# Image Factorization and Manipulation with Generative Regularizations

Zhixiang Wang
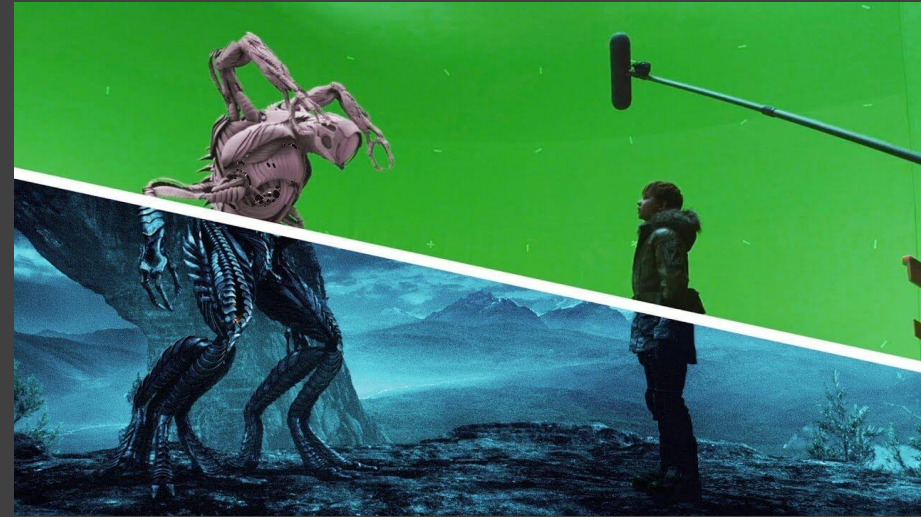
PhD candidate

The University of Tokyo

# Goal: GenAI + Advanced Cameras for VFX

Reduce actor, time, and money costs

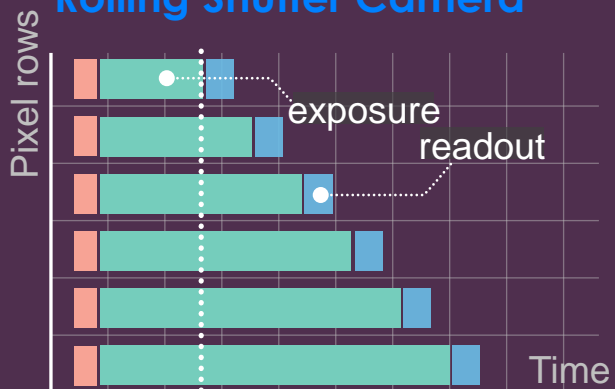# Research Works



**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**

Pixel rows
exposure
readout
Time

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

**Foggy Scene Understanding**

Foggy
Mask

Ma, **Wang** et al, CVPR 2022

**Generative Models**

**Geometric Distortion Correction**

Input
Ours

**Wang** et al, IJCV 2024

**Background Replacement**

**Wang** et al, SIGGRAPH 2024

**Style Transfer**

iPhone
DSLR

Chang, **Wang** et al, ECCV 2020

4

# Viewpoint + Lens          Background



Input

Ours

**Perspective Distortion Correction**

Wang et al, IJCV  2024

**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**
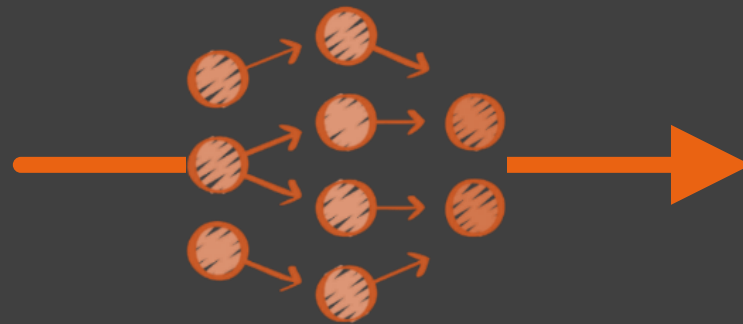
**Background**

# Difficulty in Controlling Imaging Factors



▶ Numerous factors

▶ Specialized equipment

    ▶ Inflexible

    ▶ Expensive

▶ Expertise

▶ Multiple trials

Image credit: DALLE2

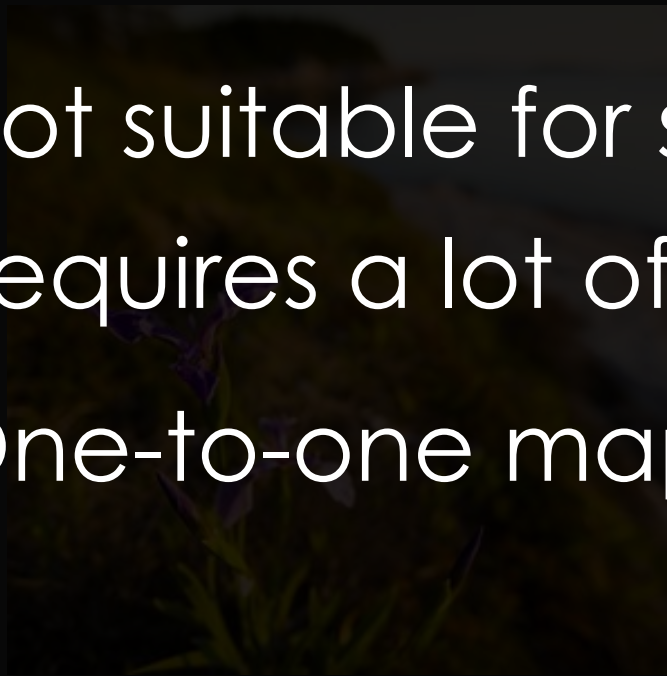# Simple yet Popular DL-based Solution



**Image-to-Image Transform**

Undesired Samplings

Desired Samplings
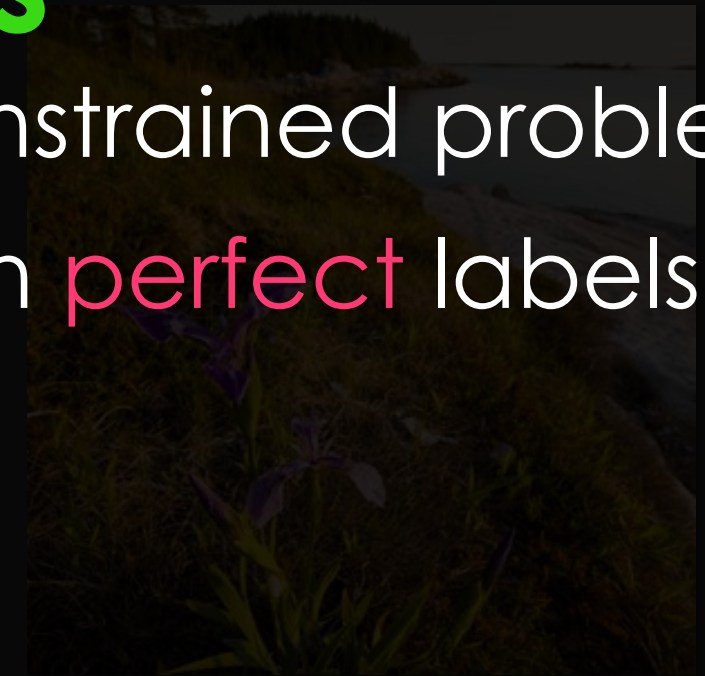
# Challenges

- Not suitable for severe under-constrained problems

- Requires a lot of paired data with perfect labels

- One-to-one mapping

Image-to-Image
Transform

Undesired
Samplings

Desired
Sampling

9

# Image Factors and Factorization



Undesired Image

**SCENE**

Material

Geometry

**IMAGING FACTORS**

Background

Camera

Viewpoint

Lighting

Motion

Sensor

Lens

Shutter

Aperture

# Image Manipulation

Undesired Image

Material

Geometry

SCENE

Background

Sensor

Camera

Lens

Viewpoint

Shutter

Aperture

Lighting

Motion

IMAGING FACTORS

Render

Re-rendering

11

# Harness Pre-trained Generative Models

**Optimization-based:** no labels required



Input

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

**Inverse** Rendering

$\boldsymbol{\theta}$

$\boldsymbol{\pi}$

*Constrained*

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

Reconstruction    Input

+Geometry Loss

**Learning with Labels:** imperfect labels

Human Annotations

$C$

Input

Generator

$\alpha$

$\hat{\alpha}$

# Viewpoint + Lens    Background



Input

Ours

**Perspective Distortion Correction**

**Wang et al, IJCV  2024**

**Matting by Generation**

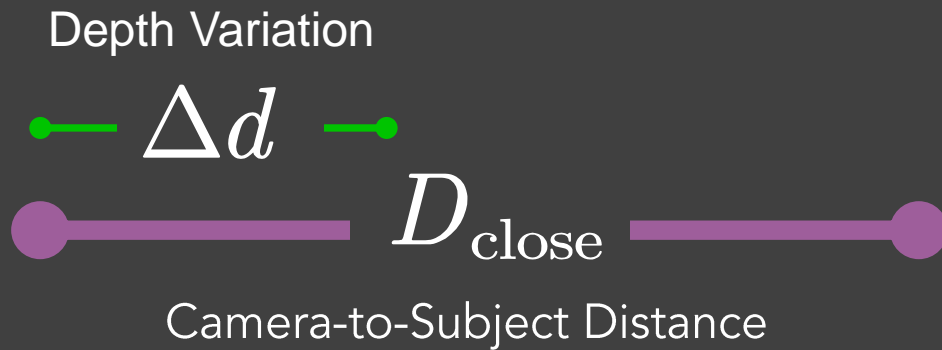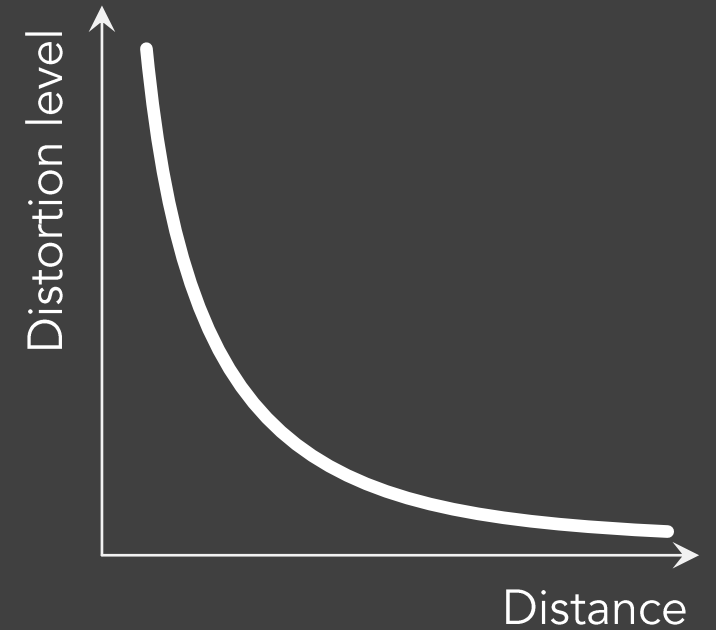**Wang et al, SIGGRAPH 2024**

Disappeared Ears

Huge Nose

Asymmetric Face

Sharp Chin

19 m

# Short Camera-to-Subject Distance



Camera-to-Subject Distance

< 60 cm

# Perspective Projection



Perspective

Depth Variation

$$\Delta d$$

$$D_{\text{close}}$$

Camera-to-Subject Distance

Distortion level

Distance

# Weak-perspective Projection



Depth Variation

$\Delta d$

$D_{\text{close}}$

Camera-to-Subject Distance

**Perspective**



**Weak-perspective**



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

$D_{\text{far}}$

# Manipulate Viewpoint and Lens

**Perspective**

**Weak-perspective**

# Existing Methods – **Warping-based**



Flow Estimation

Inpainting

Warp

Input

Correction Flow

Output

Fried et al, SIGGRAPH'16

Zhao et al, ICCV'19

# Limitations of Existing Methods



Fried+

Zhao+

Input          Output          Target

- ▶ **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

$\boldsymbol{\theta}$ $\boldsymbol{\pi}$

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

**Forward** Rendering

Reconstruction    Input

# Optimization-based Factorization

**Challenge**: ill-posed/unconstrained

Input: $I$
**Single** Image

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

$\theta$ $\pi$
*Constraints*

$\mathcal{R}(\theta, \pi)$
**Forward** Rendering

Reconstruction  Input

+Geometry Loss

23

# Ambiguity of Parameters



Many **combinations** resemble input image

Face Shape · · · Flat

Camera-to-Subject Distance — Small / Large

Focal Length — Small / Large

# 3D GAN Prior as Face Constraint



**Single** Image

Reconstruction

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

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

$$\boldsymbol{\theta} \quad \boldsymbol{\pi}$$

*Unconstrained*

**Inverse** Rendering

$$\boldsymbol{I}$$

3D GAN

Generator

Random noise

Implicit Representation

**Probabilistic Representation**

# Camera Regularization (CR)

Input: $I$
**Single** Image

Reconstruction

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

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

$\boldsymbol{\theta}$ $\boldsymbol{\pi}$

**Inverse** Rendering

*Unconstrained*

3 Strategies

# CR 1: Focal Length Re-parameterization

# CR 2: Optimization Scheduling

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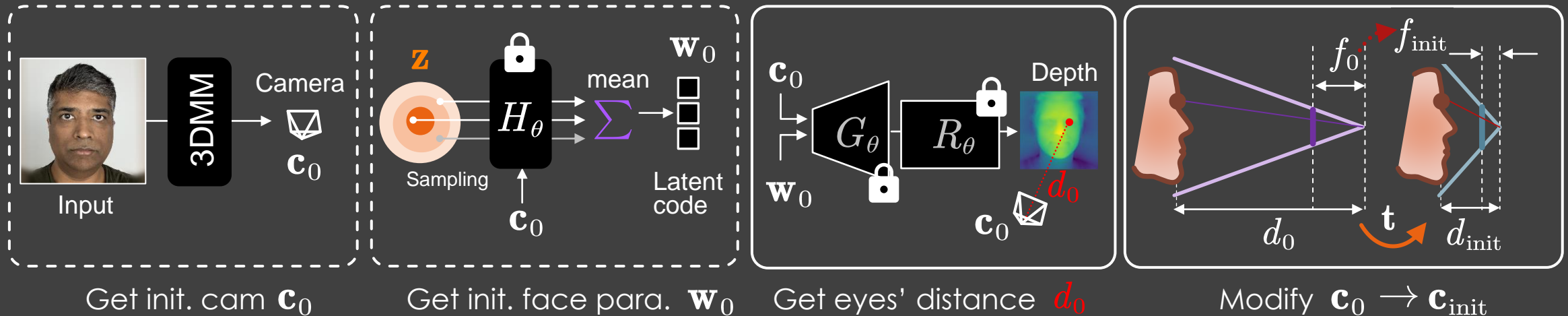


Optimization  Frozen



$\mathbf{c}_{\text{init}}$  $\mathbf{c}_{\text{opt}}$  $G_\theta$  $R_\theta$  Rendering  Input

$\mathcal{L}_{\text{landmark}} + \mathcal{L}_{\text{LPIPS}}$

$\mathbf{w}_0$

Optimize camera

$\mathbf{c}_{\text{opt}}$  $\hat{\mathbf{c}}$  $G_\theta$  $R_\theta$  Rendering  Input

$\mathcal{L}_{\text{LPIPS}}$

$\mathbf{w}_0$  $\hat{\mathbf{w}}$

Optimize face, refine camera

29

# CR 3: Better Initialization

Start from a close-up camera position



Original initialization

Get init. cam $\mathbf{c}_0$

Camera $\mathbf{c}_0$

Input

3DMM

$\mathbf{z}$

Sampling

$H_\theta$

$\mathbf{c}_0$

mean $\sum$

$\mathbf{w}_0$

Latent code

Get init. face para. $\mathbf{w}_0$

$\mathbf{c}_0$

$\mathbf{w}_0$

$G_\theta$

$R_\theta$

Depth

$d_0$

$\mathbf{c}_0$

Get eyes' distance $d_0$

Re-parameterization

$f_0$

$f_{\text{init}}$

$d_0$

$\mathbf{t}$

$d_{\text{init}}$

Modify $\mathbf{c}_0 \rightarrow \mathbf{c}_{\text{init}}$

# Ambiguity Caused by Loss

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



(a) $L2$ loss

(b) Landmark loss

# Geometric Regularization

Uncertainty-based Loss

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



prediction

input

$m_i$

$m_i'$

# Extensions for Full-frame Image

# Results – Mesh



Distorted input     HFGI3D [58]     Triplanenet [6]     Ours

Other GAN inversion methods

# Results

Input

Output

# Results – Continuous Manipulation

# Results – Comparison

Stretch-like 👎

3D **geometric consistent** 👍



Fried et al, SIGGRAPH'16



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**

41

42

# Viewpoint + Lens          # Background



Input

Ours

**Perspective Distortion Correction**

**Matting by Generation**

**Wang et al, SIGGRAPH 2024**

43

# Harness Pre-trained Generative Models



Optimization-based: no labels required

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

Inverse Rendering

$\theta$

$\pi$

Constrained

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

Reconstruction    Input

Input

+Geometry Loss

Learning with Labels: imperfect labels

Human Annotations

$C$

Input

Generator

$\alpha$

$\hat{\alpha}$

44

# Manipulate Background



Background Gallery

45

# Factorization Problem

Input: $I$
**Single** Image

$\alpha F + (1 - \alpha)B$
Reconstruction

$\alpha$

Matting

$F$

$B$

Composition

**Inverse**
Rendering

**Forward**
Rendering

$\alpha$

$F$

$B^*$

Re-rendering

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

# Learning with Labels



Human Annotations

Input

Regression

$C$

$\alpha$

$\hat{\alpha}$

Input image

Label

**Poor** label quality

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

47

# Limitations of Existing Methods



Input

ModNet

**Unnatural** boundary

Fasphoto

Ke et al, MODNet, AAAI'22

# 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



guidance

Text prompt

Latent Diffusion **U-Net**

$T$

52

# Pipeline for Processing HR Images

Get potential areas by uncertainty



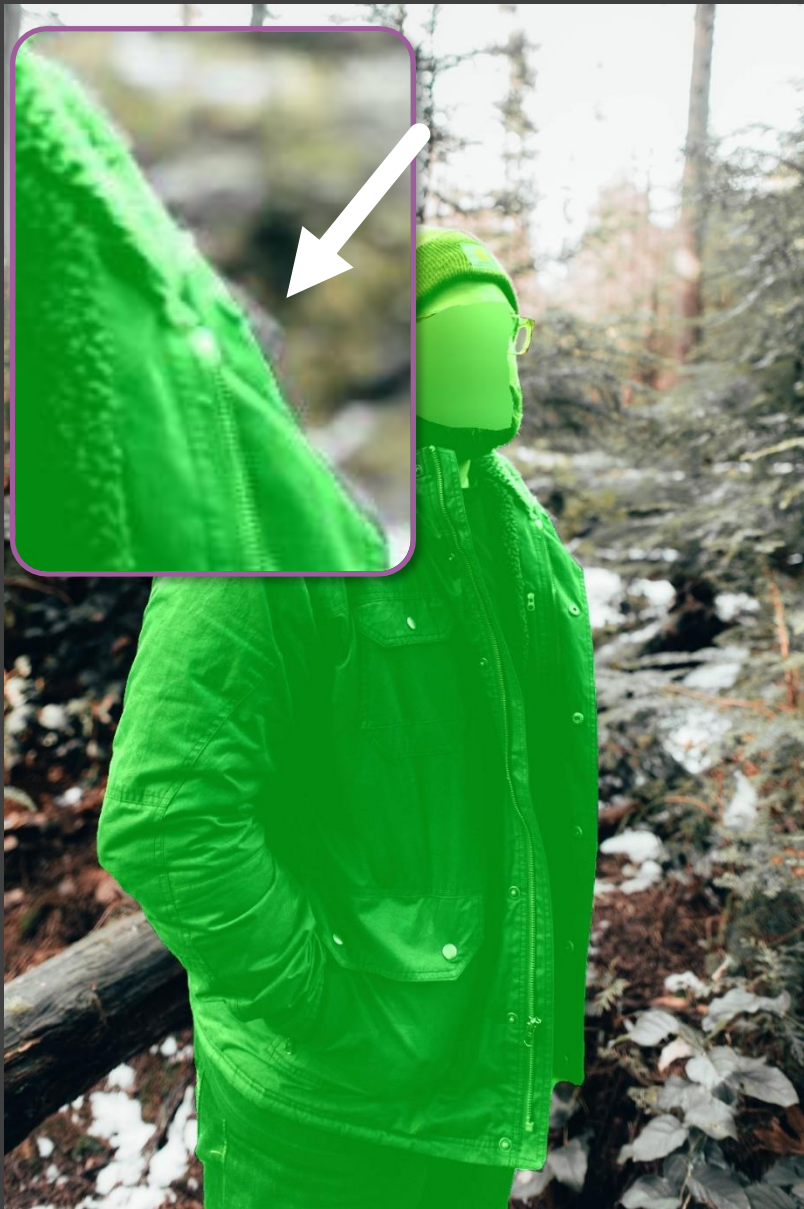Uncertainty

Candidate

# Full Pipeline

MODNet

ViTAE-S

Ours

Reference

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

55

**DiffMat**

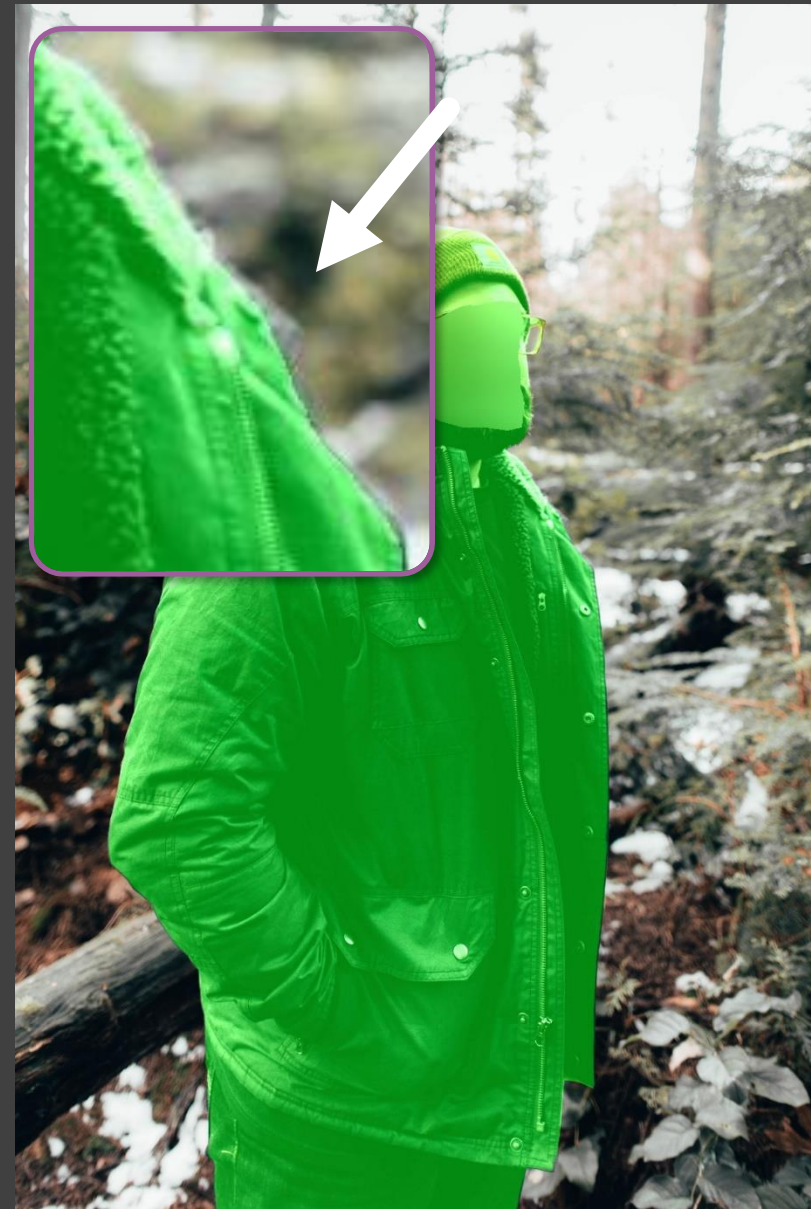**Ours**

**Human Annotation**

| Input | Ours | Human Annotation |

Input
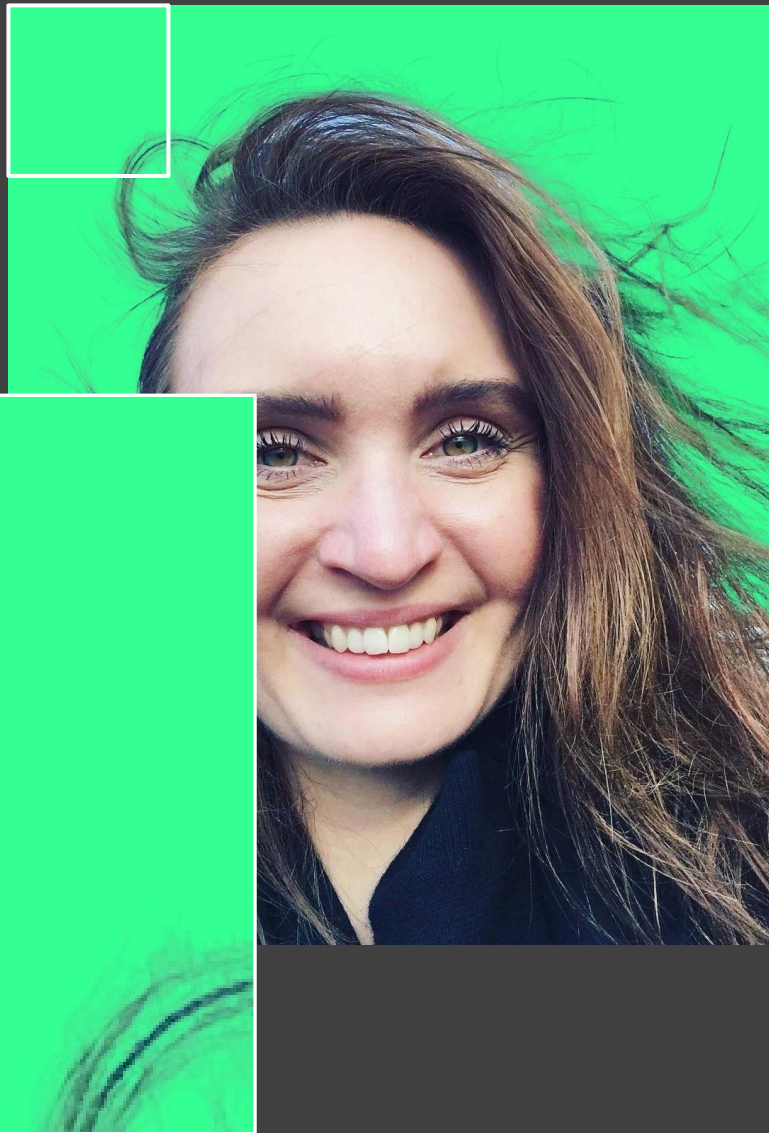
Ours

Human Annotation

ViTAE-S      Ours
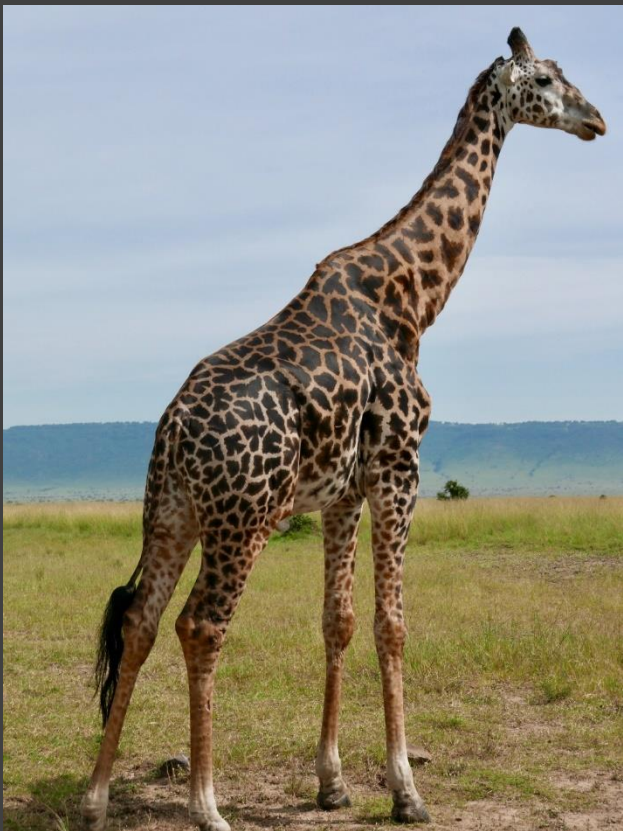
Input             ViTAE-S             Ours

ViTAE-S
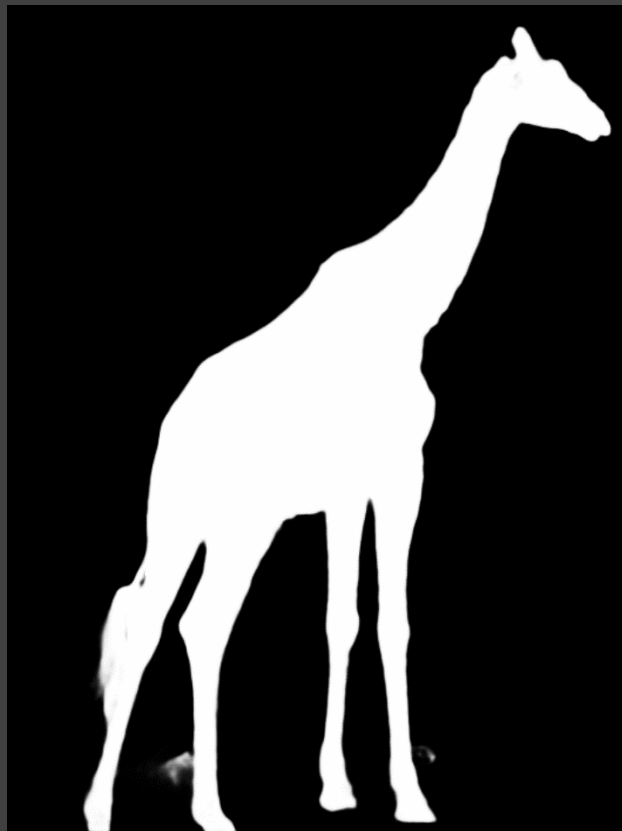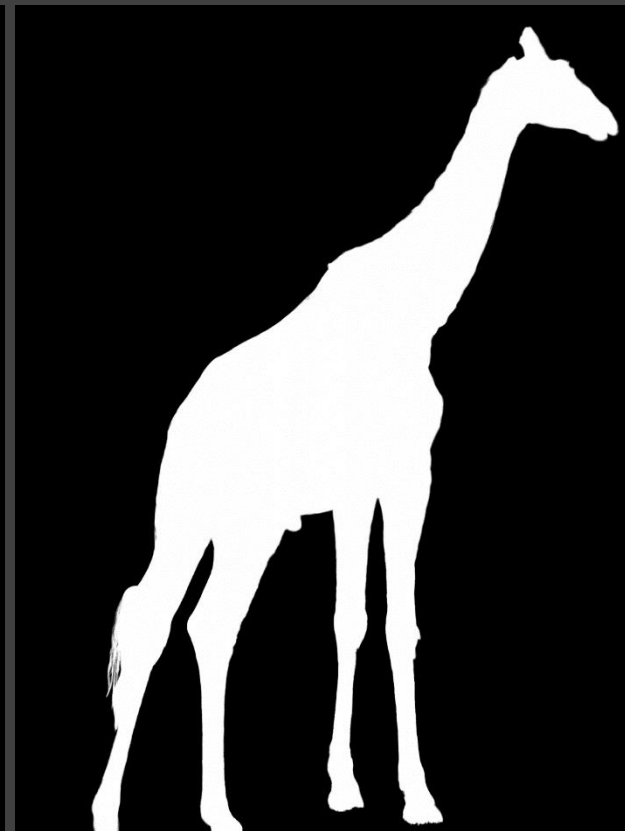
Ours

# Out-of-Distribution Matting
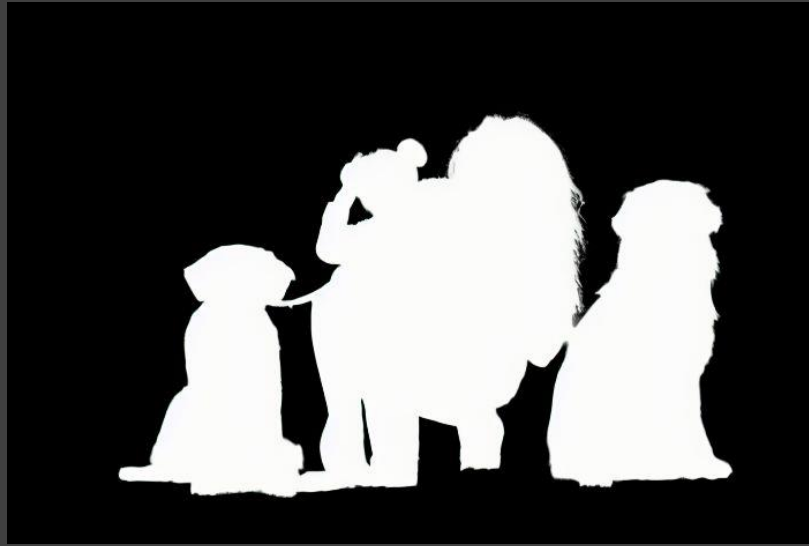


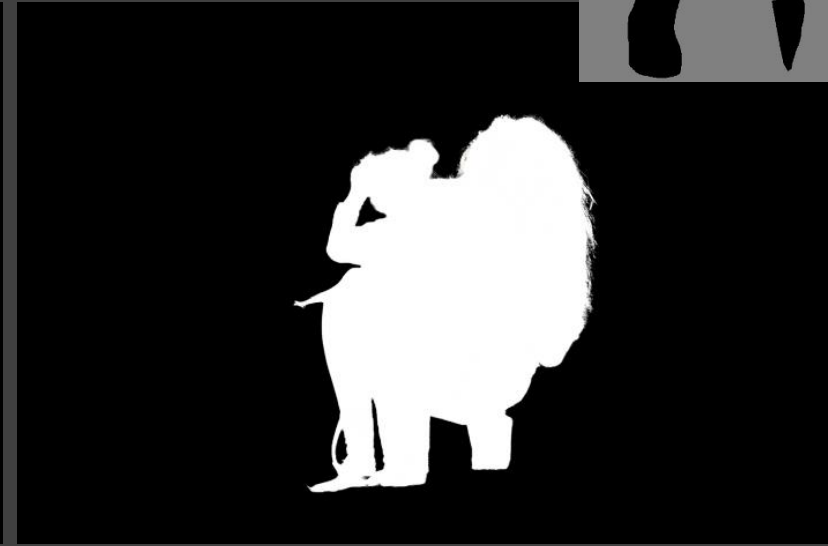Input     SAM-based     ViTAE-S     Ours
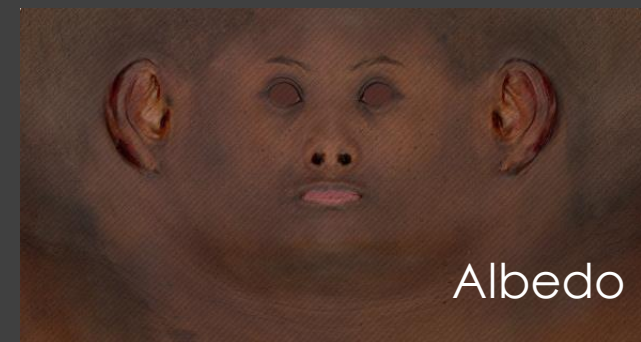
# Matting with Additional Guidance



Input

w/o guidance

w/ guidance

# Beyond Matting
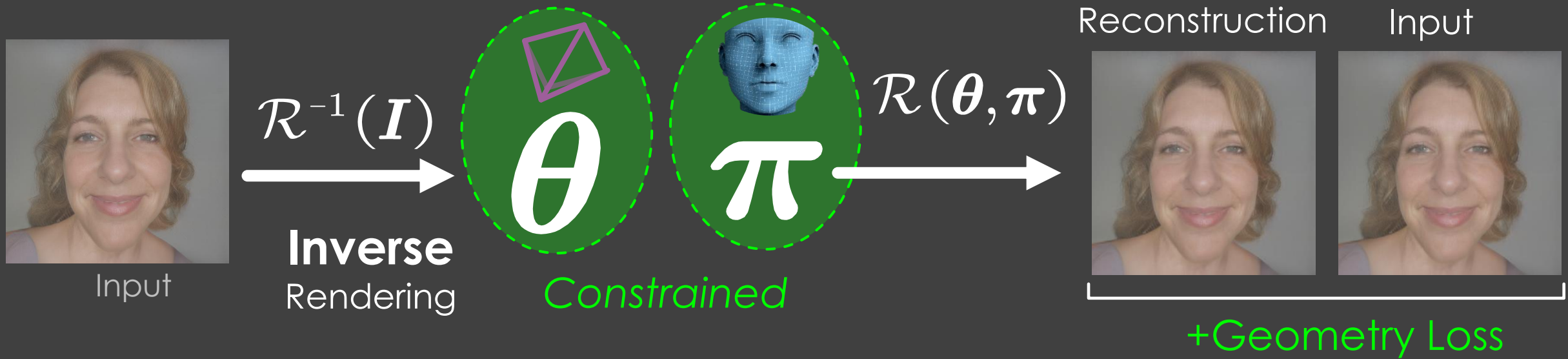
▶ Other *image-like* intermediate parameters without accurate label / real date

    ▶ Single Image Normal Map (Single Image)

    ▶ Albedo (Single Image)
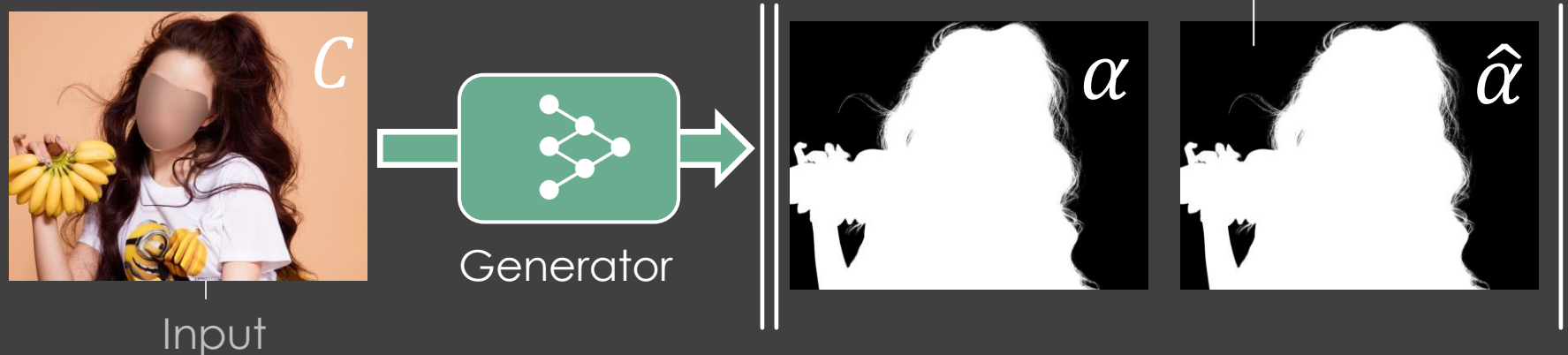
    ▶ Depth Estimation (Single Image)



Normal

Env Map

Depth

Albedo

marigold, CVPR'24

# Factorization-based Methods



**Optimization-based:** no labels required

Input

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

**Inverse** Rendering

$\boldsymbol{\theta}$  $\boldsymbol{\pi}$

*Constrained*

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

Reconstruction    Input

+Geometry Loss

**Learning Factorization with Labels:** imperfect labels

Human Annotations

$C$

Input

Generator

$\alpha$    $\hat{\alpha}$

# Thank you!
## Questions or Comments?