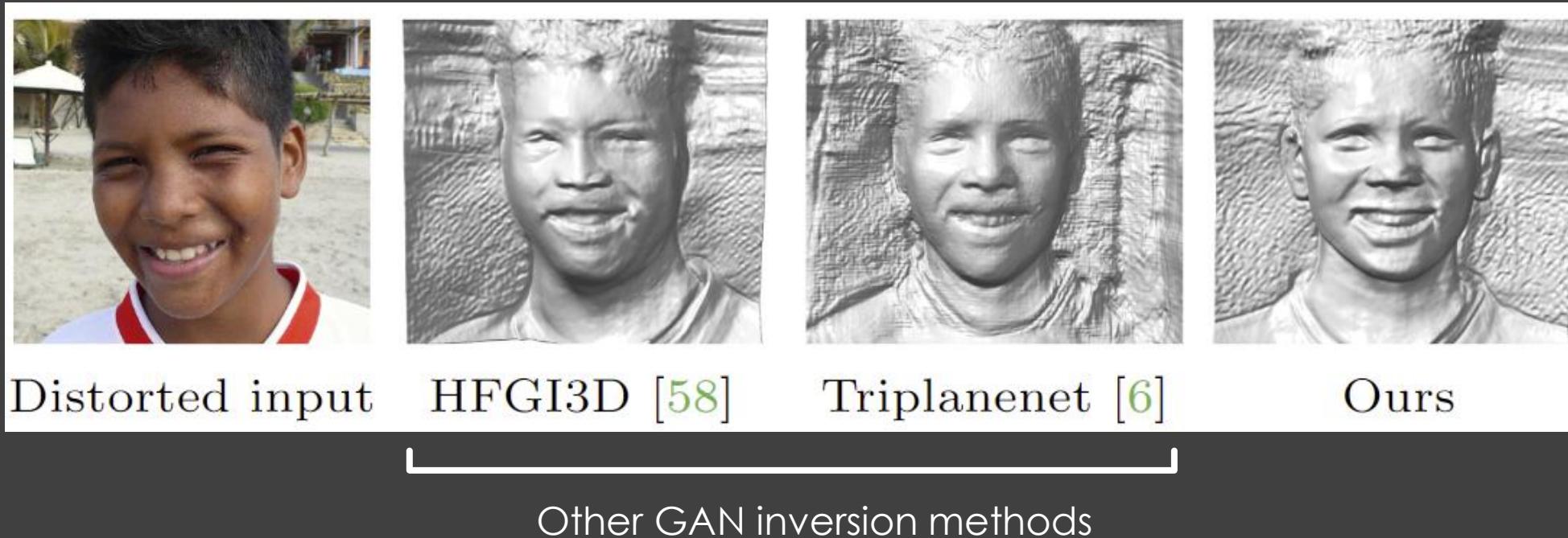
input



# Extensions for Full-frame Image



# Results – Mesh

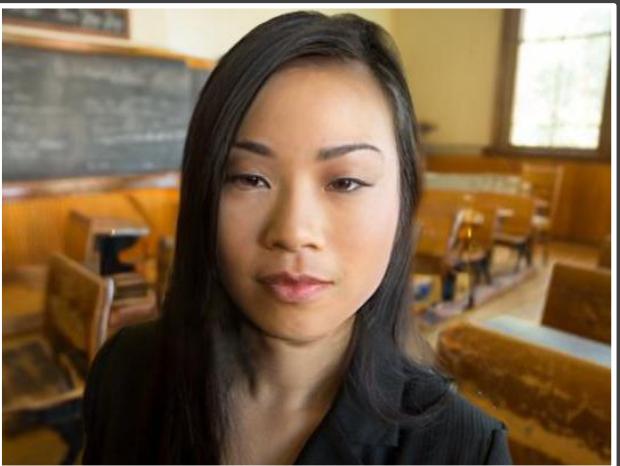


# Results

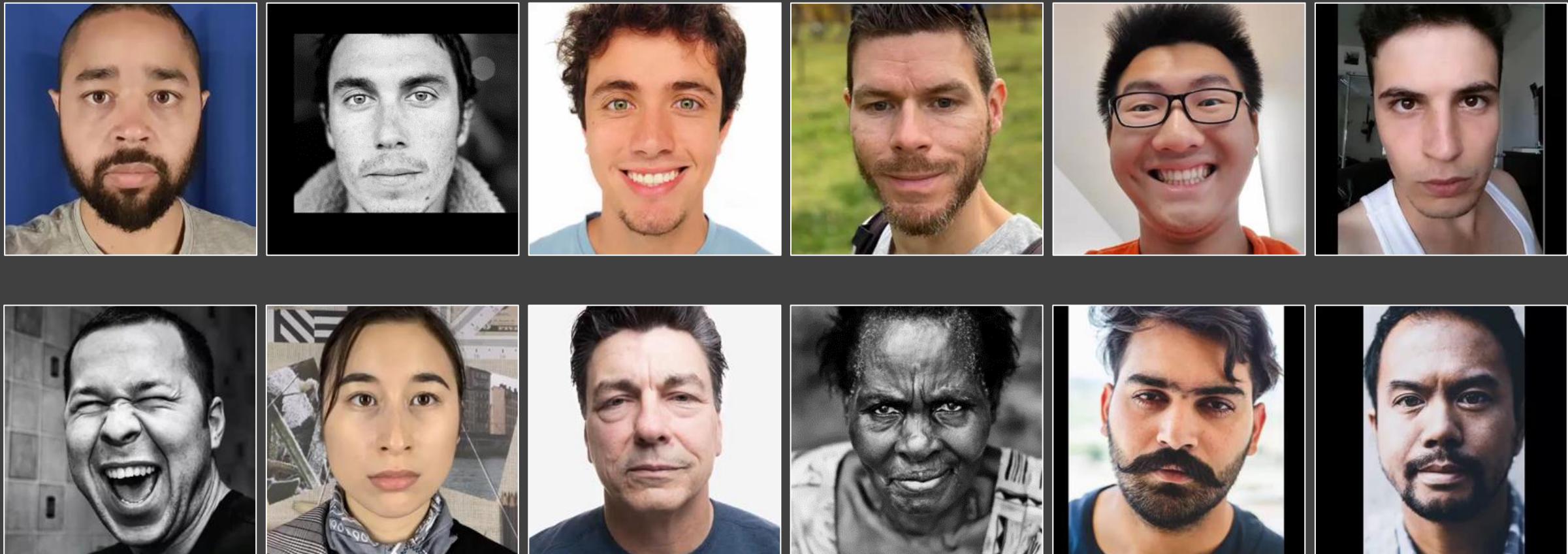
Input



Output



# Results – Continuous Manipulation



# Results – Comparison

Stretch-like



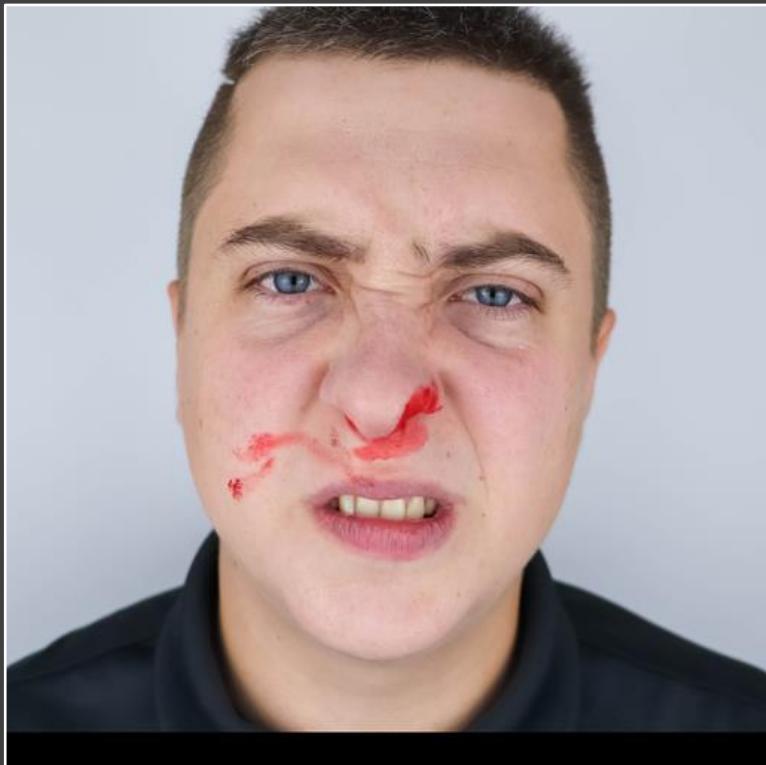
Fried et al, SIGGRAPH'16

3D geometric consistent

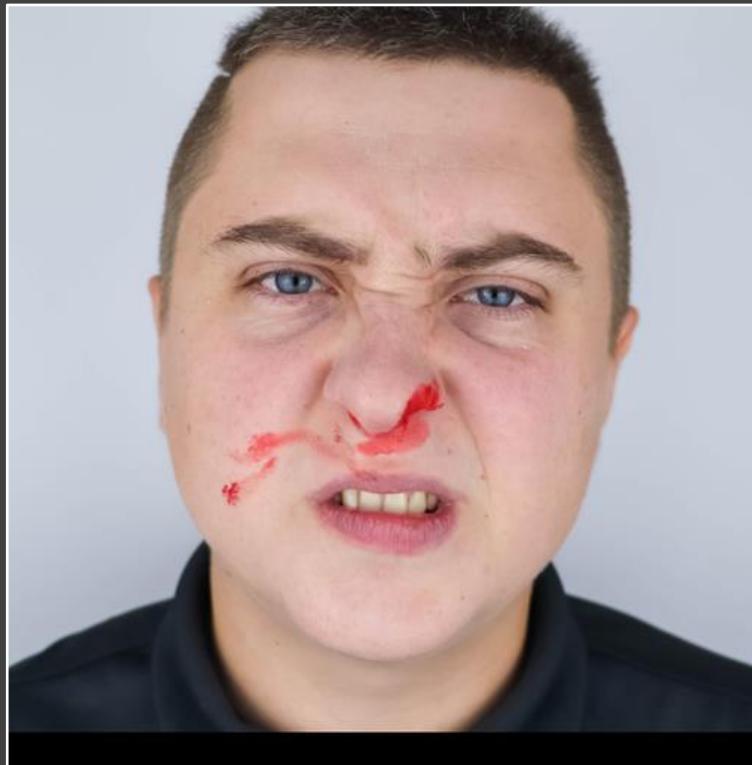


Ours

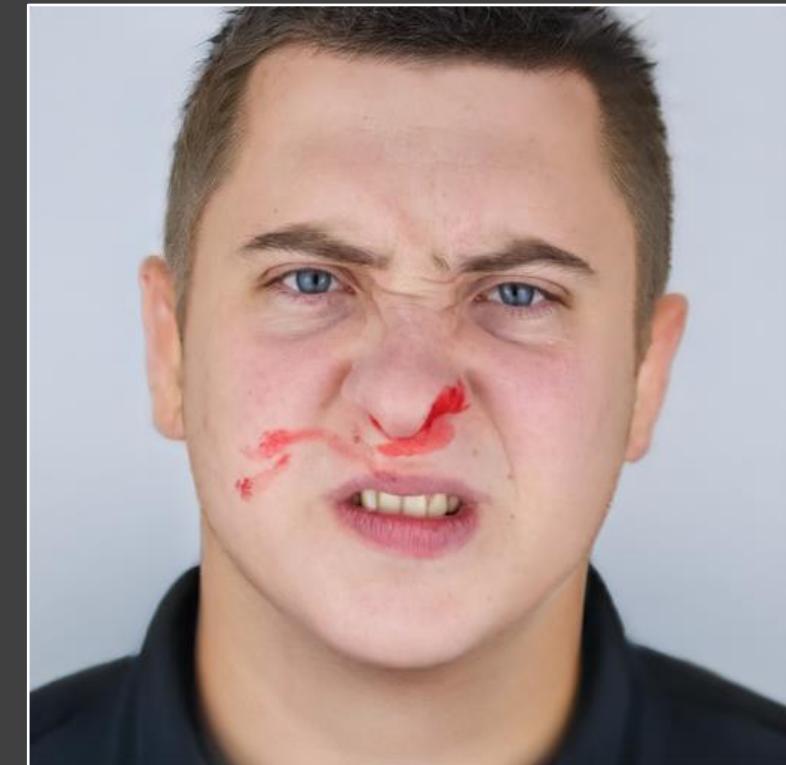
# Results – Comparison



Input



Fried et al, SIGGRAPH'16



Ours

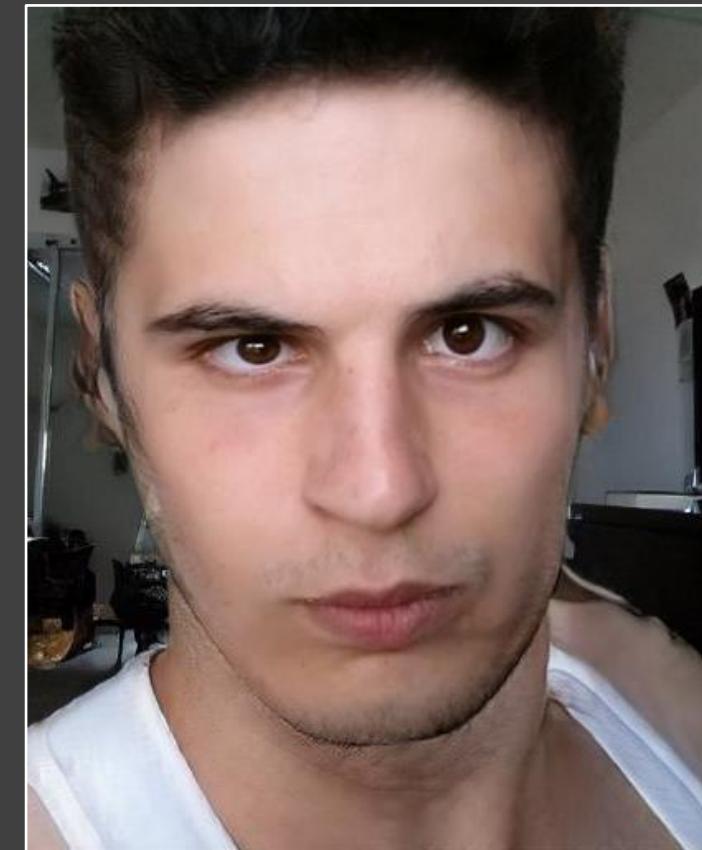
# Results – Comparison



Input



Fried et al, SIGGRAPH'16



Ours

# Results – Comparison



Input



Fried et al, SIGGRAPH'16



Ours

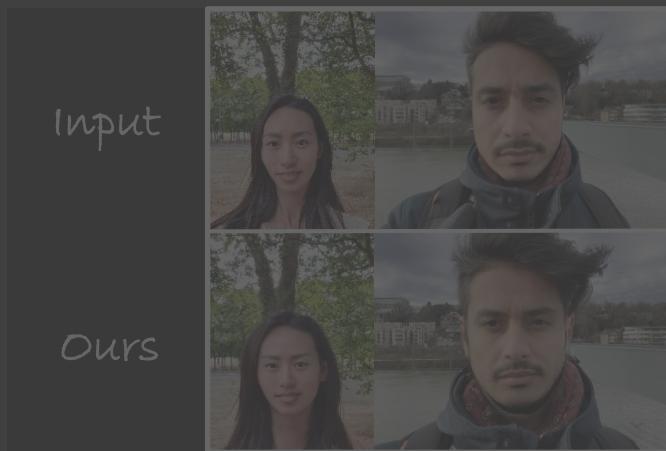
## Dolly Zoom



# Dolly Zoom



# Viewpoint + Lens



Perspective Distortion Correction

# Background Background

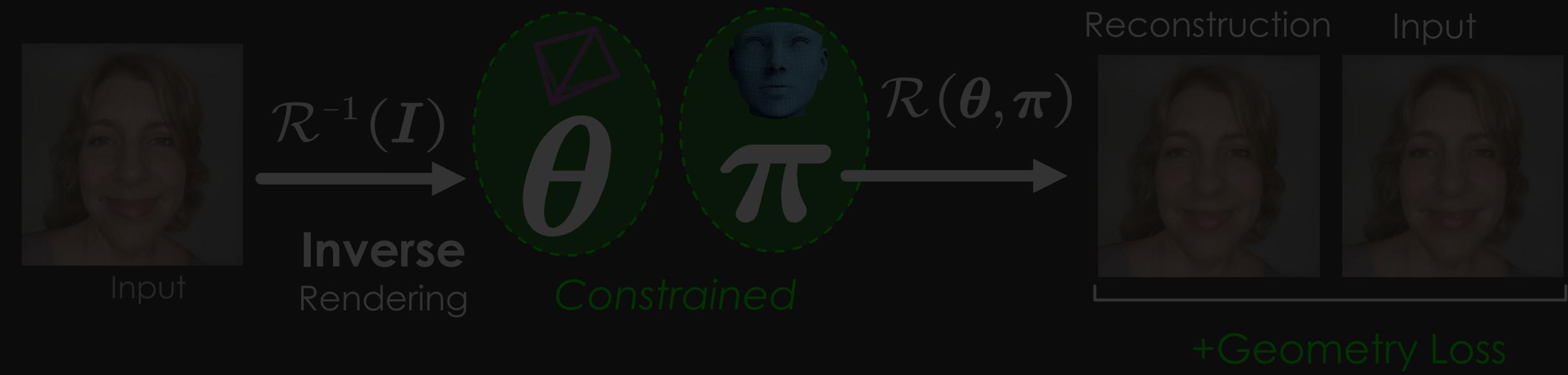


Matting by Generation

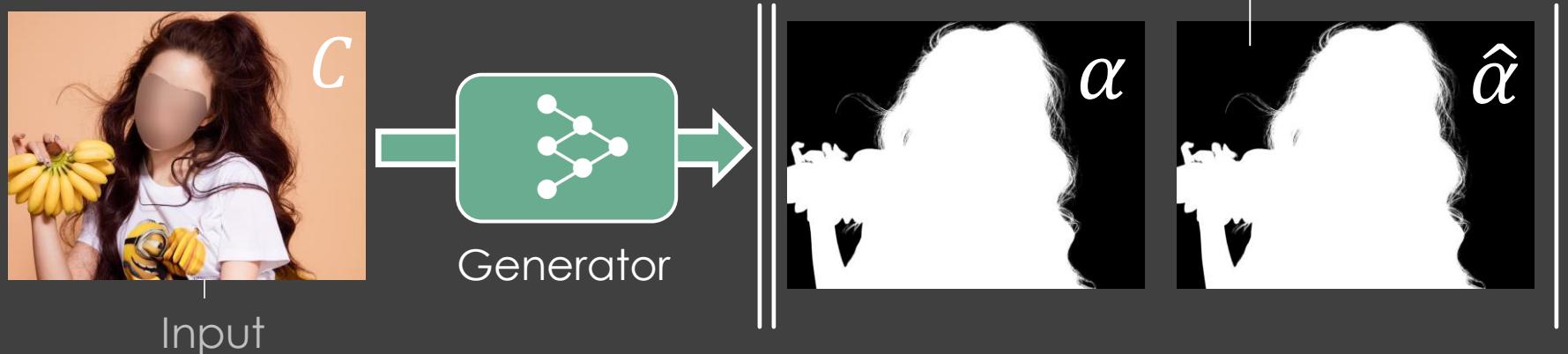
Wang et al, SIGGRAPH 2024

# Harness Pre-trained Generative Models

**Optimization-based:** no labels required



**Learning with Labels:** imperfect labels



# Manipulate Background



Background Gallery

# Factorization Problem

Input:  $I$   
Single Image

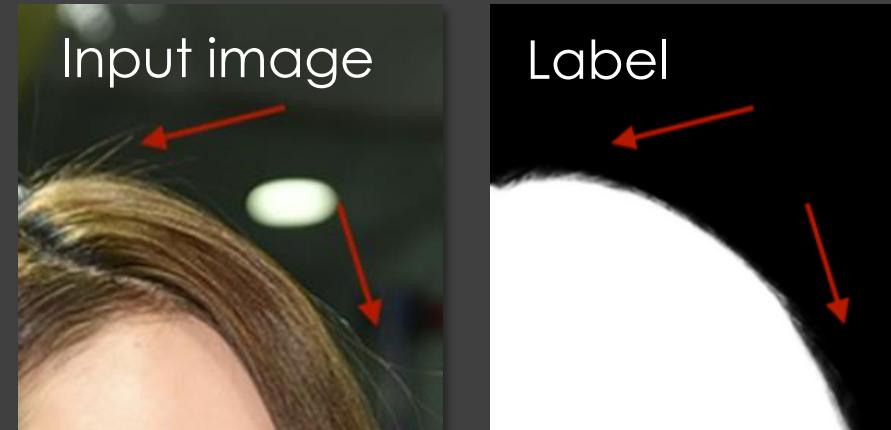


Re-rendering  
 $\alpha F + (1 - \alpha)B^*$



46

# Learning with Labels



Ke et al, MODNet, AAAI'22  
Li et al, P3M, MM'22  
Ma et al, ViTAE-S, IJCV'23

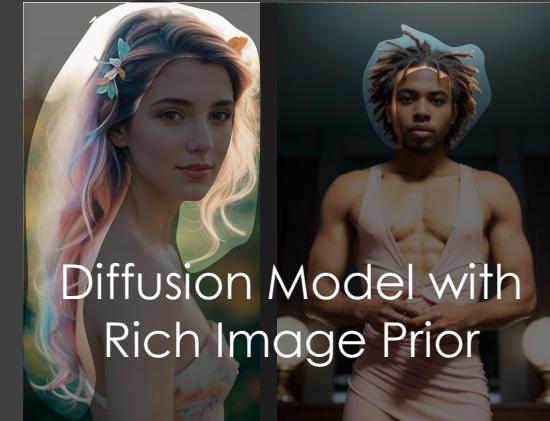
Poor label quality

# Limitations of Existing Methods

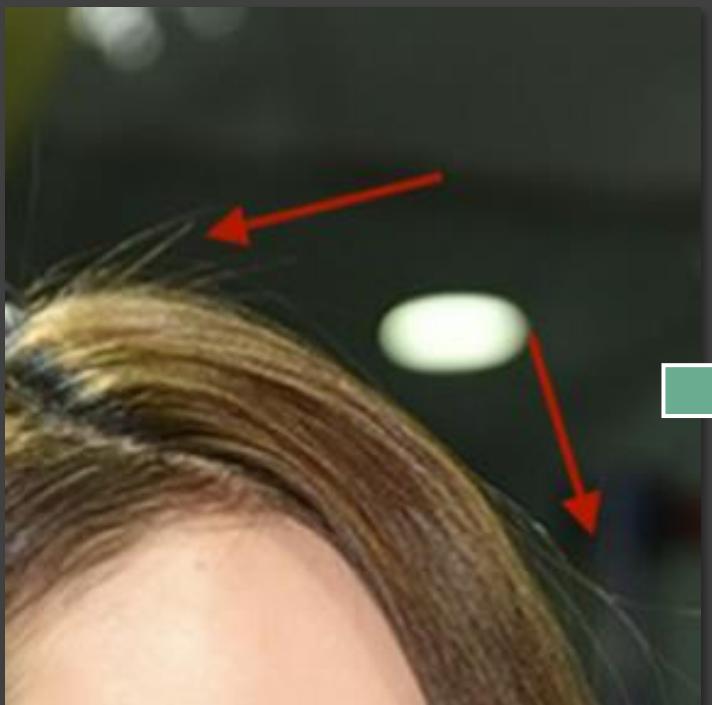


# Generative Diffusion Prior

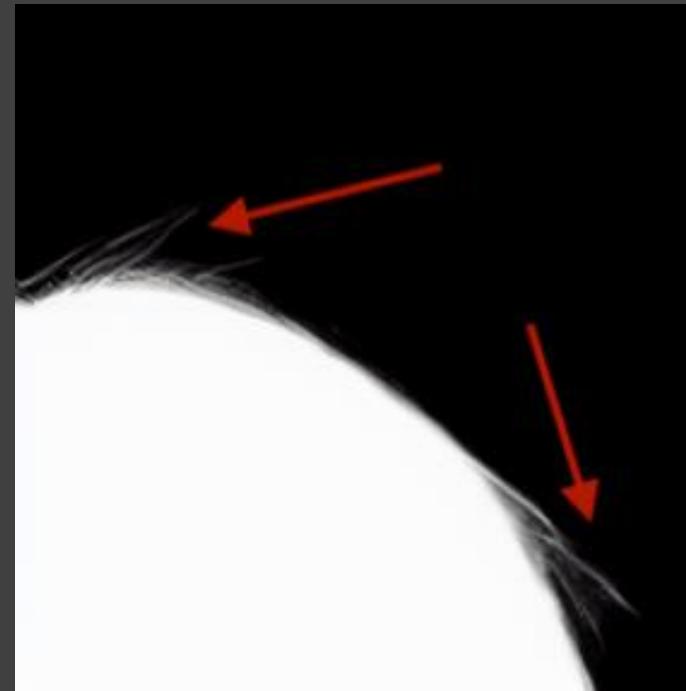
Generative Prior for Regularization



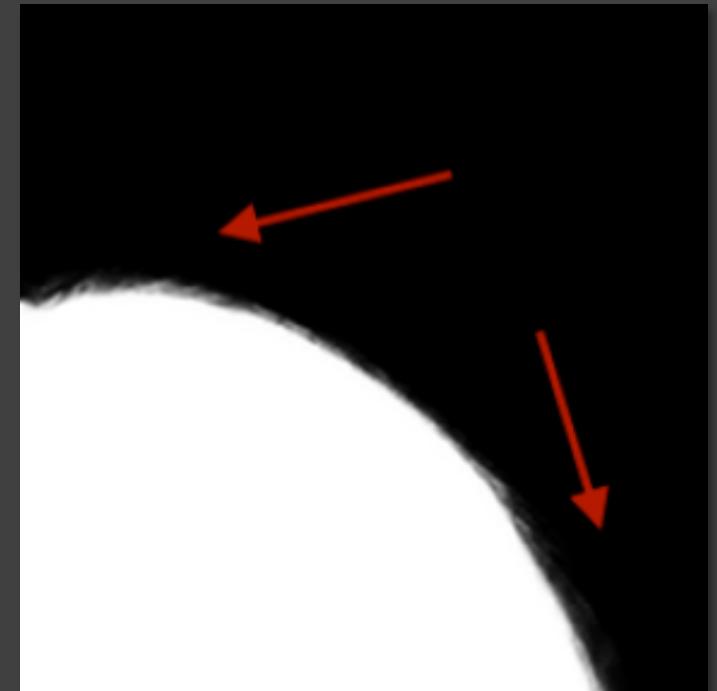
Diffusion Model with Rich Image Prior



Input training image

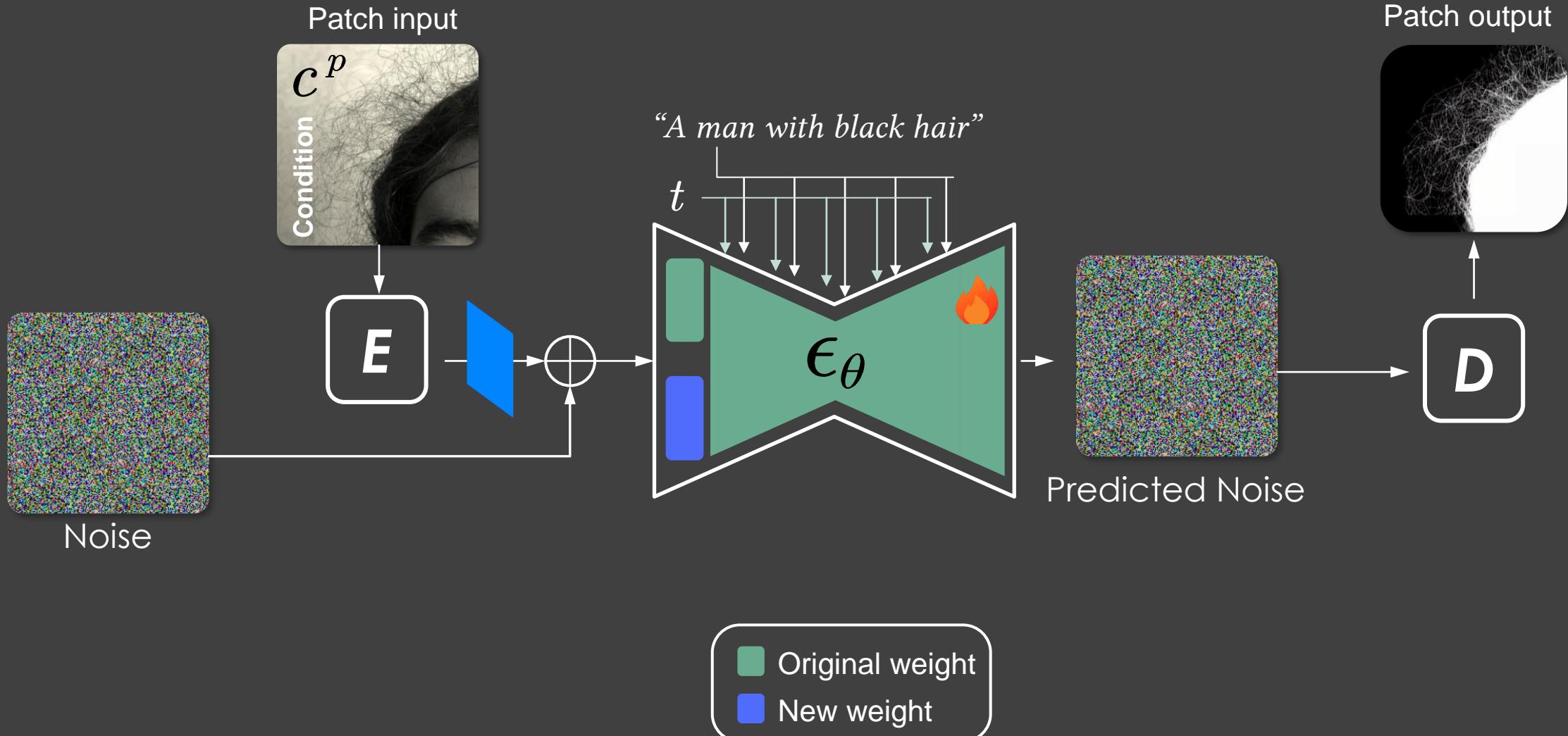


Output

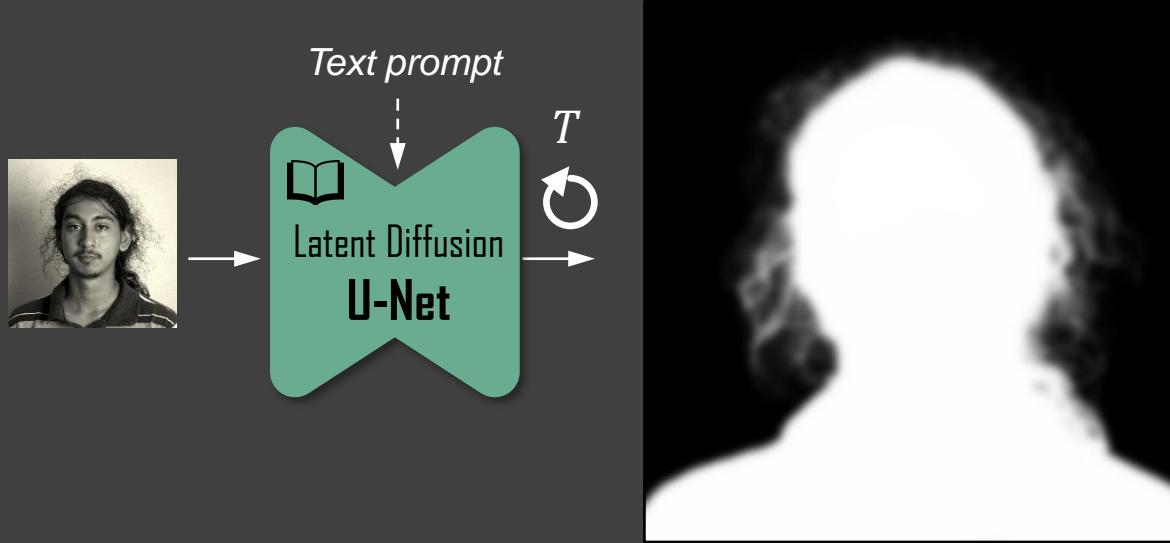


Label

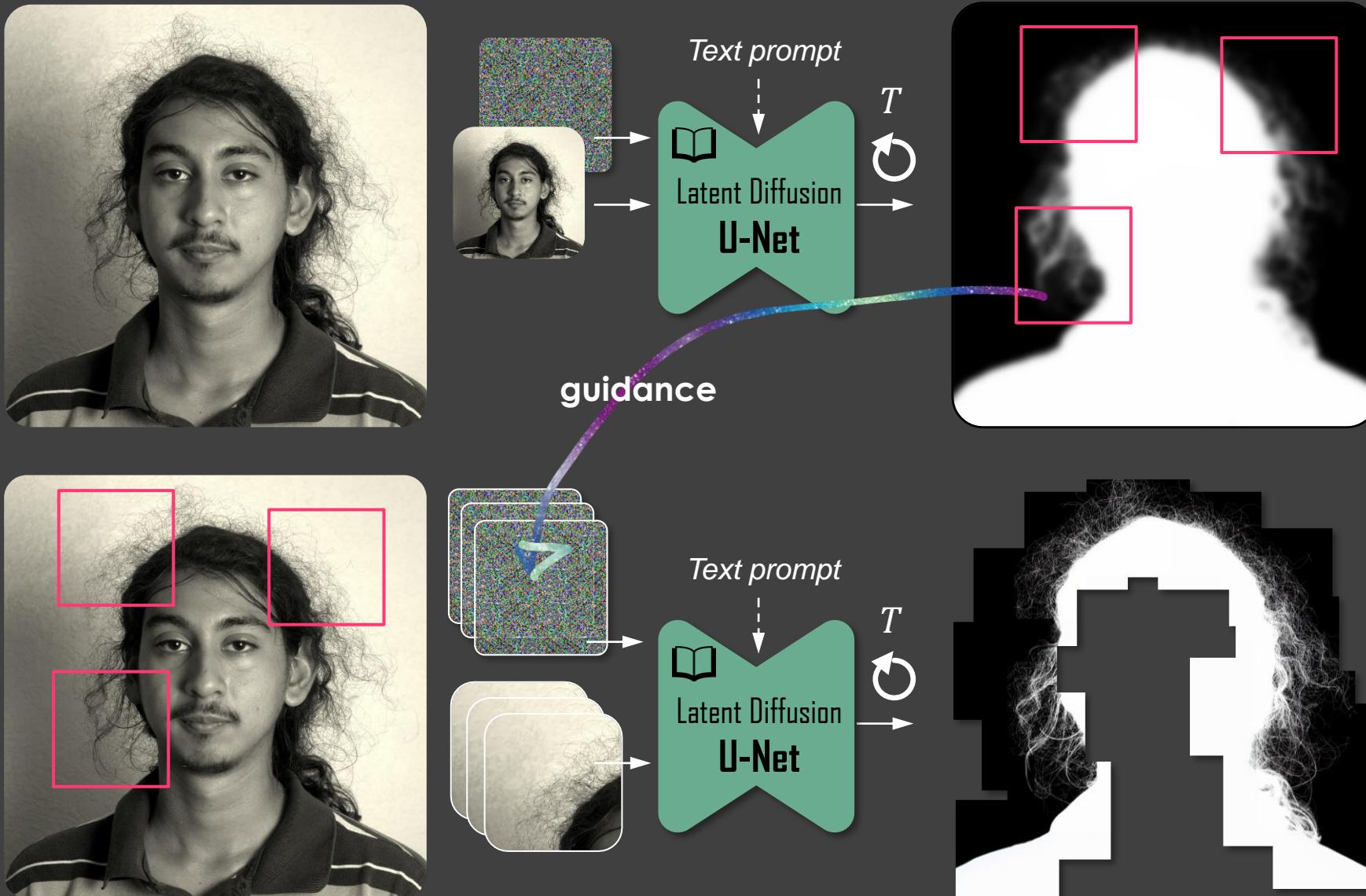
# Repurposing Latent Diffusion Model



# Challenge of Processing HR Images

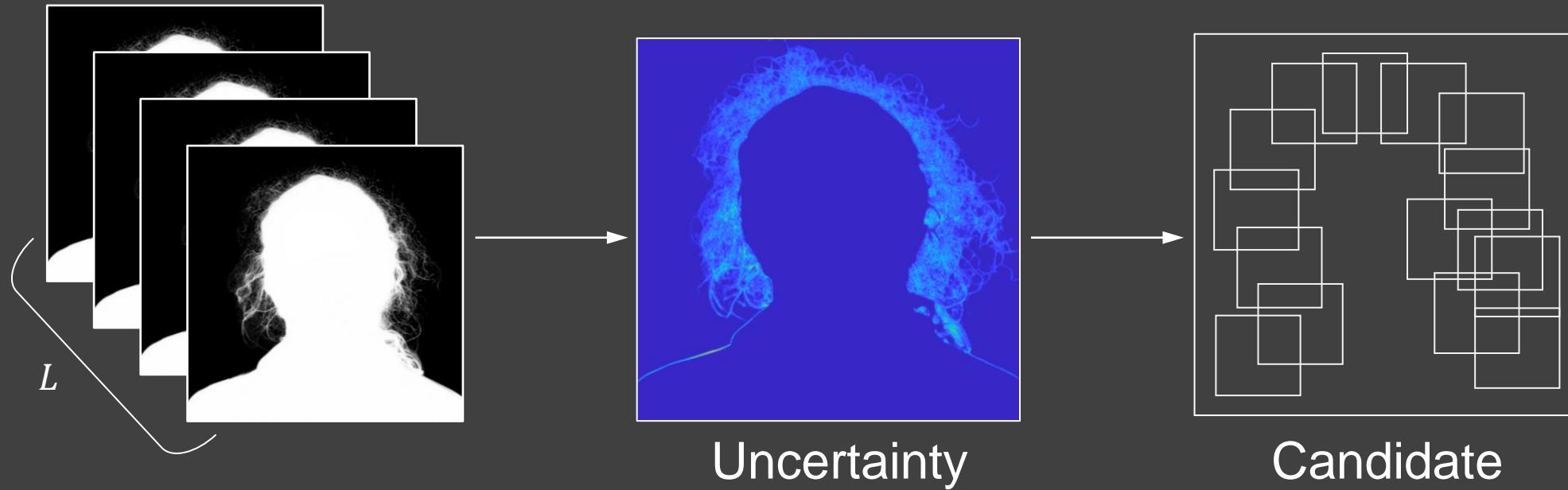


# Pipeline for Processing HR Images

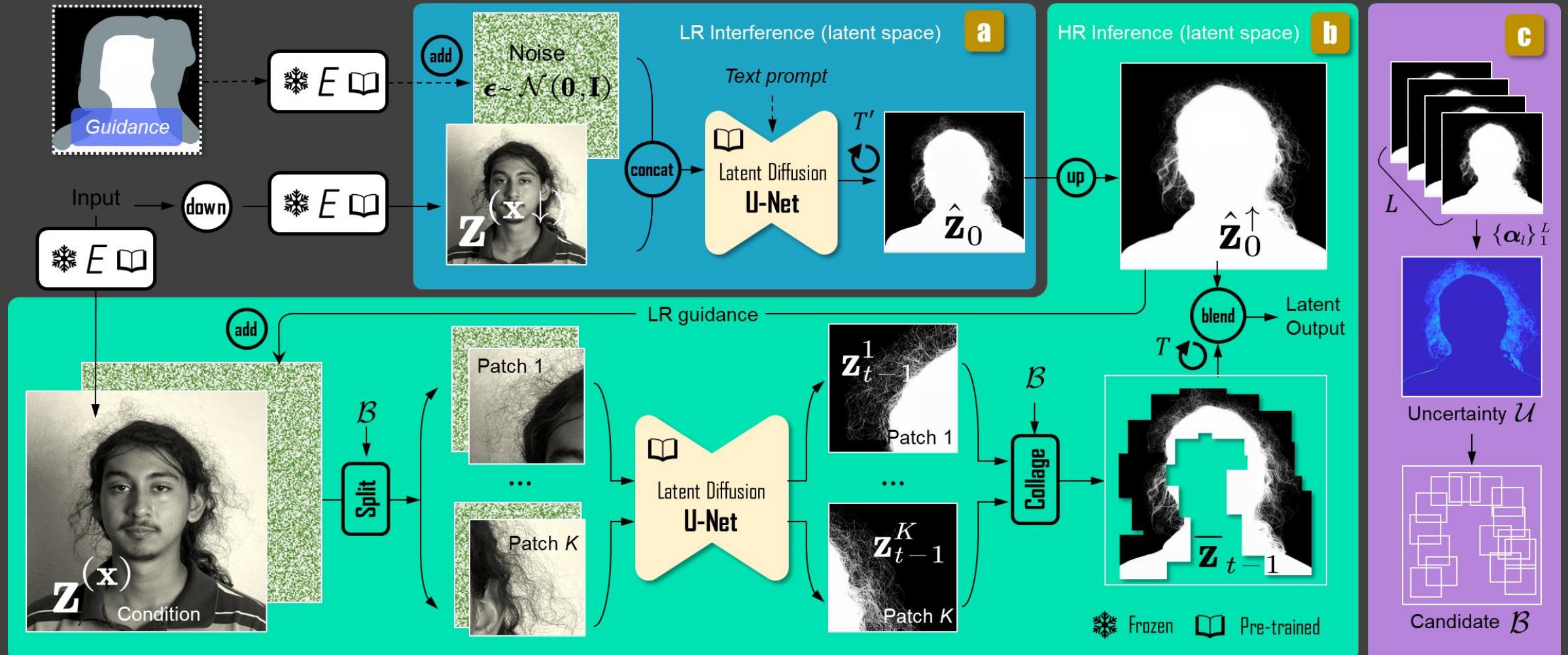


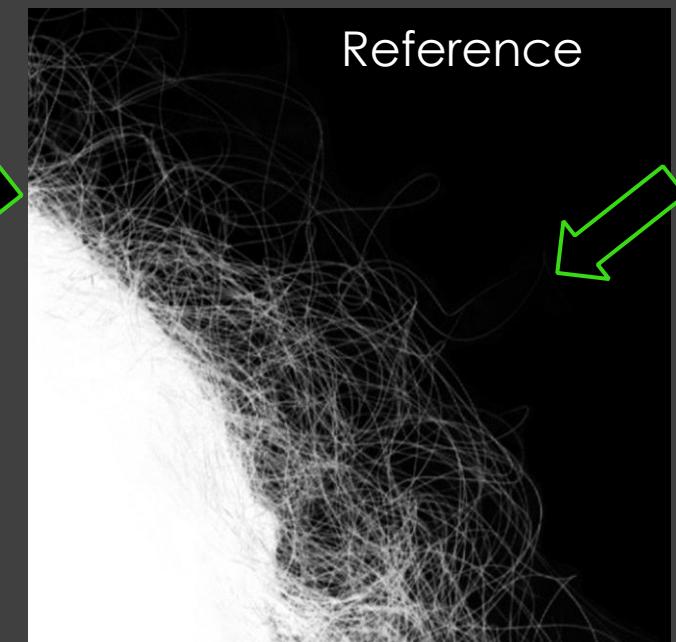
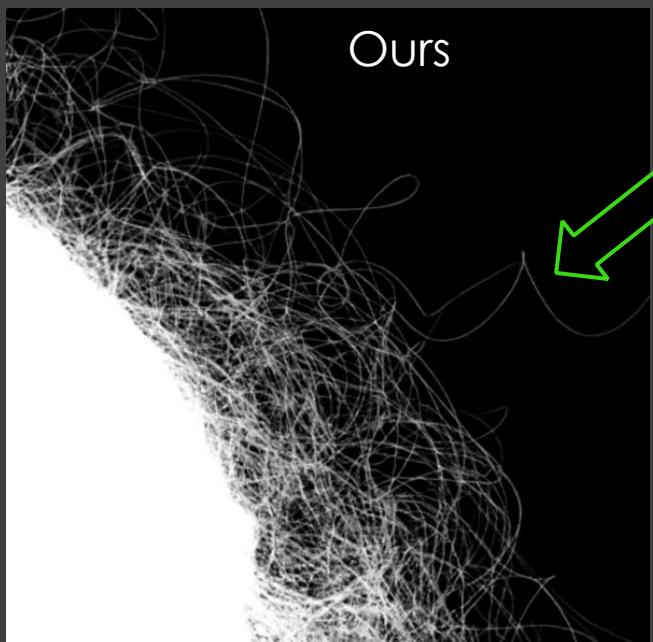
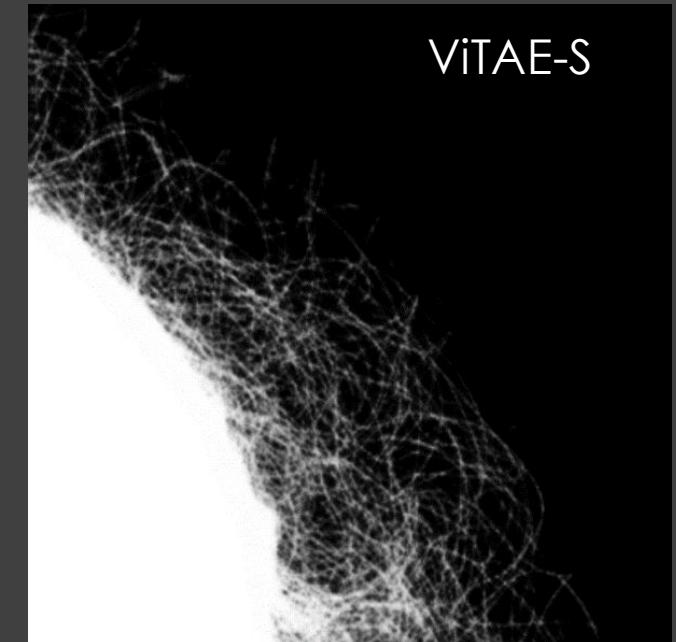
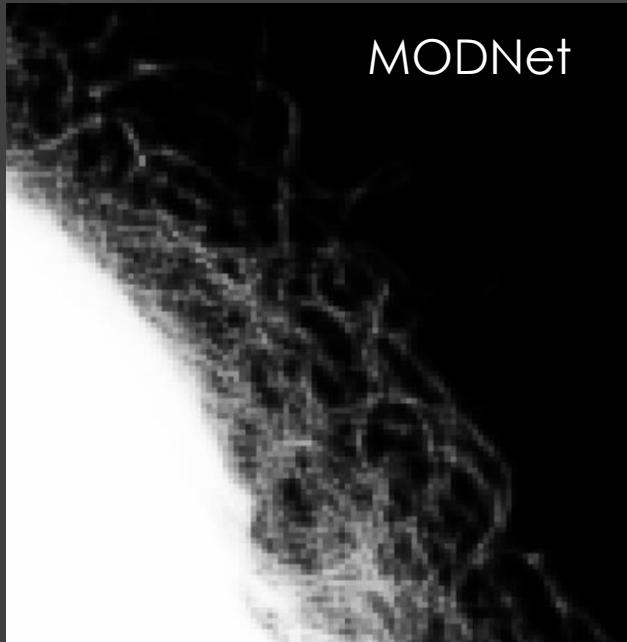
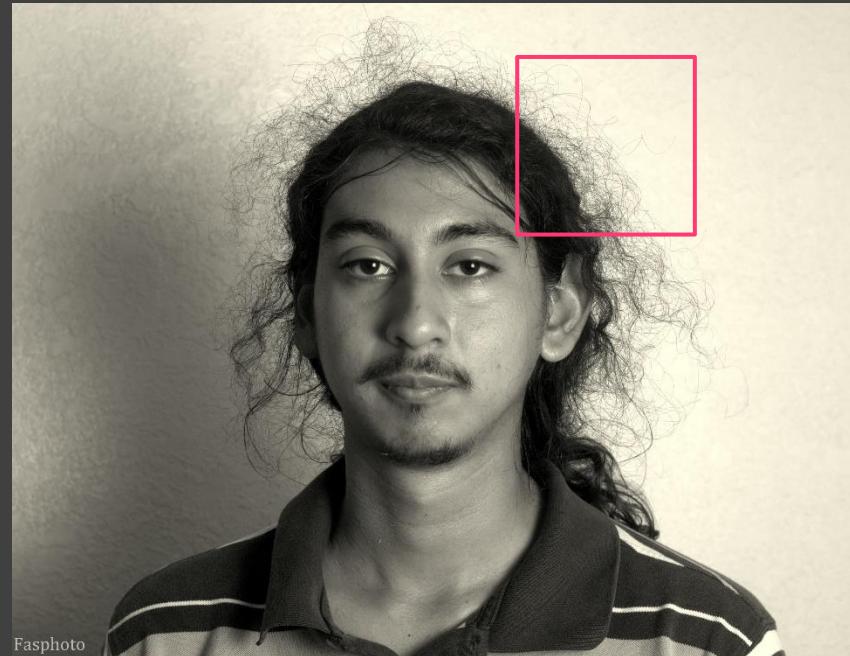
# Pipeline for Processing HR Images

Get potential areas by uncertainty



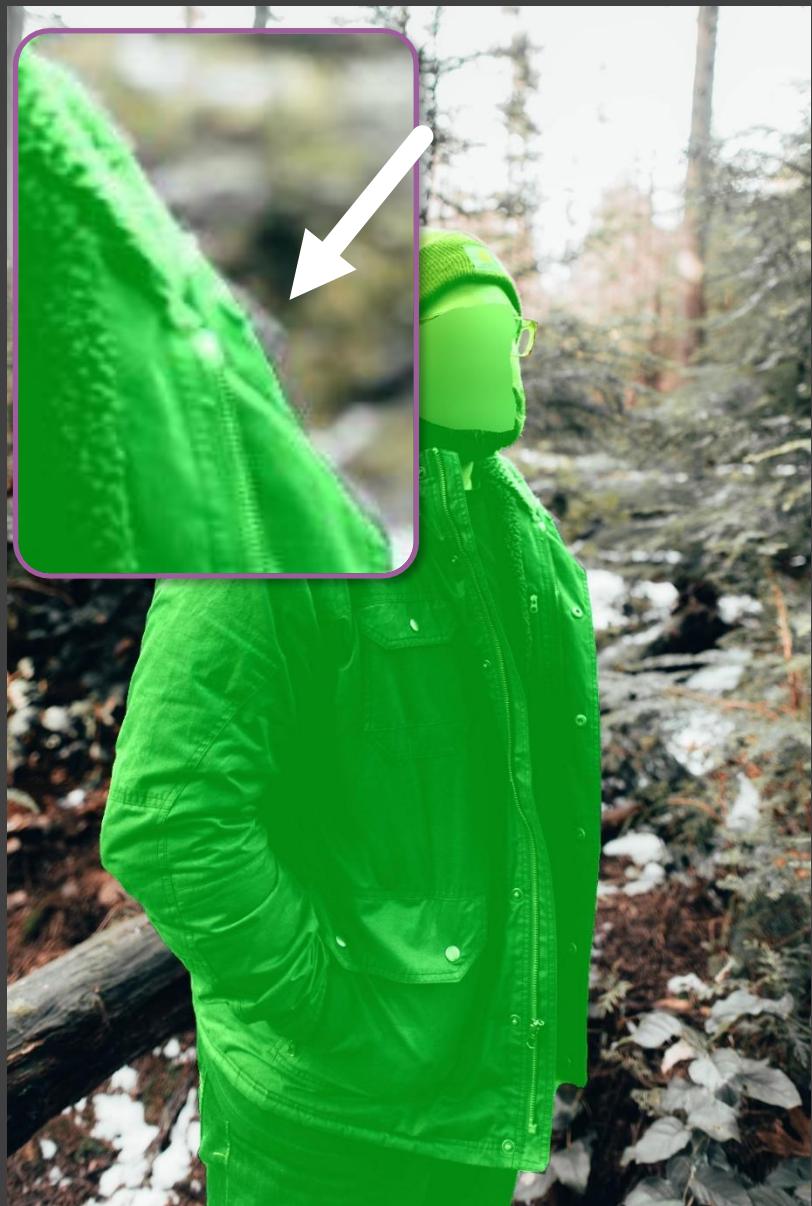
# Full Pipeline



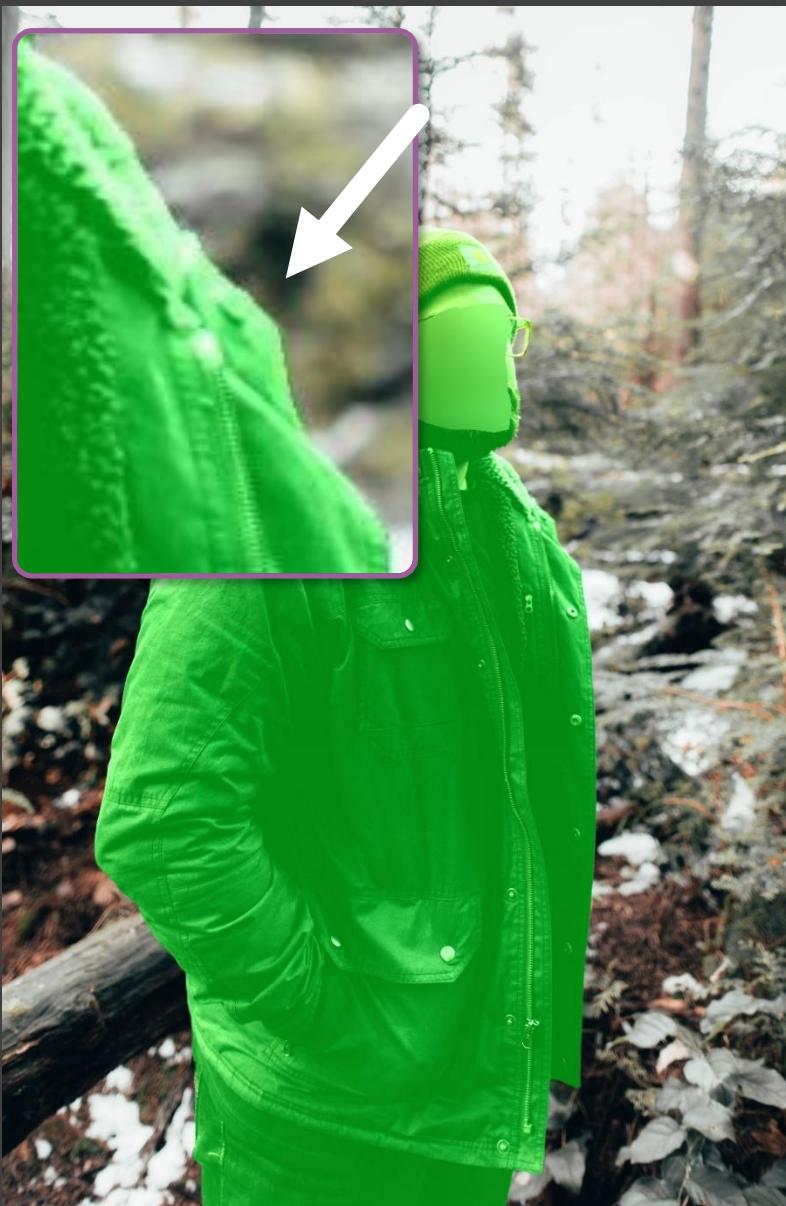


Ke et al, MODNet, AAAI'22  
Li et al, P3M, MM'22  
Ma et al, ViTAE-S, IJCV'23

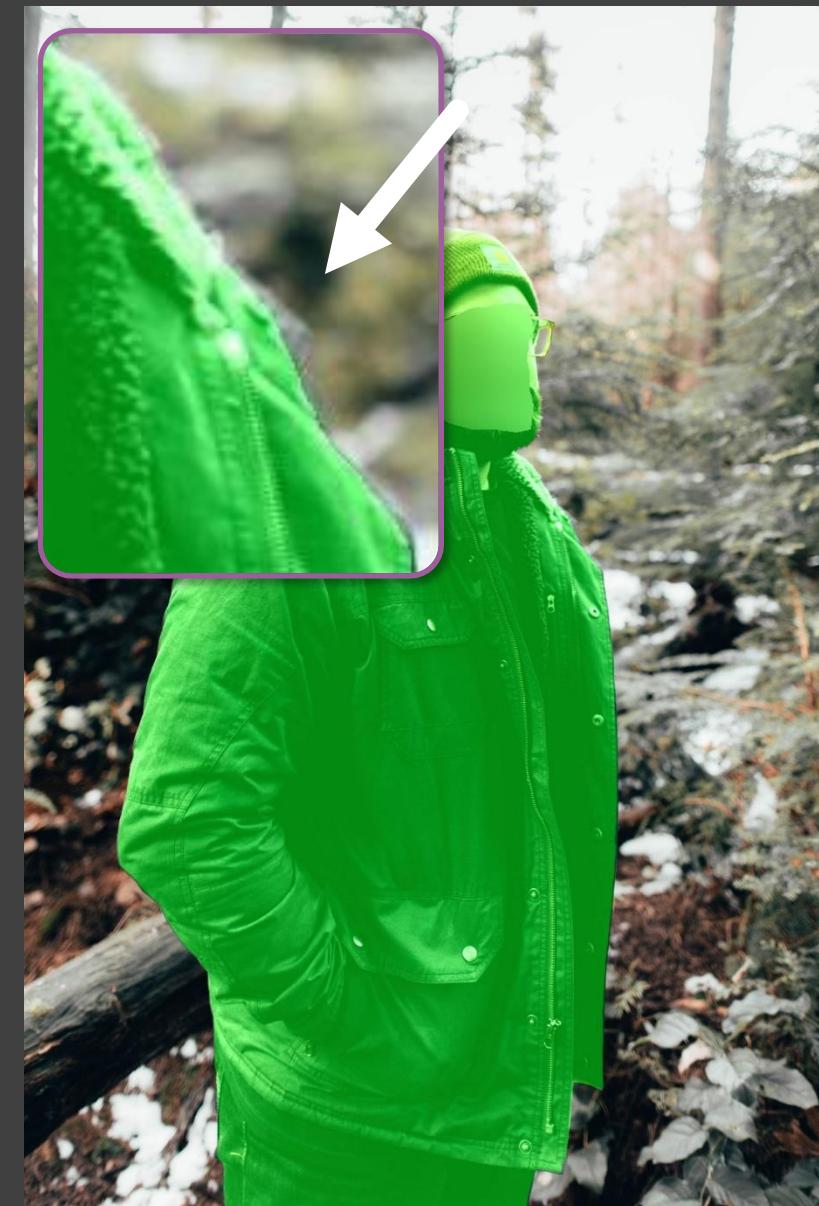
**DiffMat**



**Ours**



**Human Annotation**



**Input**



**Ours**



**Human Annotation**



Input



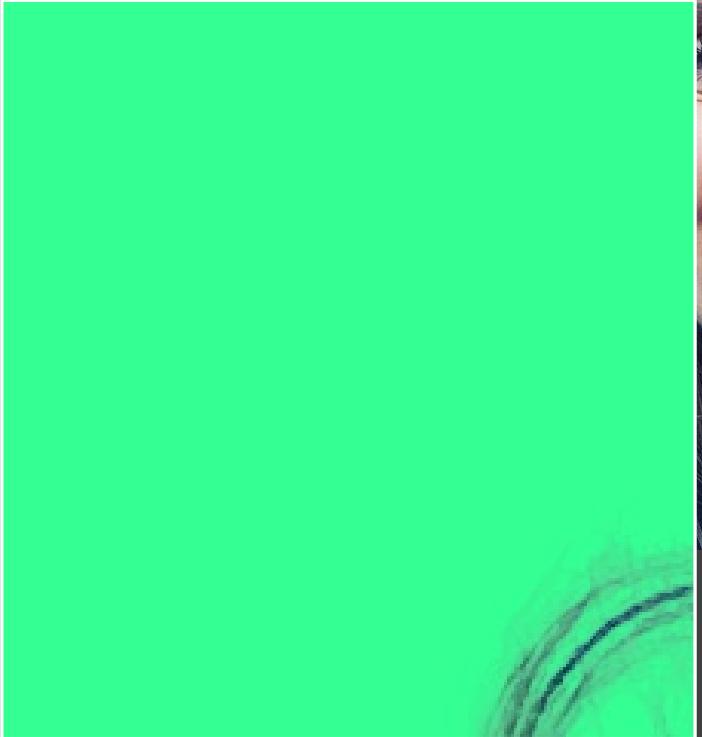
Ours



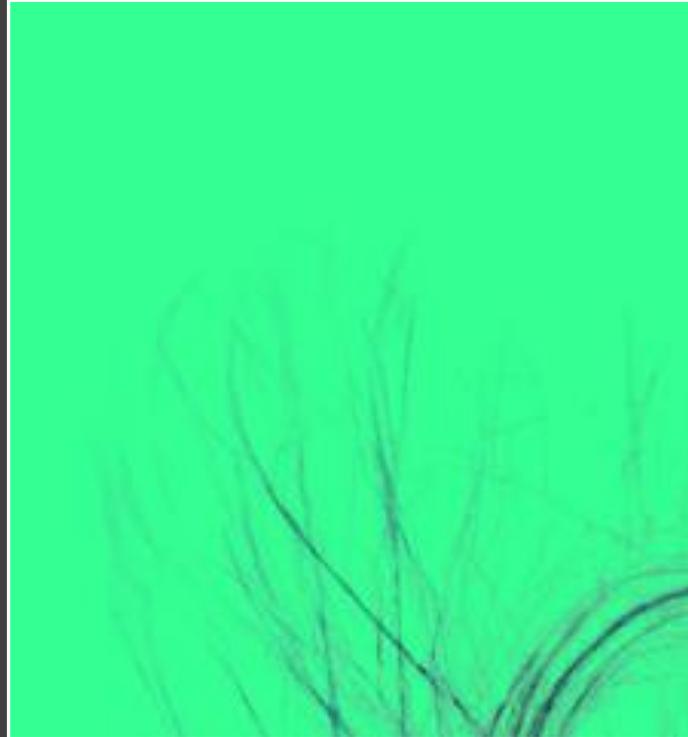
Human Annotation



**ViTAE-S**



**Ours**



**Input**



**ViTAE-S**



**Ours**



Input



ViTAE-S



Ours



ViTAE-S



Ours

# Out-of-Distribution Matting

Input



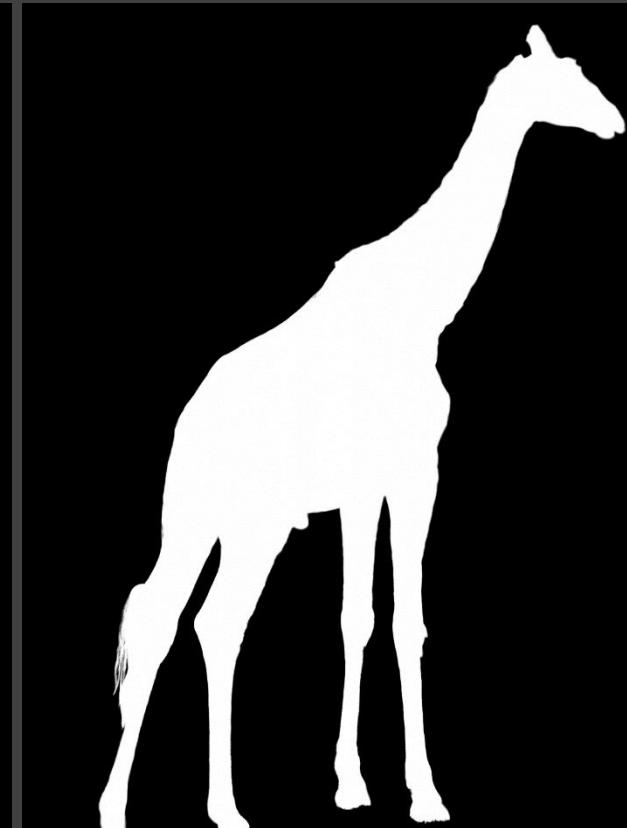
SAM-based



ViTAE-S



Ours



# Matting with Additional Guidance



Input



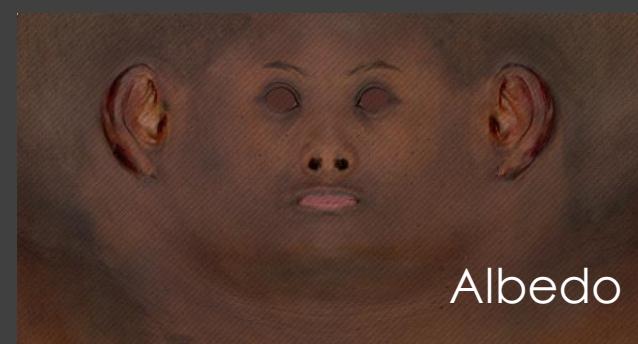
w/o guidance



w/ guidance

# Beyond Matting

- ▶ Other **image-like** intermediate parameters without accurate label / real date
  - ▶ Single Image Normal Map (Single Image)
  - ▶ Albedo (Single Image)
  - ▶ Depth Estimation (Single Image)



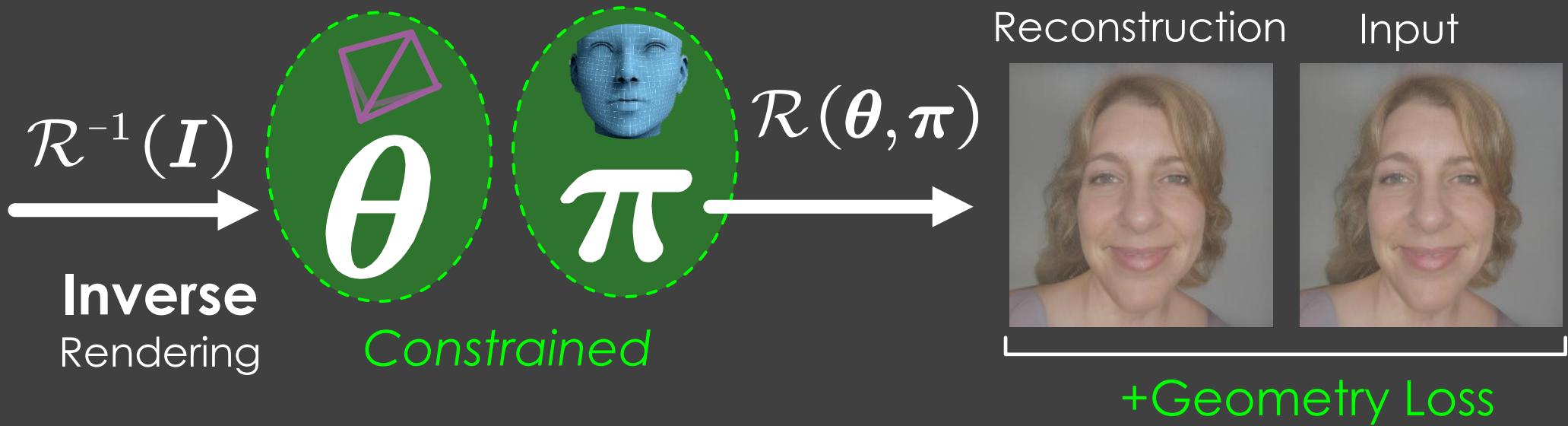
marigold, CVPR'24

# Factorization-based Methods

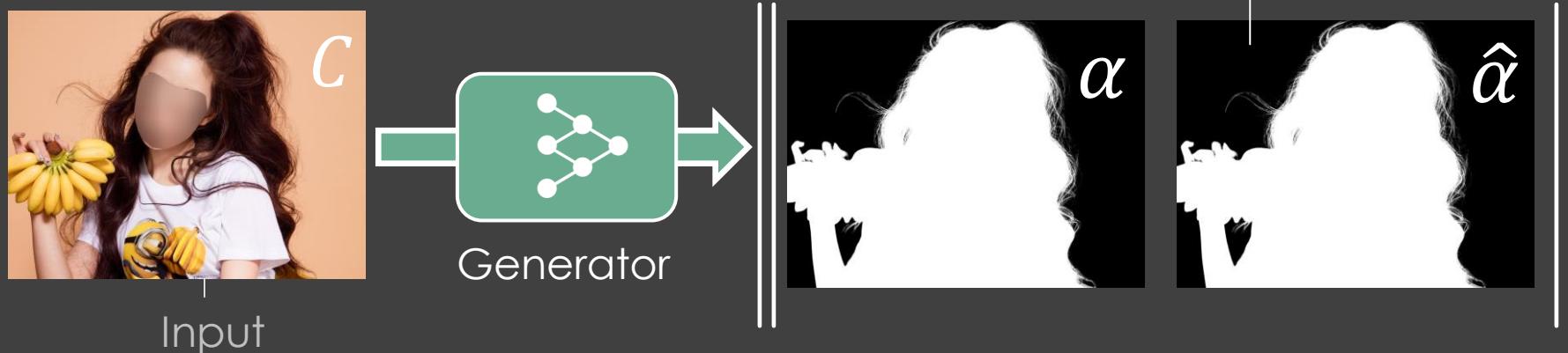
**Optimization-based:** no labels required



Input



**Learning Factorization with Labels:** imperfect labels



**Thank you!  
Questions or Comments?**