# Controllable 3D Generation

Vova Kim | Senior Research Scientist, Adobe Research

# Motivation
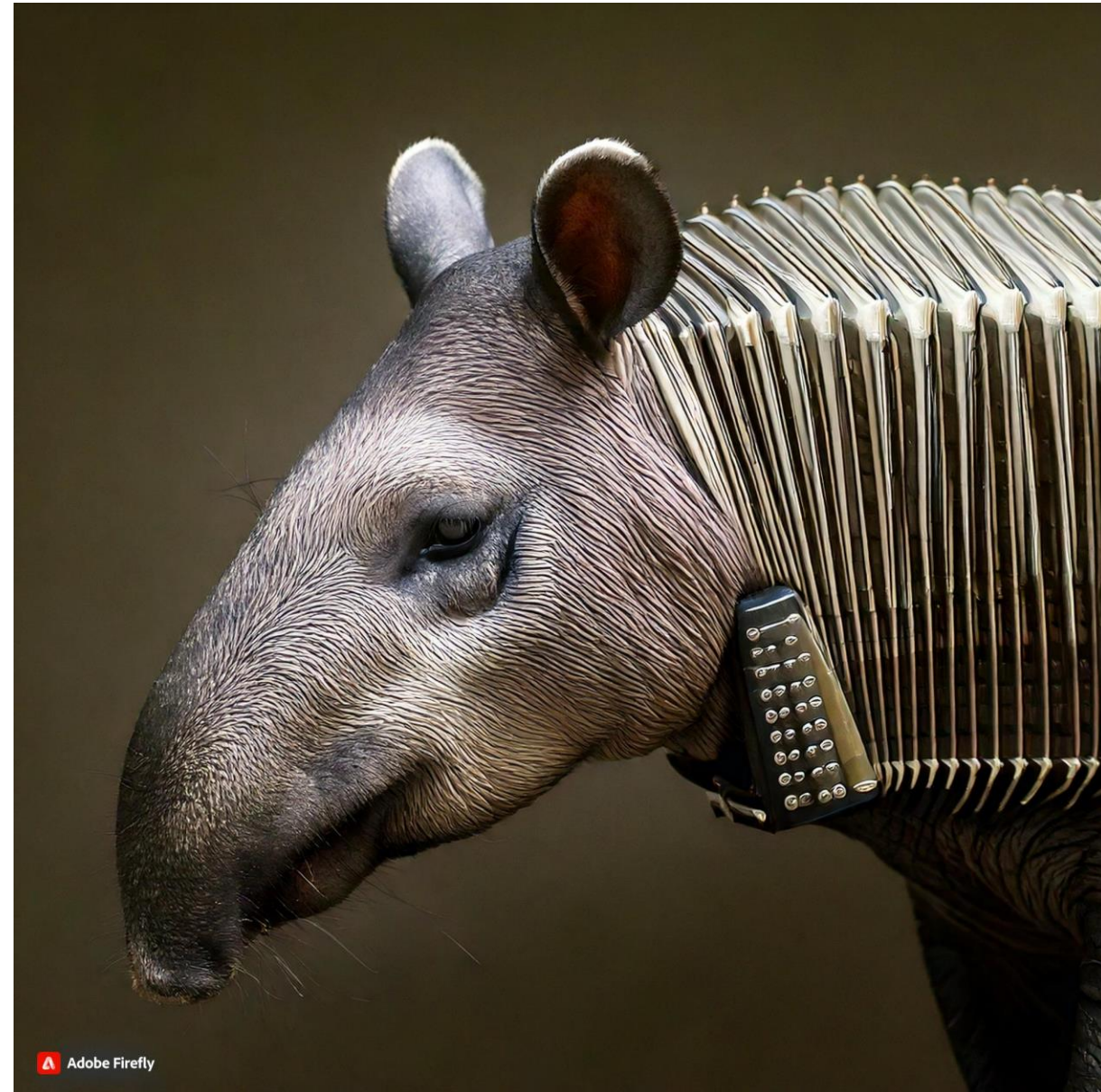
- Text-to-Image ~ amazing progress



2021



2024

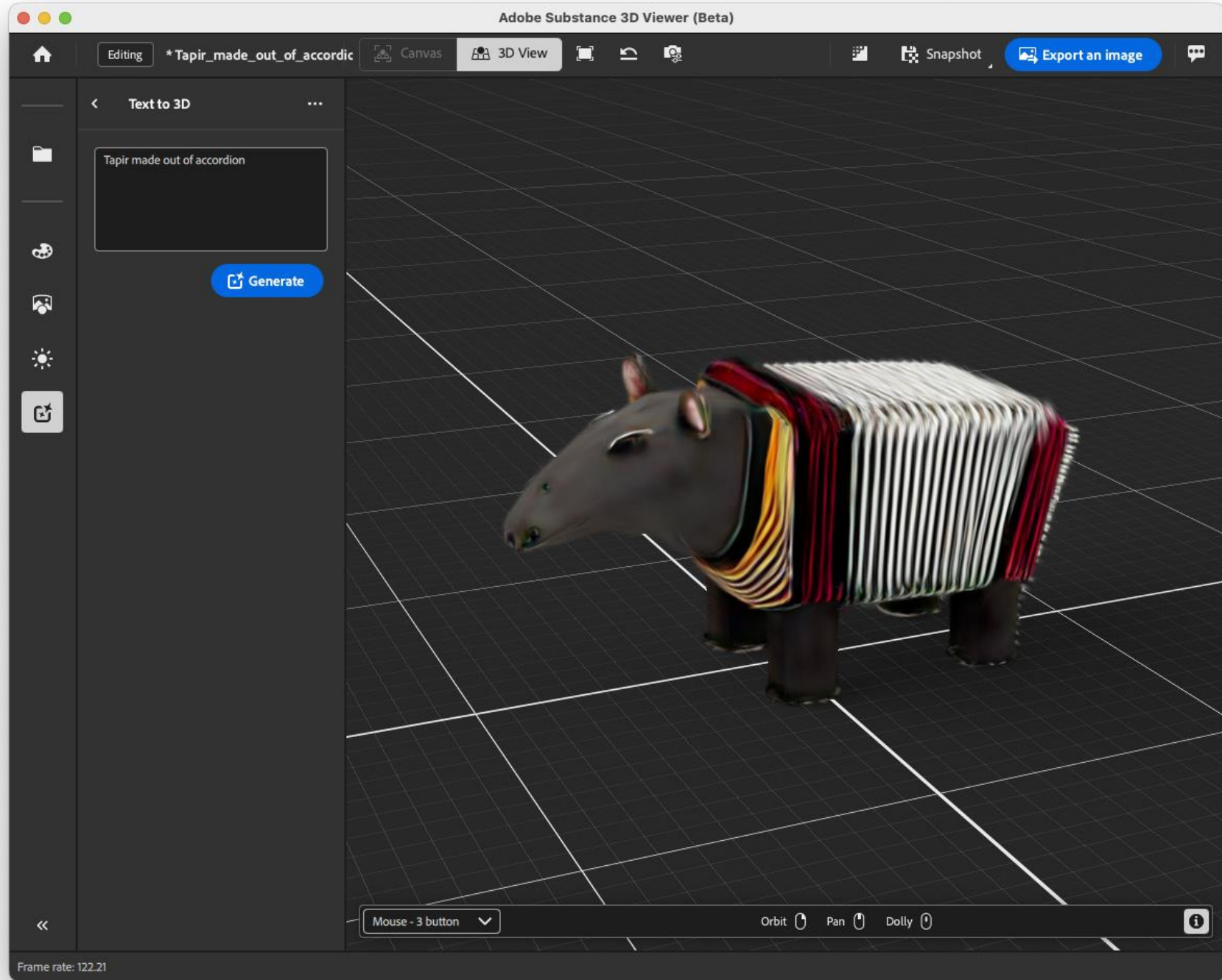"Tapir made out of accordion"

# Motivation

- Text-to-Image

- Text-to-3D



"Tapir made out of accordion"

# Motivation

- Text-to-Image

- Text-to-3D



"Tapir made out of accordion"

# Motivation

- Text-to-Image

- Text-to-3D

- 3D-to-Image



"Tapir made out of accordion"

# Motivation

- Text-to-Image

- Text-to-3D

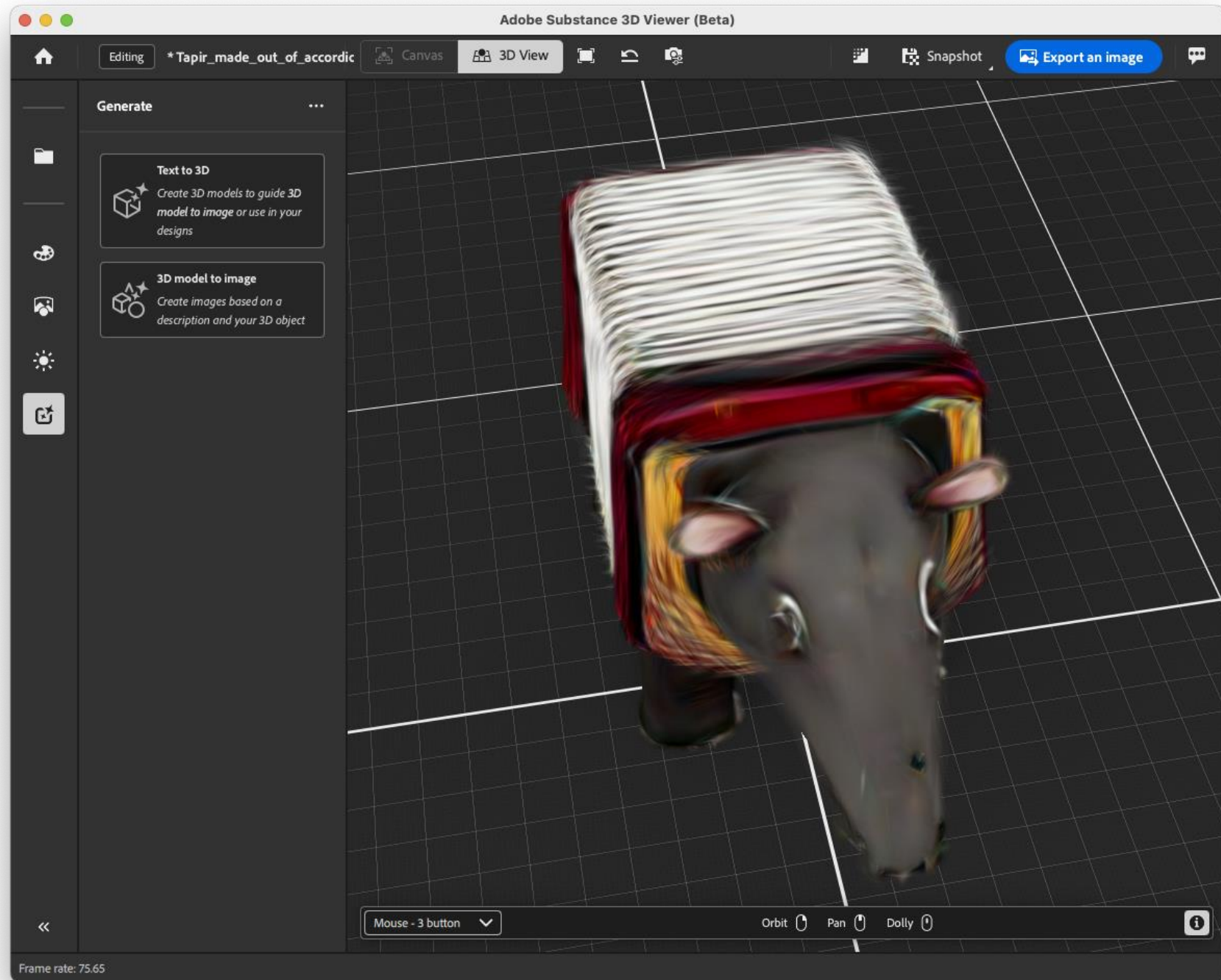- 3D-to-Image

**How do we customize it further? E.g.**

Tilt head (deform)

Make it look more like hippo (deform)

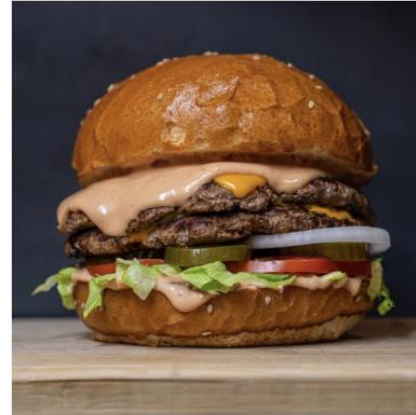Add unicorn horn (sculpt)

Make skin look more natural (detail)

Move eye to another location (detail)



"Tapir made out of accordion"

# Prior Work: Image Editing

- Entangled control:

  - Geometry

  - Materials

  - Lighting

  - Camera Intrinsics

  - Camera Extrinsics

  - Composition

- Text alone is not always well-defined
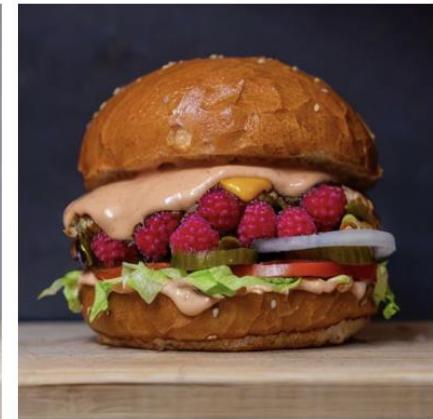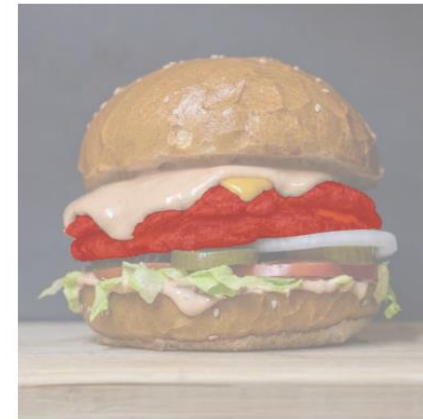


Input

"Replace the beef with raspberries"
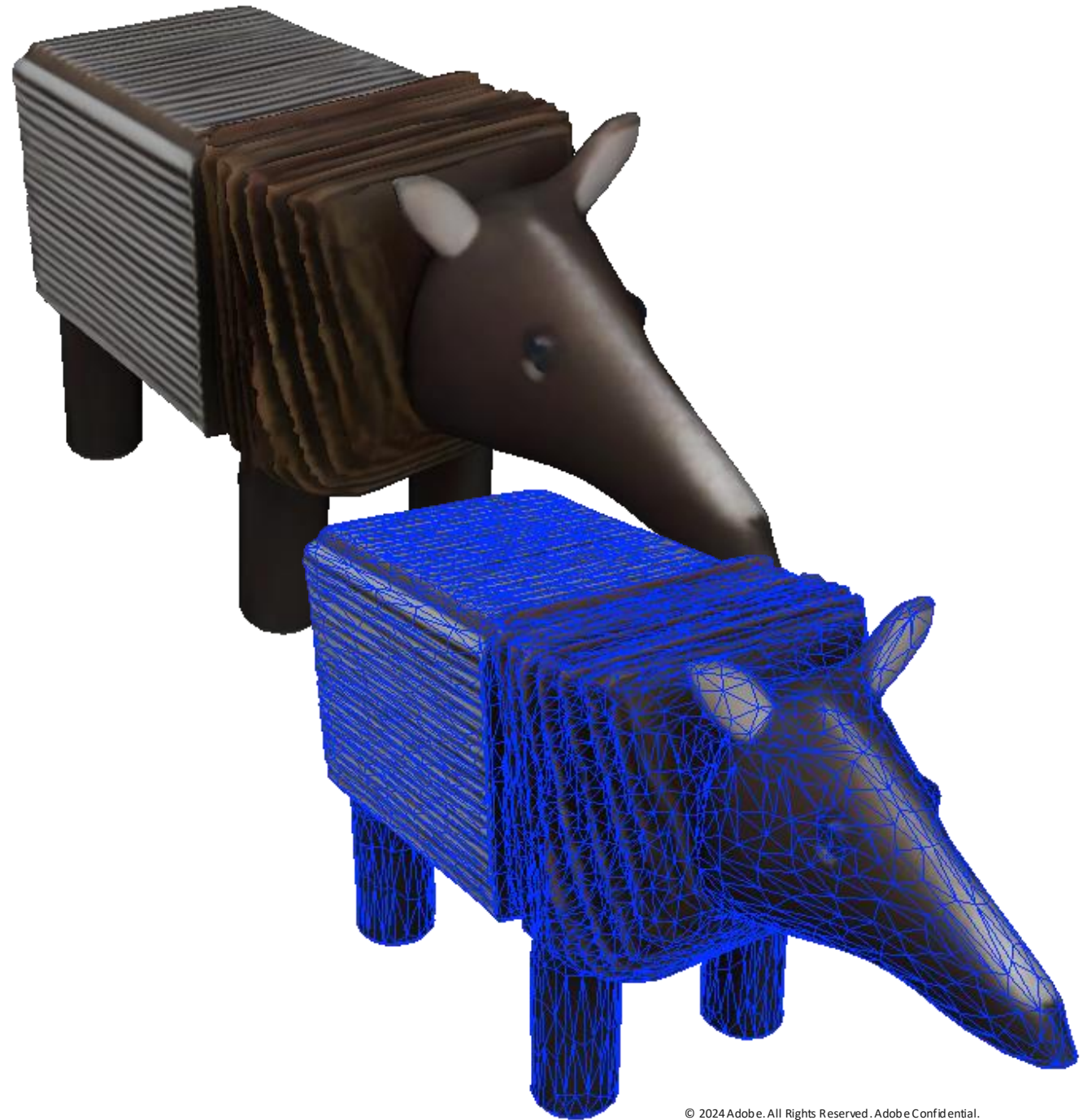
Mask image

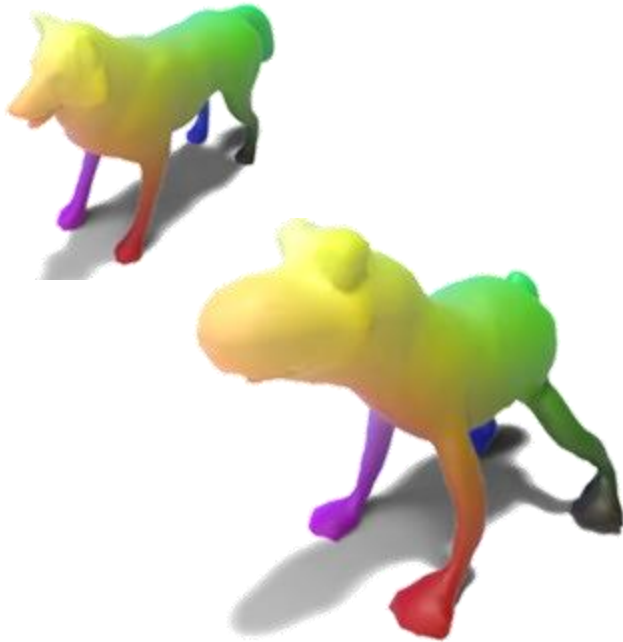Edited image

DiffEdit

InstructEdit

InstructEdit, Wang et al, 2023

# Prior Work: Representations

- NeRFs / Gaussian Splats

  - Geometry **(no surface priors)**

  - Materials **(not easy to disentangle from appearance)**

  - Lighting **(not easy to disentangle from appearance)**

  - Camera Intrinsics

  - Camera Extrinsics

  - Composition

- Meshes / Surfaces

  - Traditional CG models fully disentangle appearance

  - Easy to use with traditional tools that artists know well

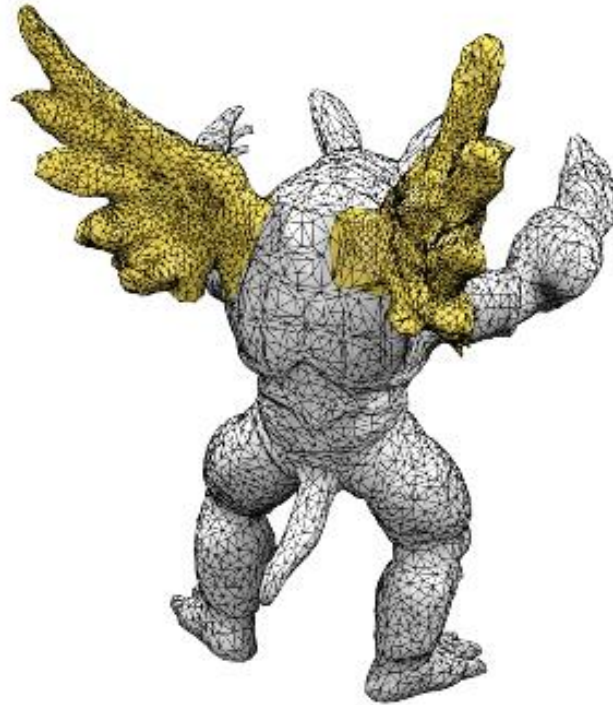  - **(hard to learn, hard to represent, poor gradients)**

# Overview

- Support mesh outputs (but use other representations as needed)

- Inspired by traditional workflows



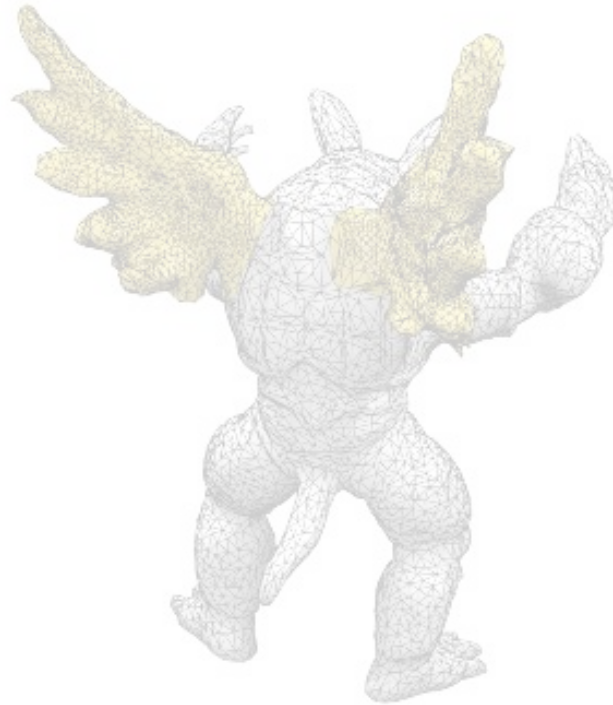**Neural Deformation**



**Generative Scultping**



**Generative Detailization**

# Overview

- Support mesh outputs (but use other representations as needed)
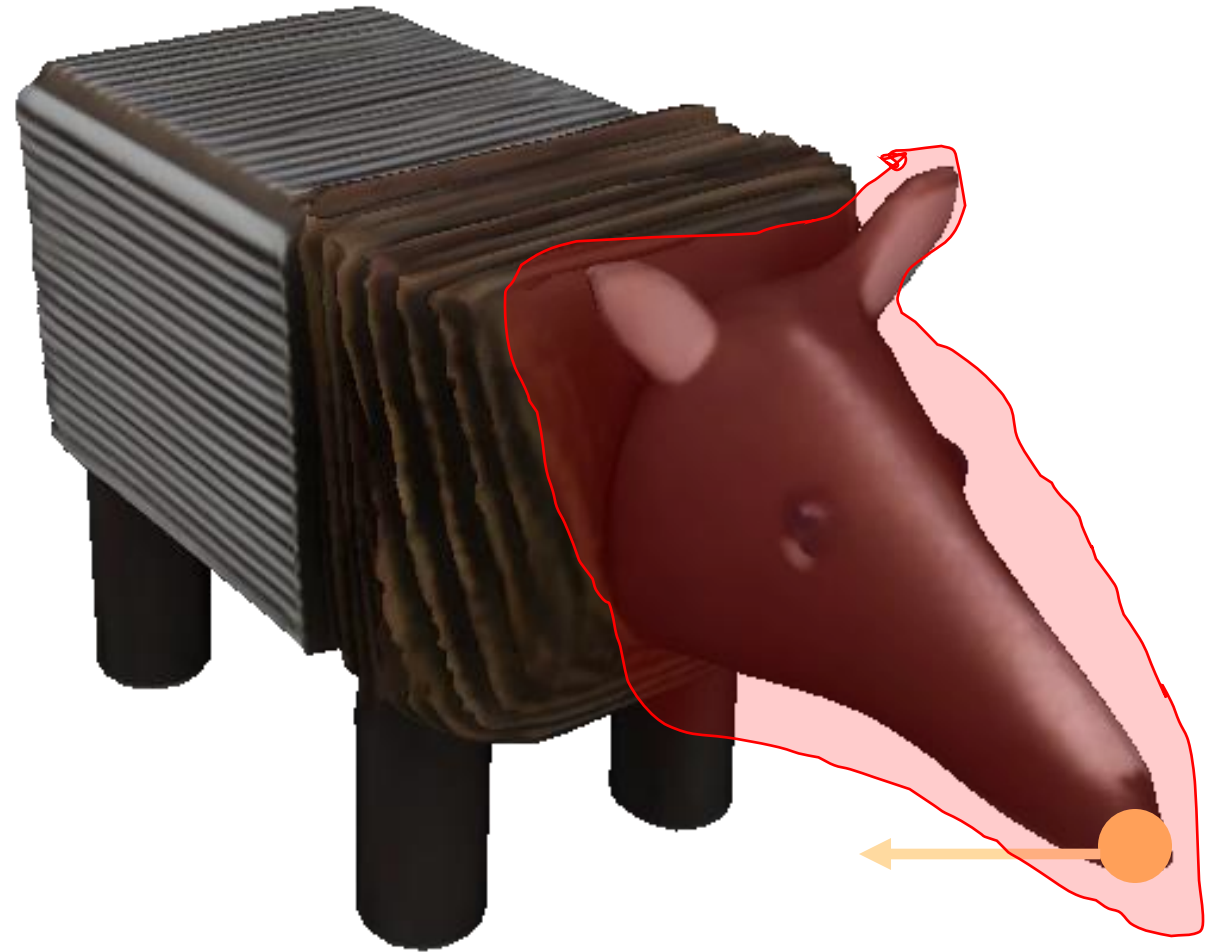
- Inspired by traditional workflows



**Neural Deformation**



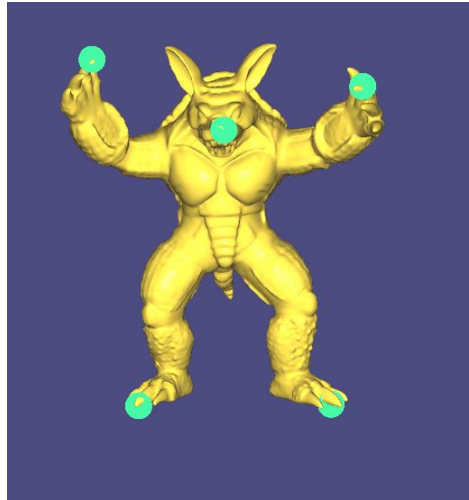Generative **Scultping**



Generative **Detailization**

# Deformation Examples

- "tilt tapir's head"

- "make it's head look more like a hippo"

# Prior Works: Deformation as Geometric Optimization

- Move control points AND preserve original geometry

- No semantics, e.g.:

  - Rubber-like behavior

  - Symmetry not preserved



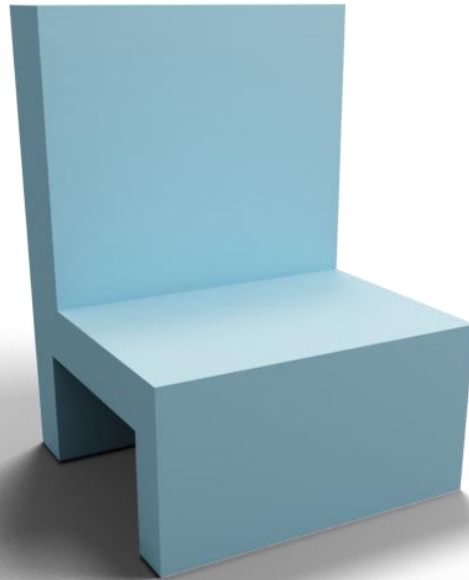**Control points**



**E.g.: As-rigid-as-possible (ARAP)**

**Adobe**

# Prior Works: Deformation as a Learnable Map

- Naïve: learn direct map ~ hard to make it smooth enough

$$f_\theta : \mathbb{R}^3 \to \mathbb{R}^3$$



Source

Target

Deep Deformation, Groueix et al. CGF 2019

# Prior Works: Deformation as a Learnable Map

- Cage-based: learn cage parameters ~ hard to predict expressive & accurate cages

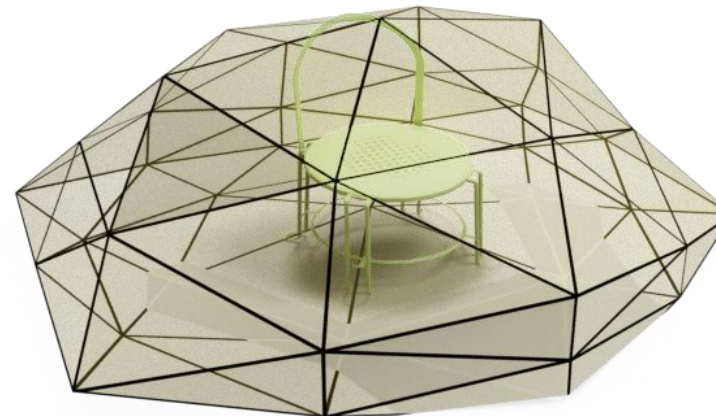$$f_\theta : \boxed{\mathcal{C}_{\text{init}} \times \mathcal{C}_{\text{deformed}}} \rightarrow \boxed{\text{MVC}} \rightarrow \mathbb{R}^3 \rightarrow \mathbb{R}^3$$

Predict cage parameters with a neural network

Use Cage-Based Deformation to define the map



Init Cage

Deformed Cage

Neural Cages, Yifan et al., CVPR 2020

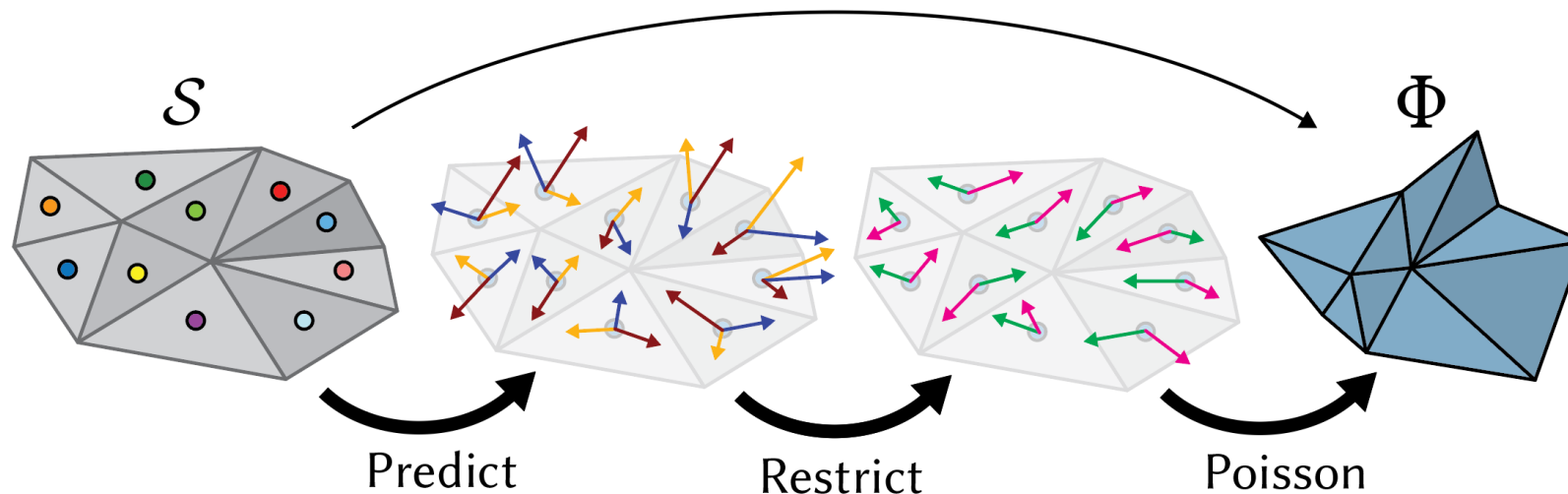# Prior Works: Deformation as a Learnable Map

- Neural Jacobian Fields

    - Smooth field (e.g., rotation is the same matrix)

    - Easy to maintain geometry details: shape-aware

$$f_\theta : \mathbb{R}^3 \rightarrow \mathbb{R}^3$$

$$f_\theta : \mathcal{S} \rightarrow \mathbb{R}^{3 \times 3}$$

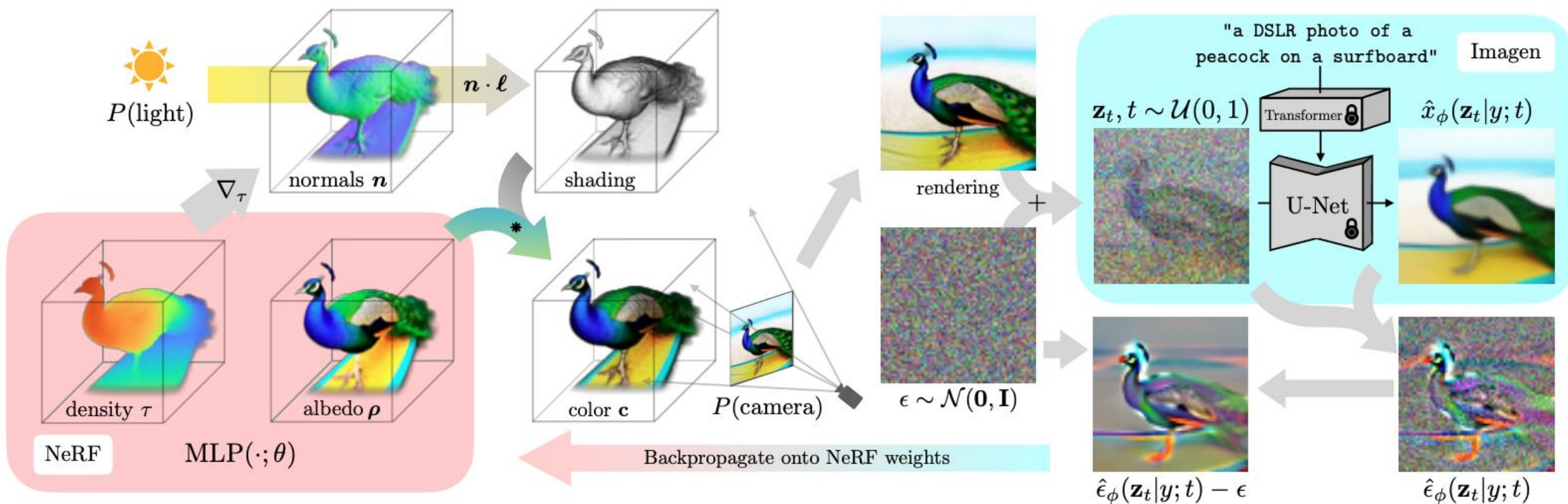Surface features

Deformation matrix



$\mathcal{S}$     $\Phi$

Predict     Restrict     Poisson

# Prior Works: Score Distillation Sampling (SDS)

- Allows interjecting priors from pre-trained 2D model



Dreamfusion, Poole et al., 2022

# Prior Works: Score Distillation Sampling (SDS)

- Allows interjecting priors from pre-trained 2D model

**Differentiable Rendering**

# Prior Works: Score Distillation Sampling (SDS)

- Allows interjecting priors from pre-trained 2D model



Dreamfusion, Poole et al., 2022

# Prior Works: Score Distillation Sampling (SDS)

- Allows interjecting priors from pre-trained 2D model
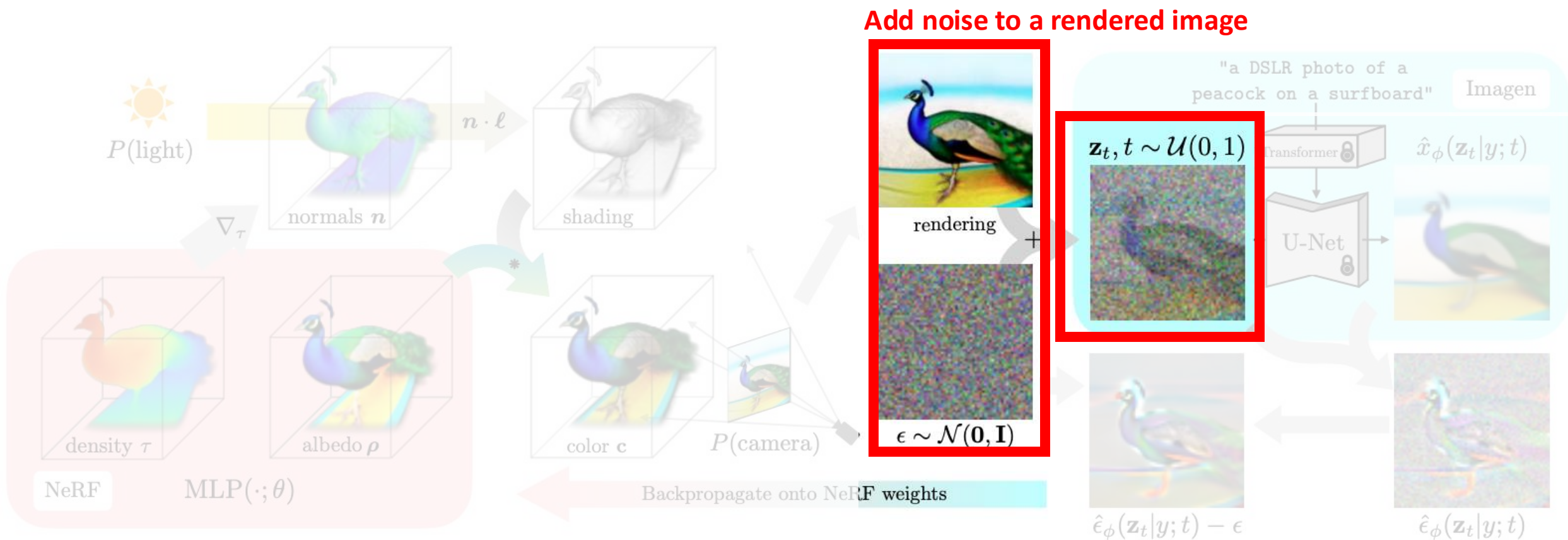
**Denoise via a pre-trained**

**Image Diffusion Model**



Dreamfusion, Poole et al., 2022

# Prior Works: Score Distillation Sampling (SDS)

- Allows interjecting priors from pre-trained 2D model



**Propagate Gradients to the 3D representation**

Dreamfusion, Poole et al., 2022

# As-<u>Plausible</u>-As-Possible Deformation



**Input**

# As-Plausible-As-Possible Deformation



**Mesh + Neural Jacobian Fields to represent deformable shape**

*Jacobian Fields*

**Deformable Input**

"A photo of [category] on a white background"

*Stable Diffusion + LORA*

$\mathcal{L}_{SDS}$

$\mathcal{L}_{handles}$

**Differentiable Rendering**

# As-Plausible-As-Possible Deformation

▪ STAGE 1: Deform via Geometric Optimization only



*Jacobian Fields*

**Mesh + Neural Jacobian Fields to represent deformable shape**

" A photo of [category] on a white background "

*Stable Diffusion + LORA*

$\mathcal{L}_{SDS}$

$\mathcal{L}_{handles}$

**Deformable Input**

# As-<u>Plausible</u>-As-Possible Deformation

- STAGE 1: Deform via Geometric Optimization only



*Jacobian Fields*

**Deformable Input**

**STAGE 1**
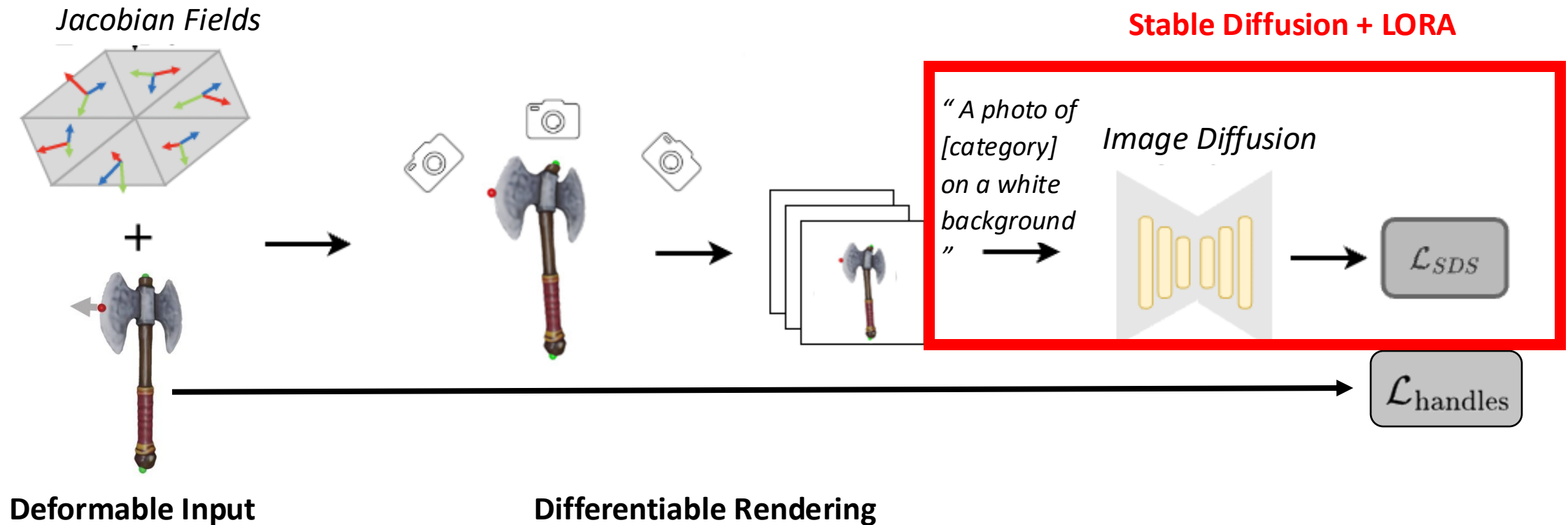
$\mathcal{L}_{\text{handles}}$

**Input**

**STAGE 1**

# As-<u>Plausible</u>-As-Possible Deformation

- STAGE 1: Deform via Geometric Optimization only

- STAGE 2: Project to "plausible" using SDS



*Jacobian Fields*

**NVDiffRast**

" A photo of [category] on a white background "

*Stable Diffusion + LORA*

$\mathcal{L}_{SDS}$

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

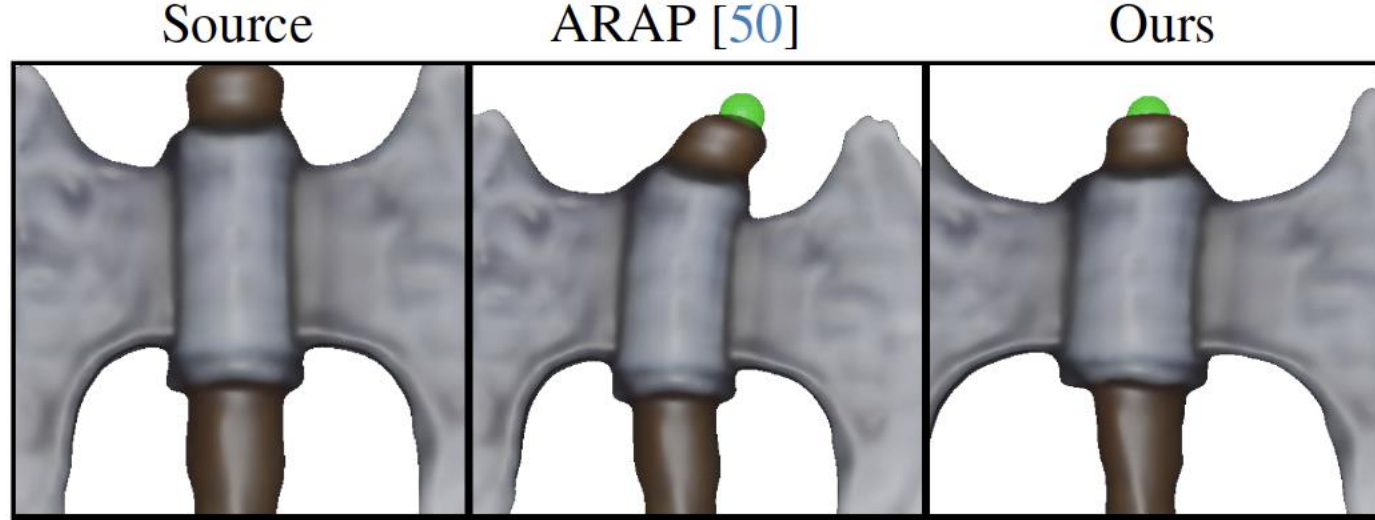**Deformable Input**                **Differentiable Rendering**

# As-Plausible-As-Possible Deformation

- STAGE 1: Deform via Geometric Optimization only

- STAGE 2: Project to "plausible" using SDS

# As-Plausible-As-Possible Deform



- STAGE 1: Deform via Geometric Optimization
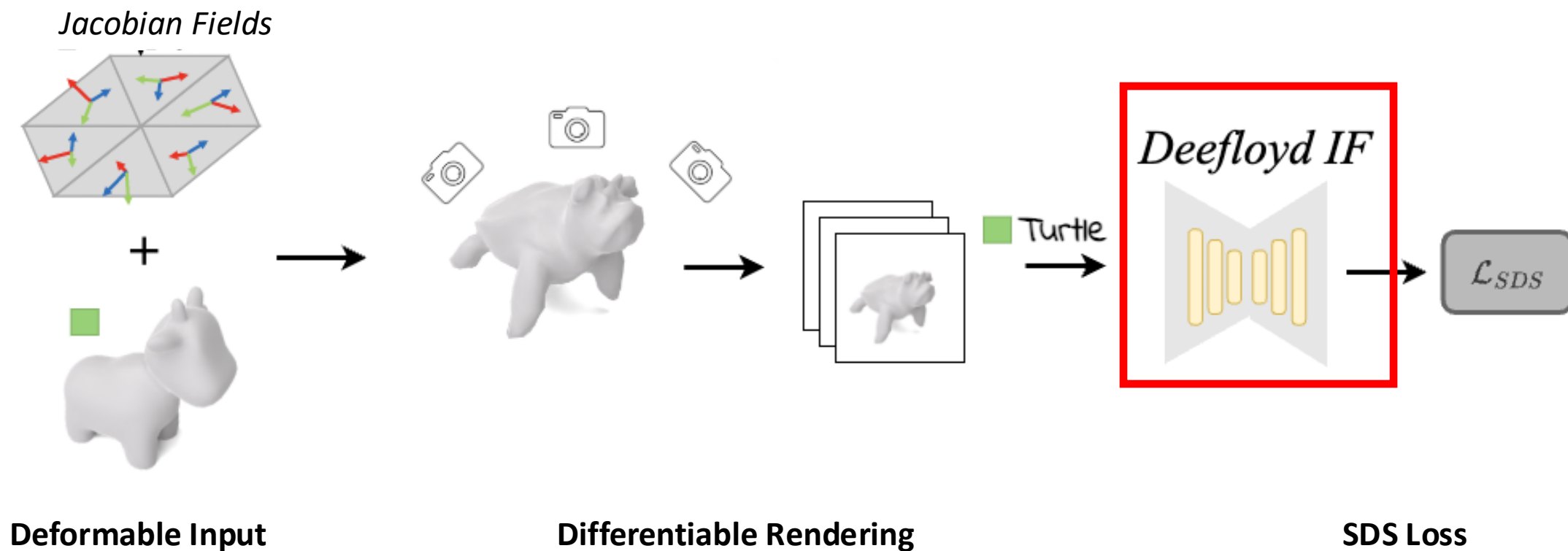
- STAGE 2: Project to "plausible" using SDS

Source   ARAP [50]   Ours

Adobe

# Source Mesh

# Deforming with Language Controls

- Condition on text



*Jacobian Fields*

**Deformable Input**     **Differentiable Rendering**     **SDS Loss**
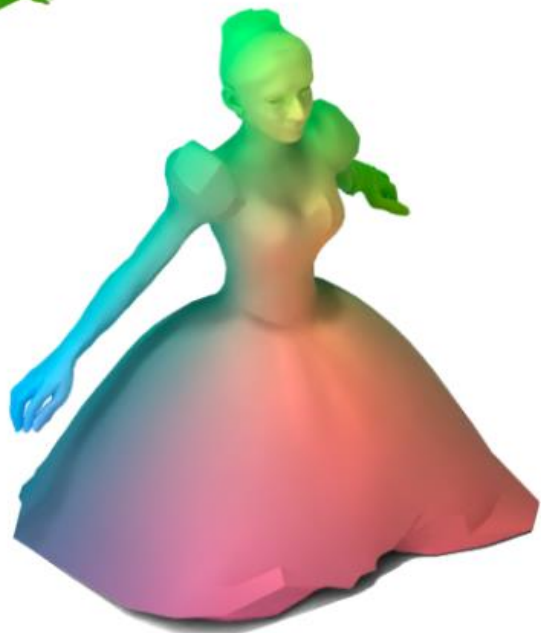
# Deforming with Language Controls



Source

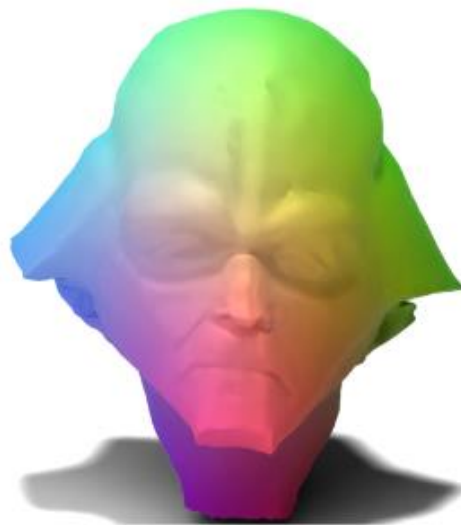Cinderella

Marge Simpson

Text

Sheep

# Deforming with Language Controls



Source

Satyr

vader
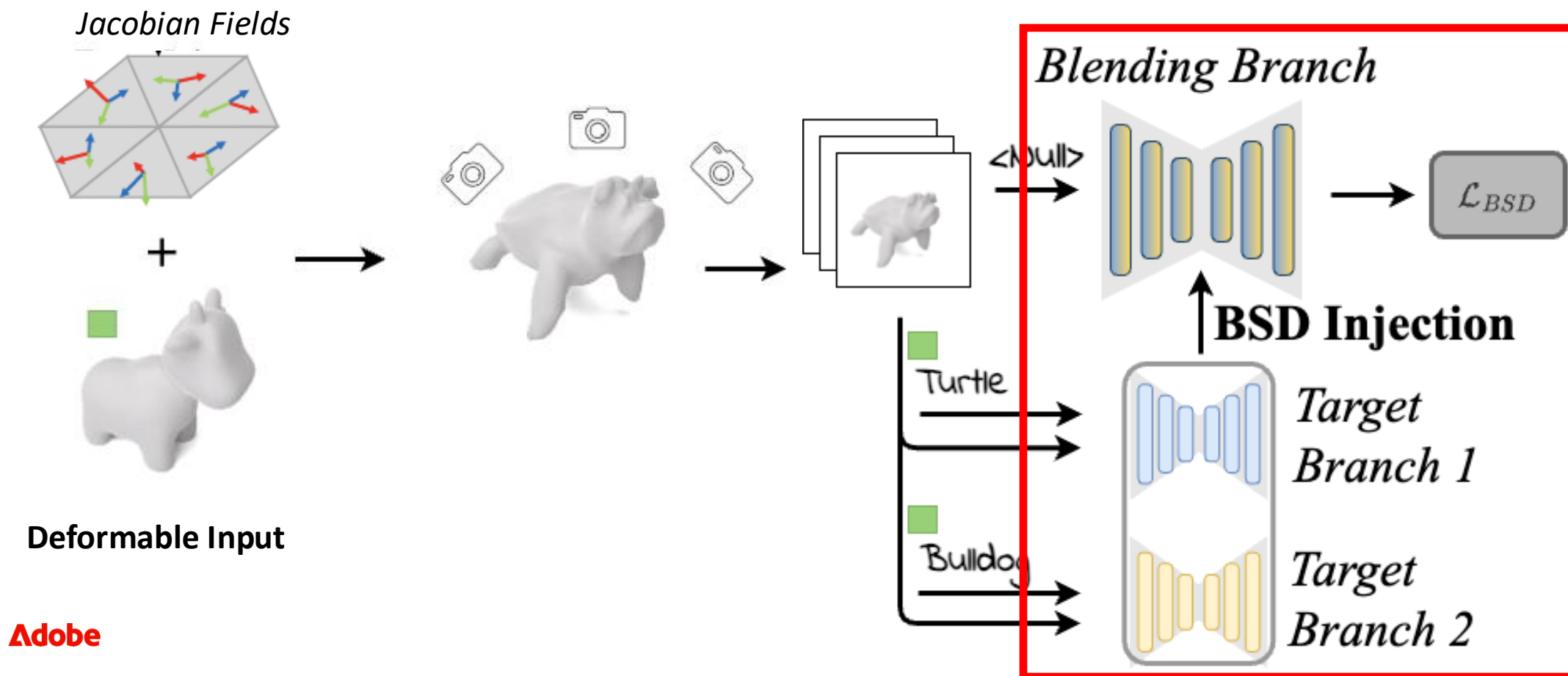
Elk

# Deforming with Language Controls: Multi-Target

- Inject weighted features into attention module



et al., 3DV 2025

# Deforming with Language Controls: Multi-Target

# Deforming with Language Controls: Multi-Target



WALLE-E 60%
Eve 30%

Cat woman 60%
Wonder Woman 40%

Jack Sparrow 100%

Luke Skywalker 60% Marge Simpson 100%
Dath vader 40%

Baku 100%

Gazelle 70%
Springbok 30%

Tiger 70%
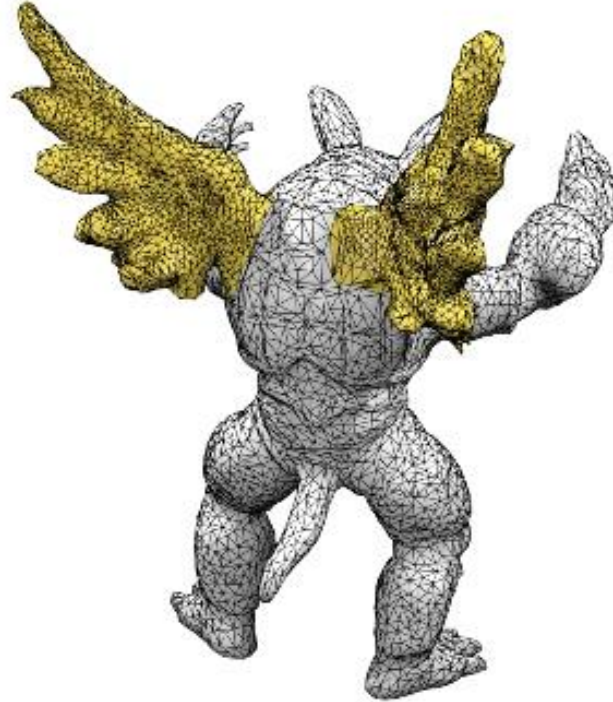Leopard 50%

Ra 70%
Anubis 30%

Warthog 70%
wildebeast 30%

# Overview

- Support mesh outputs (but use other representations as needed)

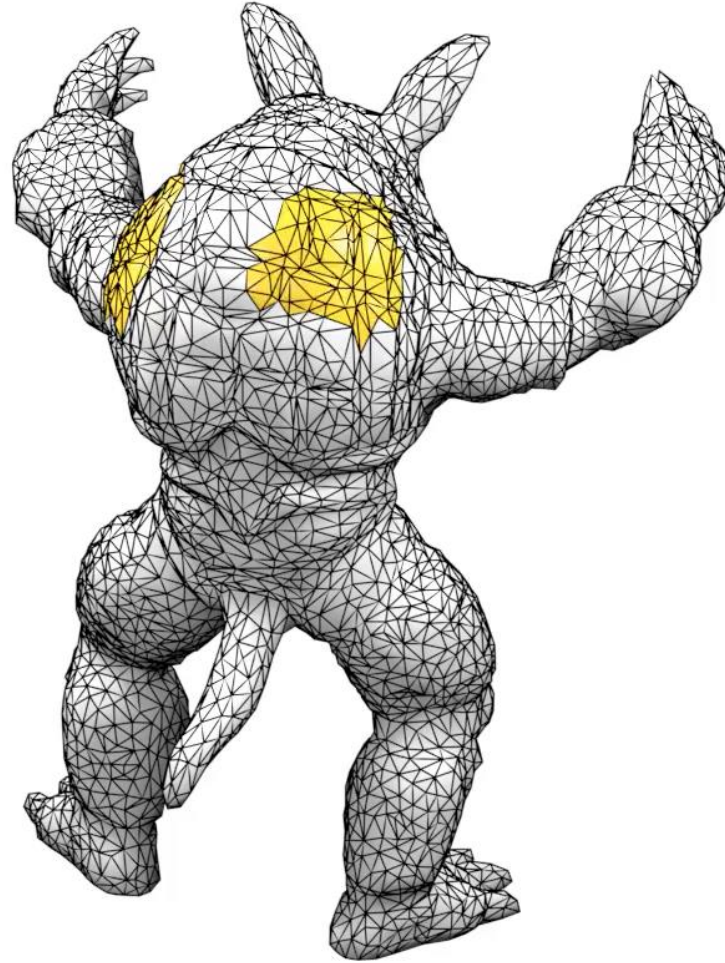- Inspired by traditional workflows
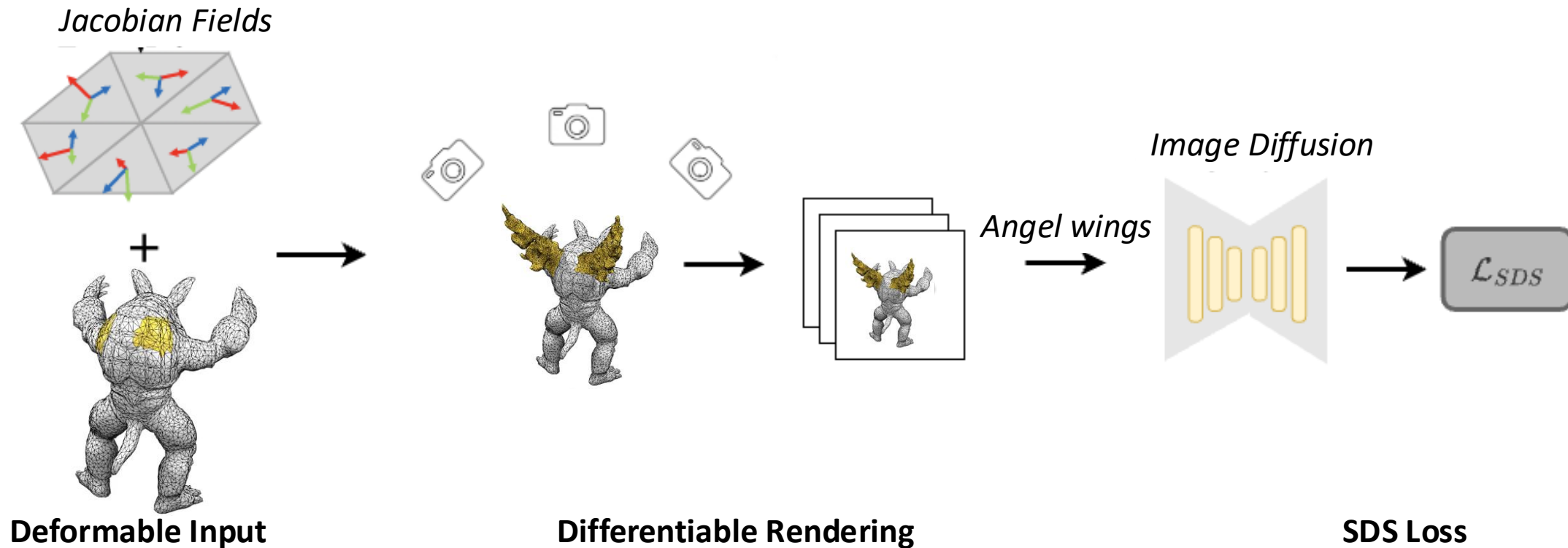


Neural **Deformation**

**Generative Scultping**
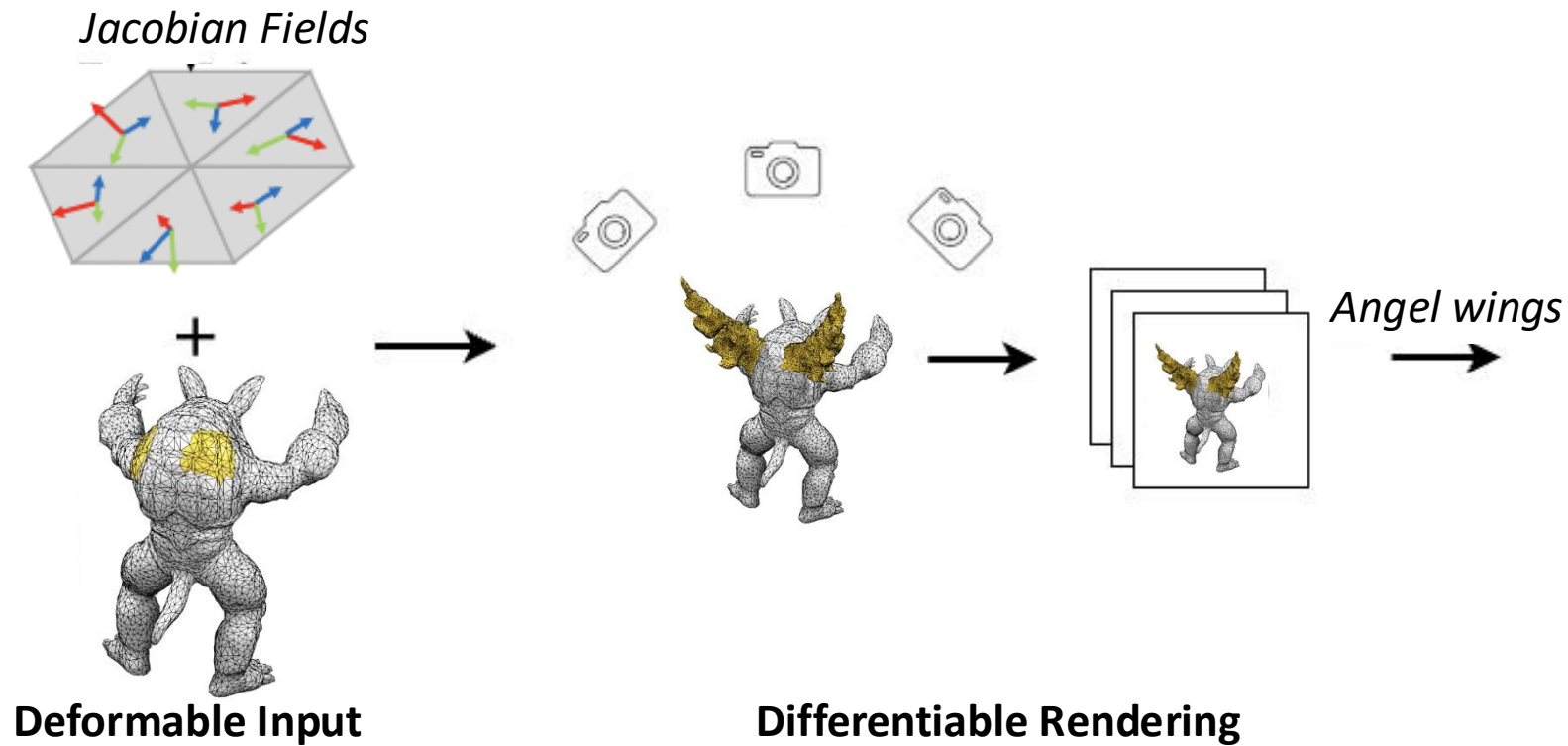
Generative **Detailization**

# Generative Sculpting



*"Man with angel wings"*

# Generative Sculpting via Deformation

- Only deform highlighted area



*Jacobian Fields*

**+**

*Image Diffusion*

*Angel wings*

$\mathcal{L}_{SDS}$

**Deformable Input**

**Differentiable Rendering**

**SDS Loss**

# Generative Sculpting via Deformation

▪ Only deform highlighted area



*Jacobian Fields*

**+**

**Deformable Input**

**Differentiable Rendering**

*Angel wings*

**Adobe**

# Prior Work: Continuous Remeshing

- Dynamic remeshing: allows to add more details



Jacobian Fields

+ −

Deformable Input

# Prior Work: Continuous Remeshing

- Dynamic remeshing: allows to add more details

- Remeshing + SDS
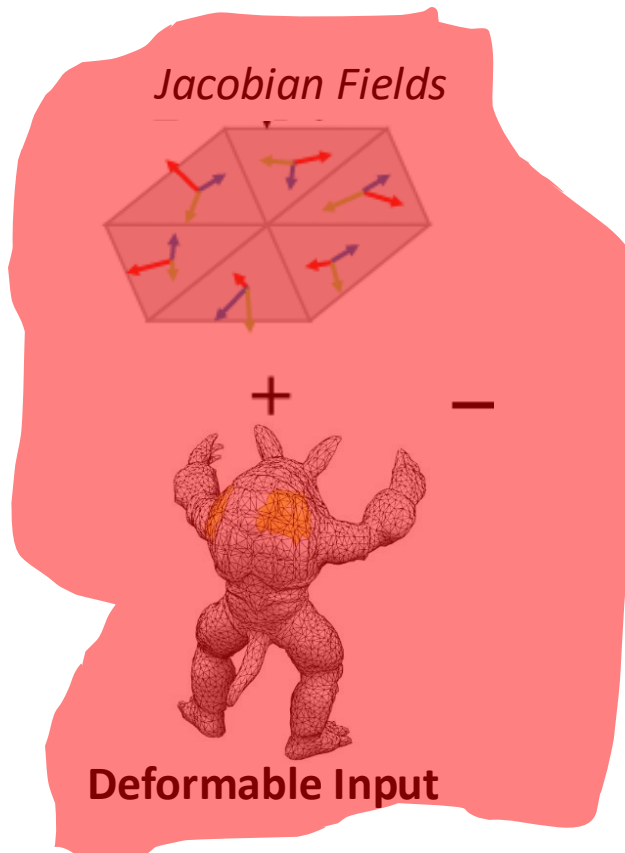


*Jacobian Fields*

\+   \-

**Deformable Input**

Prompt: "A deer"

# ~~Prior Work: Continuous Remeshing~~

- Use continuous remeshing instead?

- ~~Remeshing~~ + SDF + SDS



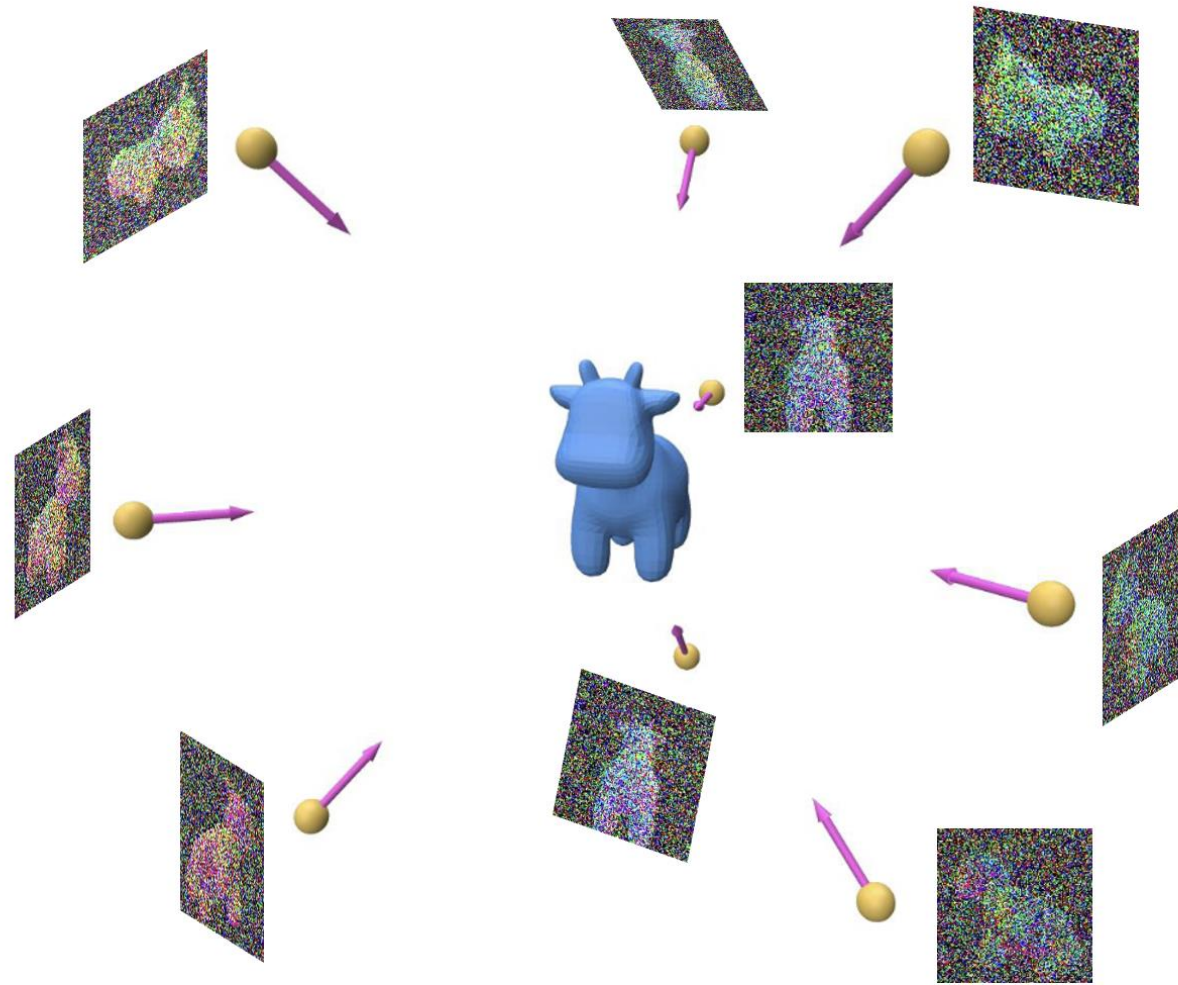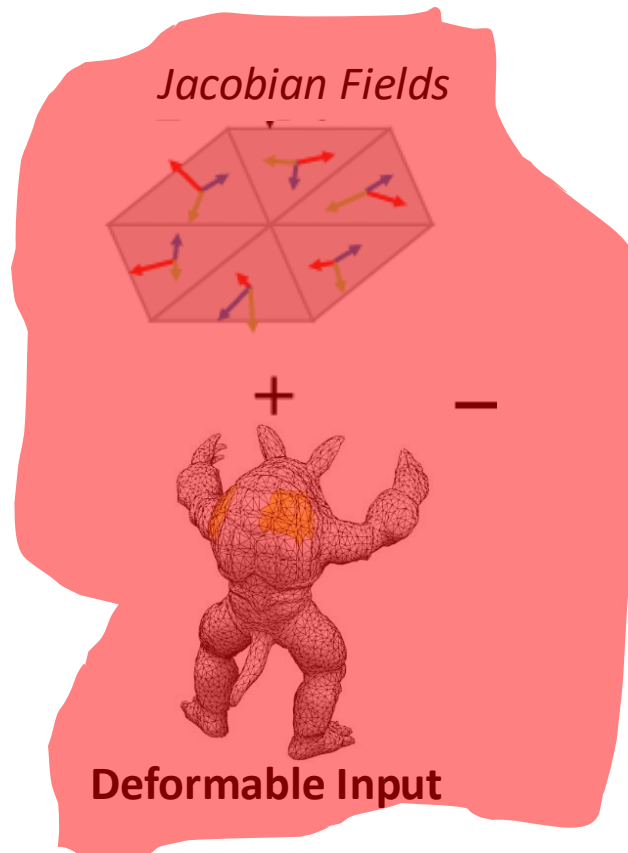*Jacobian Fields*

+ −

**Deformable Input**

Prompt: "A deer"



Adobe

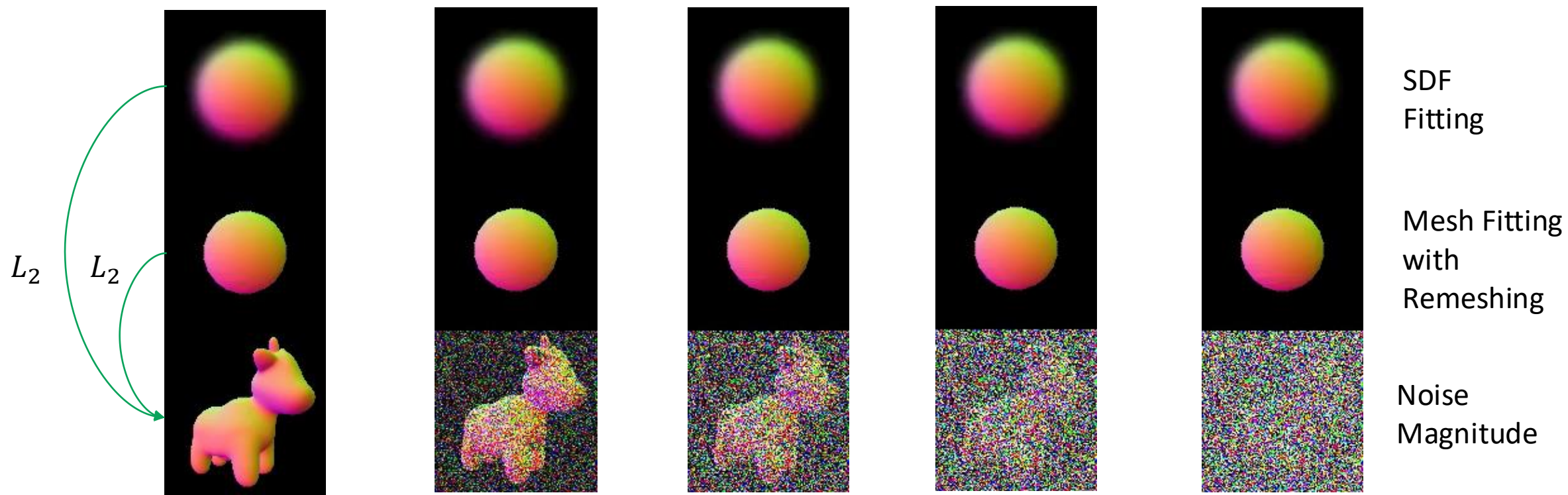# Controlled Experiment: Continuous Remeshing vs SDF

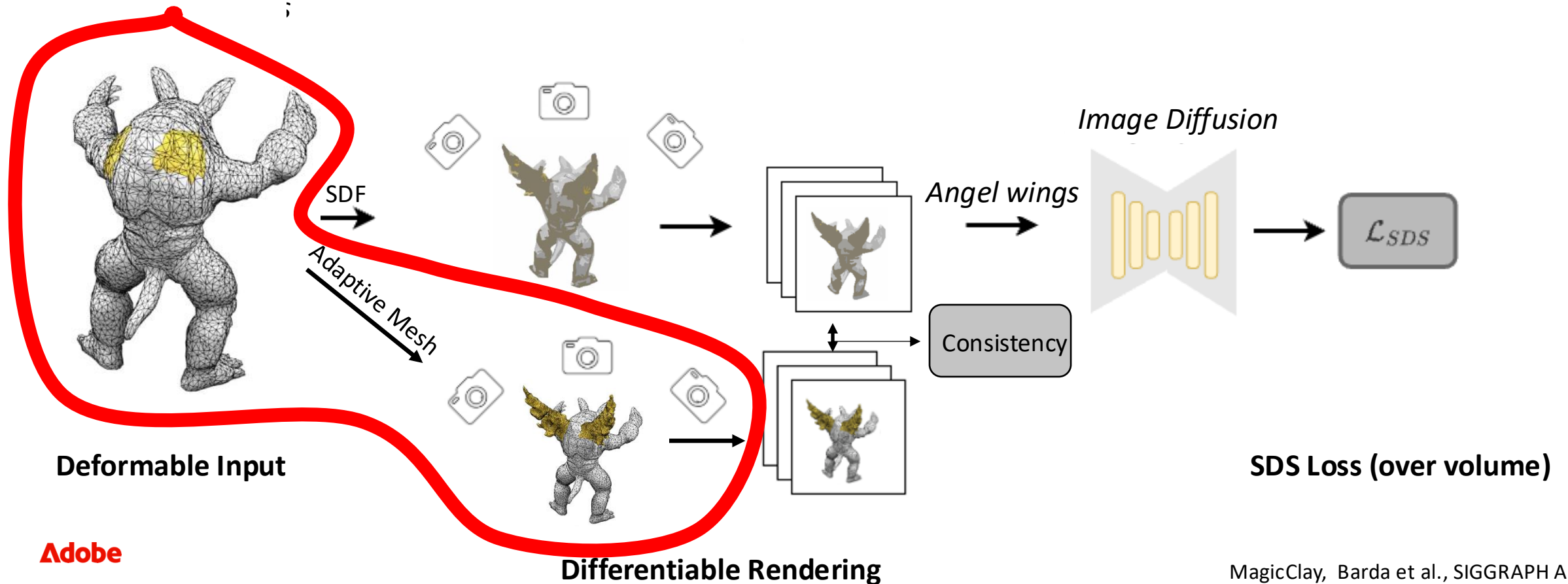- Reconstruct from renderings with different levels of noise



*Jacobian Fields*

+ −

**Deformable Input**

# Controlled Experiment: Continuous Remeshing vs SDF

▪ Brittle



$L_2$ $L_2$

SDF Fitting

Mesh Fitting with Remeshing

Noise Magnitude

# Generative Sculpting via Deformation: MagicClay

- **Dynamic remeshing: allows to add more details**

- Hybrid (SDF+Mesh) representation



SDF

Adaptive Mesh

Angel wings

*Image Diffusion*

Consistency

$\mathcal{L}_{SDS}$

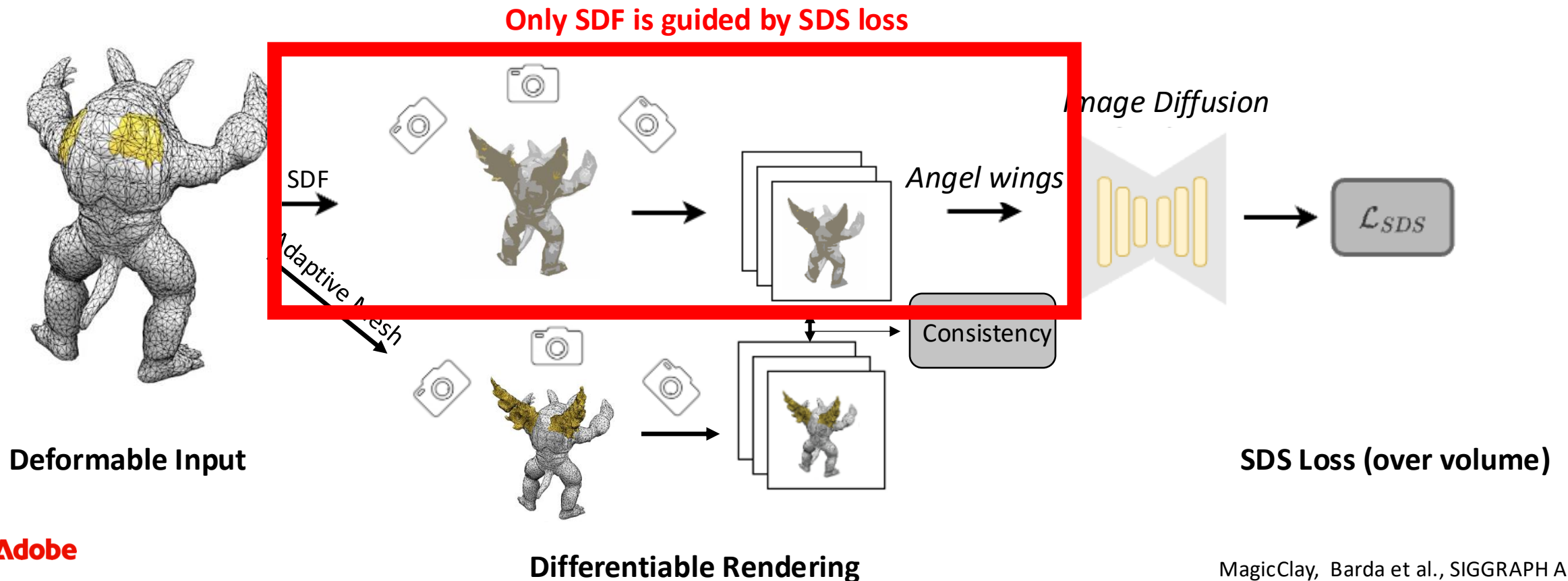**Deformable Input**
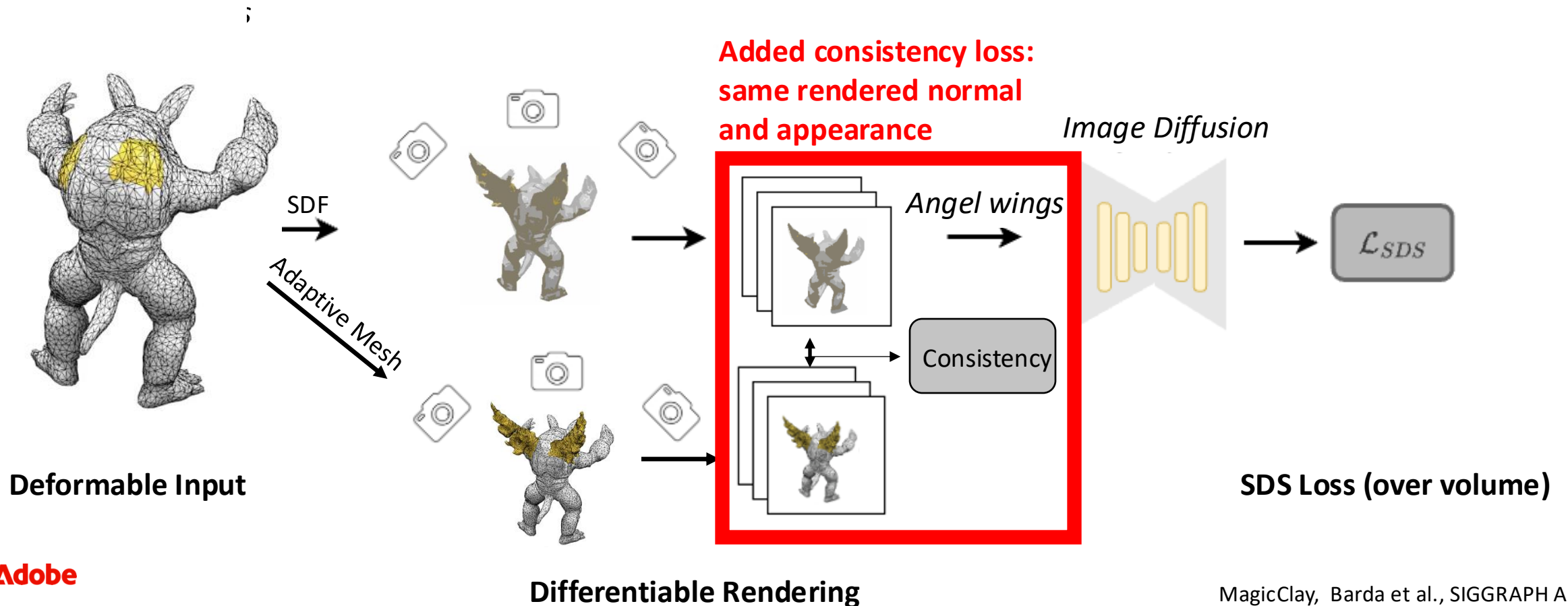
**Differentiable Rendering**

**SDS Loss (over volume)**

# Generative Sculpting via Deformation: MagicClay

- Dynamic remeshing: allows to add more details

- **Hybrid (SDF+Mesh) representation**



**Only SDF is guided by SDS loss**

SDF

*Adaptive Mesh*

*Image Diffusion*

*Angel wings*

Consistency

$\mathcal{L}_{SDS}$

**Deformable Input**

**Differentiable Rendering**

**SDS Loss (over volume)**

# Generative Sculpting via Deformation: MagicClay

- Dynamic remeshing: allows to add more details

- **Hybrid (SDF+Mesh) representation**



**Added consistency loss: same rendered normal and appearance**

*Image Diffusion*

*Angel wings*

Consistency

$\mathcal{L}_{SDS}$

SDF

Adaptive Mesh

**Deformable Input**

**Differentiable Rendering**

**SDS Loss (over volume)**

# Generative Sculpting via Deformation: MagicClay



"Man holding a...                    ...knight's sword"                    ...Wizard staff"                    ...maraca"
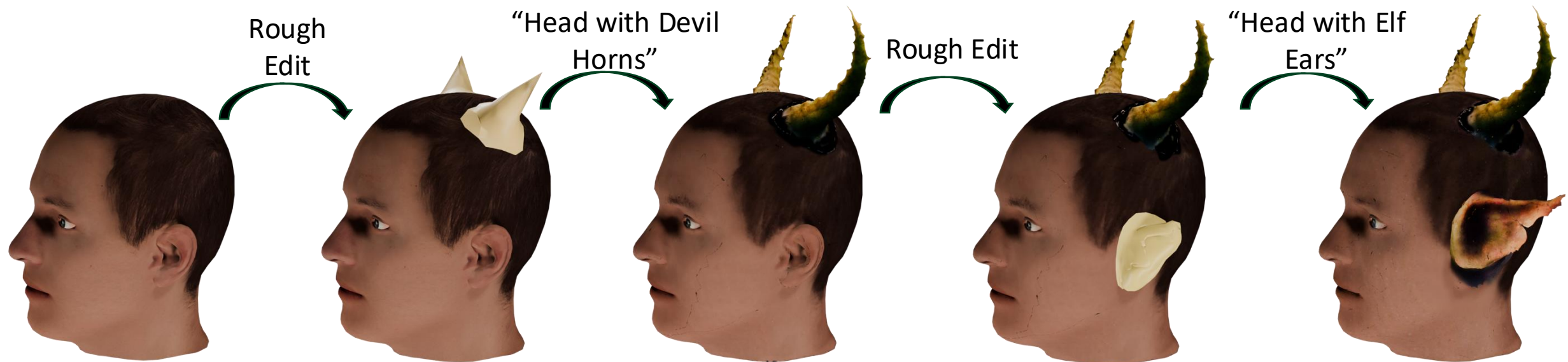
# SDS Limitations

- Brittle (requires careful tuning of hyperparameters)
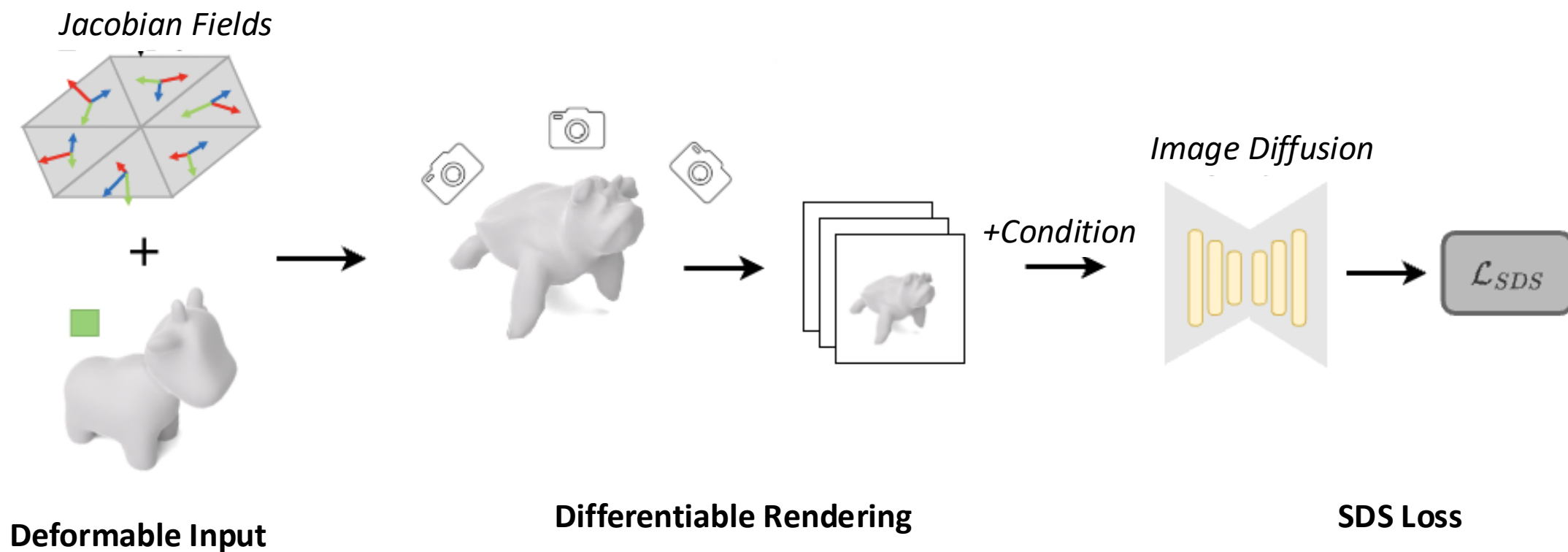
- **Inconsistent guidance from different views**

# SDS Limitations

- Brittle (requires careful tuning of hyperparameters)

- **Inconsistent guidance from different views**
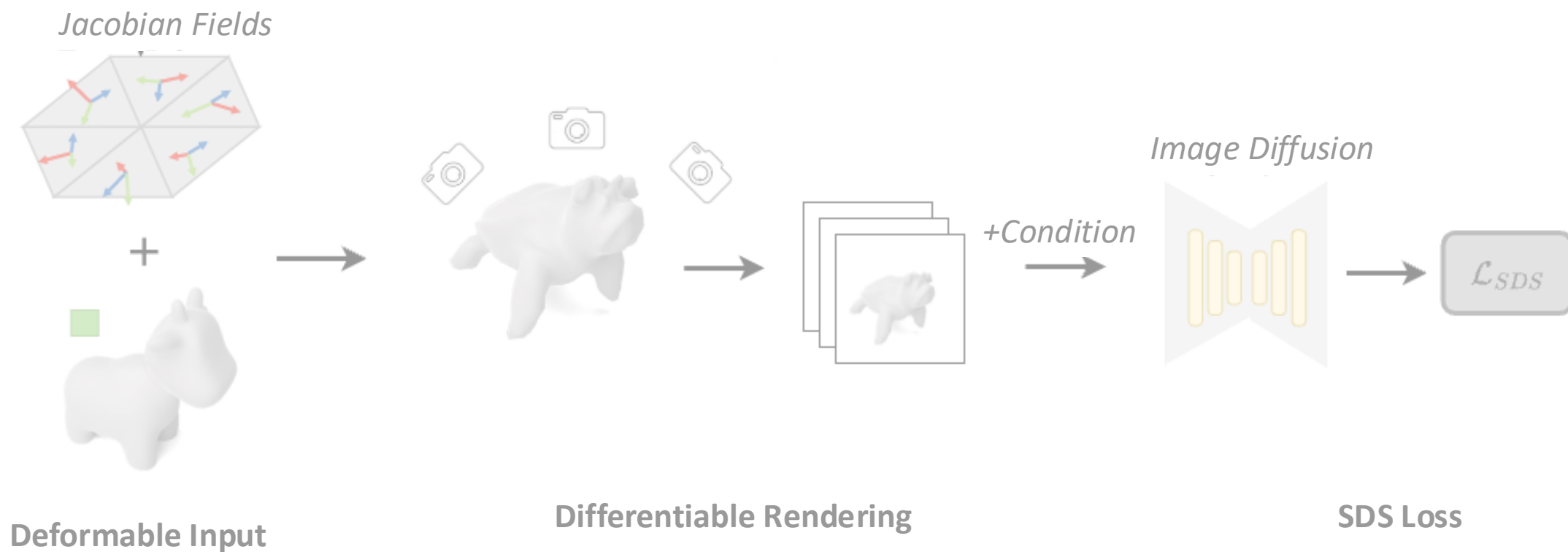
- **Slow (hours per iteration)**



Rough Edit → "Head with Devil Horns" → Rough Edit → "Head with Elf Ears"

MagicClay, Barda et al., SIGGRAPH A 2024

# Getting Rid of SDS

- Reminder: SDS pipeline



*Jacobian Fields*

*Image Diffusion*

+Condition

$\mathcal{L}_{SDS}$

**Deformable Input**

**Differentiable Rendering**

**SDS Loss**

# Generative Sculpting via Multi-view Inpainting

- ~~Reminder: SDS pipeline~~



*Jacobian Fields*

*Image Diffusion*

*+Condition*

$\mathcal{L}_{SDS}$

**Deformable Input**

**Differentiable Rendering**

**SDS Loss**

# Generative Sculpting via Multi-view Inpainting

- In-paint multi-view

- Reconstruct



*Image Diffusion*

$\mathcal{L}_{SDS}$

**Input**

~~Differentiable~~ **Rendering (masks)**

**SDS Loss**

# Generative Sculpting via Multi-view Inpainting

- In-paint multi-view

- Reconstruct



"An astronaut riding a rocking horse"

*Image Diffusion*

**Input**

~~Differentiable~~ **Rendering (masks)**

~~SDS Loss~~

Instant3dit, Barda et al., In submission, 2025

# Generative Sculpting via Multi-view Inpainting

- In-paint multi-view

- Reconstruct

"An astronaut riding a rocking horse"



**Input**          **Rendering (masks)**          **Fine-tune**          **Multi-view**

*Image Diffusion*

# Generative Sculpting via Multi-view Inpainting

- In-paint multi-view

- Reconstruct



"An astronaut riding a rocking horse"

Image Diffusion

Option 2: Reconstruct via LRM

**Input**          **Rendering (masks)**          **Fine-tune**          **Multi-view**

Option 1: Reconstruct via Differentiable Renderer optimization

# Generative Sculpting via Multi-view Inpainting

- In-paint multi-view

- Reconstruct



"An astronaut riding a rocking horse"

*Image Diffusion*

Option 2: Reconstruct via LRM

**Input**          **Rendering (masks)**          **Fine-tune**          **Multi-view**

Option 1: Reconstruct via Differentiable Renderer optimization

No off-she shelf multi-view inpainting

"An astronaut riding a rocking horse"

Option 2: Reconstruct via LRM

*Image Diffusion*

**Input**

**Rendering (masks)**

**Fine-tune**

**Multi-view**

Option 1: Reconstruct via Differentiable Renderer optimization

Instant3dit, Barda et al., In submission, 2025

# Training Data

- 5000 high quality meshes from the Objaverse dataset

- Each rendered from 4 canonical views

- Detailed Captions created using LLaVa



"A bronze sculpture of a mythical creature with intricate detailing, featuring a dragon-like creature with scales and a flowing mane, perched atop a round base adorned with additional decorative elements. The surface is textured and ornate patterns embellish the pedestal."



"A low-poly cannon model with a textured surface, featuring a cylindrical barrel and a wooden carriage with large wheels. The cannon is mounted on a stone base, which has a small round platform in front of it. The overall design suggests an old-fashioned or historical style."

# I: Coarse Edit

- Scenario: "user creates a large chunk with a coarse proxy"

# I: Coarse Edit

- Scenario: "user creates a large chunk with a coarse proxy"

- 3D mask is a polyhedron fully containing the region to be inpainted

# II: Mesh Sculpting

- Scenario: "user provides a more precise approximation to the target"

# II: Mesh Sculpting

- Scenario: "user provides a more precise approximation to the target"

- 3D mask is the exact mesh that needs to be inpainted

# III: Surface Edit

- Scenario: "user provides wants local geometry and texture modifications"

# III: Surface Edit

- Scenario: "user provides wants local geometry and texture modifications"

- 3D mask is a small region of the exact mesh

# Training Multi-View Inpainting

- Train from scratch on 5B Multi-view inpainted images?

# Training Multi-View Inpainting

- ~~Train from scratch on 5B Multi-view inpainted images?~~

- Option 1: fine-tune multi-view image generation to teach it to inpaint



Opt 1: Instant3D

# Training Multi-View Inpainting

- ~~Train from scratch on 5B Multi-view inpainted images?~~

- Option 1: fine-tune multi-view image generation to teach it to inpaint

- Option 2: fine-tune generative image inpainting to teach it to create multi-view



Opt 1: Instant3D

Opt 2: SDXL-Inpaint

# Training Multi-View Inpainting

**Which one is a better base model?**

- ~~Train from scratch on 5B Multi-view inpainted images?~~

- Option 1: fine-tune multi-view image generation to teach it to inpaint

- Option 2: fine-tune generative image inpainting to teach it to create multi-view



Opt 1: Instant3D

Opt 2: SDXL-Inpaint

# Training Multi-View Inpainting

## Which one is a better base model?

- ~~Train from scratch on 5B Multi-view inpainted images?~~

- ~~Option 1: fine-tune multi-view image generation to teach it to inpaint~~

- Option 2: fine-tune generative image inpainting to teach it to create multi-view



Opt 1: Instant3D

Opt 2: SDXL-Inpaint

Instant3dit, Barda et al., In submission, 2025

# Ablation: choose the right backbone

- What if we start with multi-view backbone instead of inpainting backbone?



Input

Mask

"A whale wrapped in a pink bow"
Prompt

Instant3D Finetuned

SDXL IP Finetuned
(Ours)

be Confidential.

# Ablation: choose the right backbone

- What if we start with multi-view backbone instead of inpainting backbone?



| | Input |
| Mask |
| "A whale wrapped in a pink bow" Prompt |

**Benchmark Stats**

| | | | |
|---|---|---|---|
| Prompt Adherance (CLIPL): | 28.57 | vs | 29.01 |
| Multi-View Consistency (DreamSim): | 0.97 | vs | 0.100 |
| Visual Quality (FID): | 121.1 | vs | 118.4 |

Instant3D Finetuned

**SDXL IP Finetuned (Ours)**

# Ablation: use multiple masks

- What if we train only on one type of mask?



| Input | Random 2d | Type I only | Type II only | Type III only | All (ours) |
|-------|-----------|-------------|--------------|---------------|------------|
| Mask  |           |             |              |               |            |

**Benchmark Stats (FID)**

"A cow with a baseball cap"
Prompt

| | Random 2d | Type I only | Type II only | Type III only | All (ours) |
|---|---|---|---|---|---|
| | 131.1 | 121.3 | 128.2 | 142.2 | 118.4 |

# Interactive Generative Sculpting

speed 5x

Instant3dit demo.

Adobe

# Results: support all representations

- Just swap different LRM models (or optimize mesh)



Input

Mask

"A bear holding a honey pot" Prompt

Gaussian Splats (3s)

Mesh (6s)

Adaptive Remeshing (25s)

# Results: support all representations
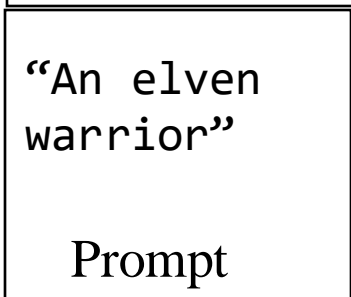
▪ Just swap different LRM models (or optimize mesh)



Input

Mask

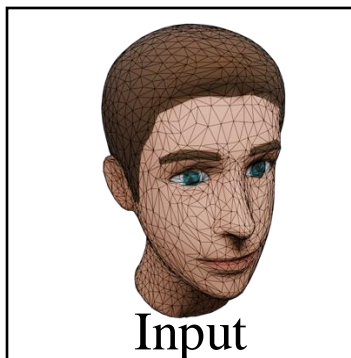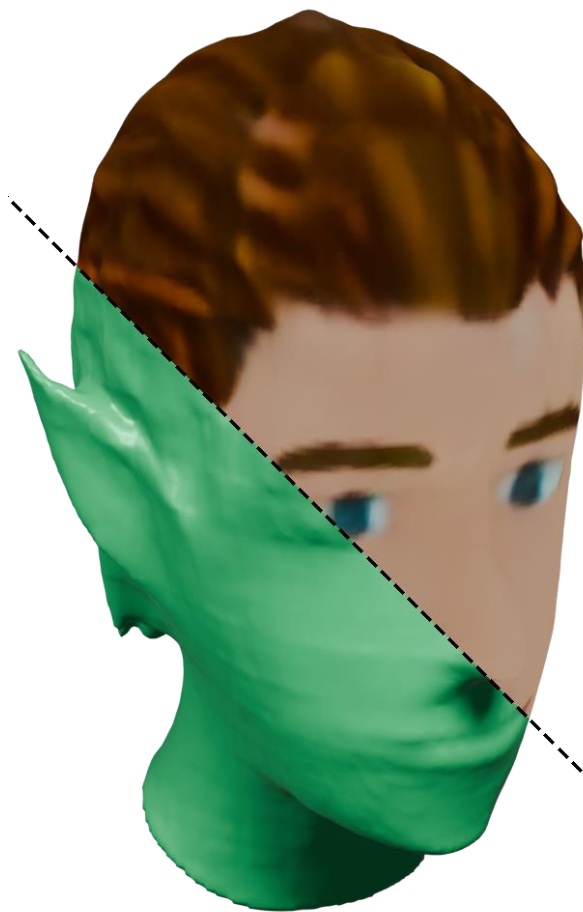"A bear
with wings"

Prompt

Gaussian Splats (3s)

Mesh (6s)

Adaptive Remeshing (25s)

# Results: support all representations

▪ Just swap different LRM models (or optimize mesh)



Input

Mask

"Man wearing a medieval helmet" Prompt

NeRF (3s)

Mesh (6s)

Adaptive Remeshing (25s)

# Results: support all representations

- Just swap different LRM models (or optimize mesh)



| Input | Mask | "An elven warrior" Prompt |

NeRF (3s)          Mesh (6s)          Adaptive Remeshing (25s)

Instant3dit, Barda et al., In submission, 2025

# Overview

- Support mesh outputs (but use other representations as needed)

- Inspired by traditional workflows



Neural **Deformation**

Generative **Scultping**

Generative **Detailization**

# Generate or Change Details



>_: A teddy bear

Input 3D mesh             Input text             Output 3D mesh

# Multi-view Detail Generation
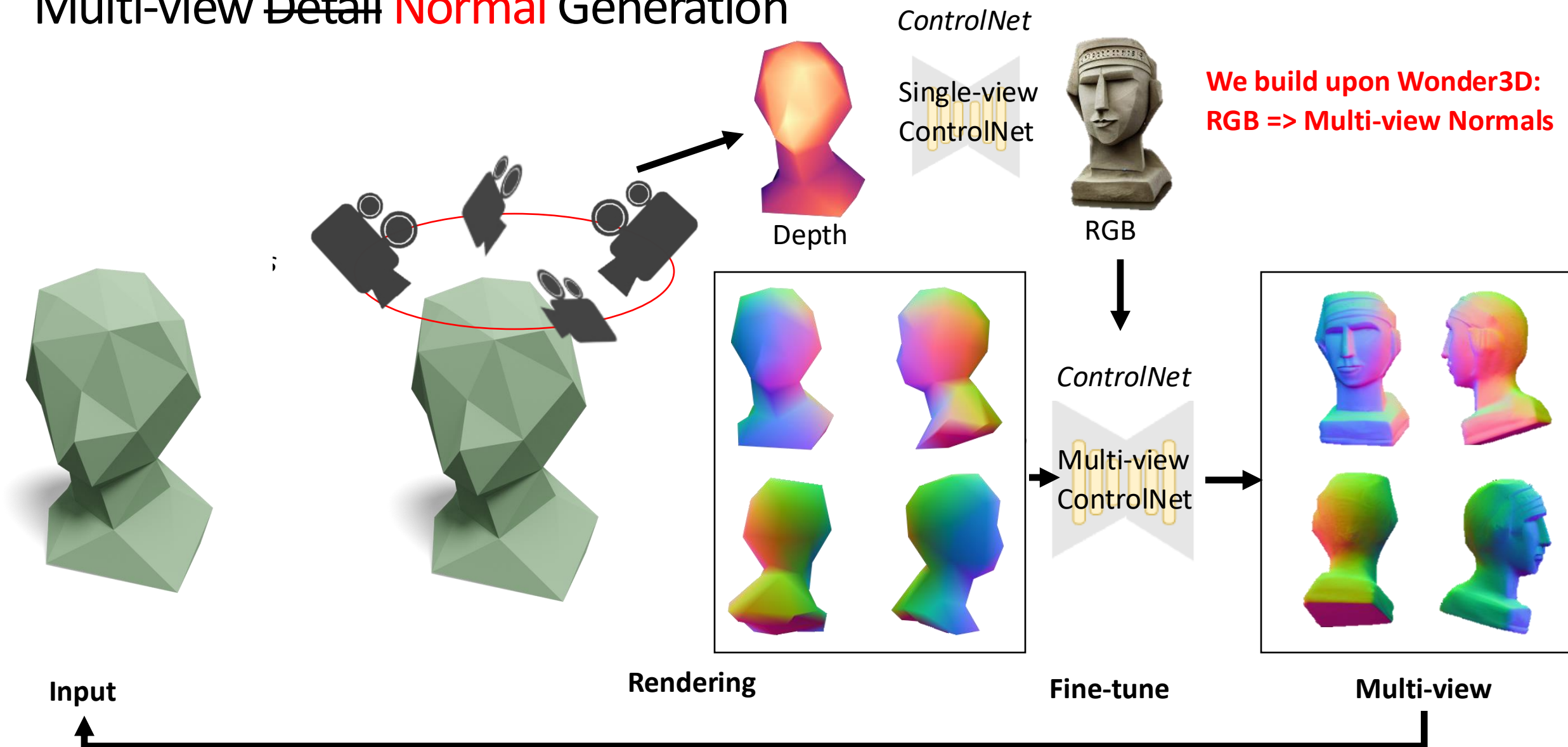


"A stone figure"

*Image Diffusion*

**Input**  **Rendering**  **Fine-tune**  **Multi-view**

Option 1: Reconstruct via Differentiable Renderer optimization

# Multi-view ~~Detail~~ Normal Generation



Input

Rendering

"A stone figure"

Image Diffusion

Fine-tune

Multi-view

Option 1: Reconstruct via Differentiable Renderer optimization

Text-Guided Refinement, Yun-Chun Chen et al., SIGGRAPH A, 2024

# Multi-view ~~Detail~~ Normal Generation



**Multi-view normal sharpening does not exist**

"A stone figure"

*Image Diffusion*

**Input**   **Rendering**   **Fine-tune**   **Multi-view**

Option 1: Reconstruct via Differentiable Renderer optimization

# Multi-view ~~Detail~~ Normal Generation



ControlNet

Single-view ControlNet

Depth

RGB

**We build upon Wonder3D: RGB => Multi-view Normals**

ControlNet

Multi-view ControlNet

**Input**     **Rendering**     **Fine-tune**     **Multi-view**

Option 1: Reconstruct via Differentiable Renderer optimization

Adobe

# Multi-view ~~Detail~~ Normal Generation



**Vanilla ControlNet**

*ControlNet*

Single-view ControlNet

Depth    RGB

*ControlNet*

Multi-view ControlNet

Input    Rendering    Fine-tune    Multi-view

Option 1: Reconstruct via Differentiable Renderer optimization

Text-Guided Refinement,  Yun-Chun Chen et al., SIGGRAPH A, 2024

Adobe

# Multi-view ~~Detail~~ Normal Generation

**Fine-tuned Wonder3D with additional input (blured normal)**

Single-view ControlNet

Depth

RGB

*ControlNet*

Multi-view ControlNet

**Blur Normals**

**Input**

**Rendering**

**Fine-tune**

**Multi-view**

Option 1: Reconstruct via Differentiable Renderer optimization

Text-Guided Refinement, Yun-Chun Chen et al., SIGGRAPH A, 2024

Adobe

# Mesh Optimization



Diff renderer            Normal rendering

Loss

# Mesh Optimization



Diff renderer          Normal rendering

**Gradient Descent**

# Mesh Optimization



Diff renderer       Normal rendering

Gradient Descent

>_ : A cartoon figure    >_ : A teddy bear    >_ : A teddy bear

Input 3D mesh    Input 3D mesh    Input 3D mesh

Output 3D mesh    Output 3D mesh    Output 3D mesh

# Mesh Texturing

>_ : A cartoon cat head



Input 3D mesh

Geometry

Texture

# More on Texturing

- Textures are not just colors



Mesh without textures

Generated texture

# More on Texturing

- Textures are not just colors



Mesh without textures

Material maps

(works with different illumination conditions)

# More on Texturing

- Multi-view generation

# More on Texturing

- VLLM-assisted retrieval



"cover": leather

"strap": fabric

"paper": paper

# More on Texturing – But How Do We Segment?

- VLLM-assisted retrieval

# Material-aware 3D Segmentation

- Fine-tune SAM2 (object segmentation in video) to segment materials

# Material-aware 3D Segmentation

- Fine-tune SAM2 (object segmentation in video) to segment materials



**Pre-trained VIT Image Encoder (frozen)**

# Material-aware 3D Segmentation

- Fine-tune SAM2 (object segmentation in video) to segment materials



**Fine-tune mask decoder with point selection conditioning (train on material data)**

# Material-aware 3D Segmentation

- Fine-tune SAM2 (object segmentation in video) to segment materials



**Fine-tune memory attention (improves consistency)**

# Synthetic Training Data for Material-aware Segmentation

- Per-frame masks

# Synthetic Training Data for Material-aware Segmentation

SAMa, Fischer et al., In Submission, 2025

# Object vs Material Segmentation

▪ SAM2 is not suitable for material-aware segmentation without fine-tuning
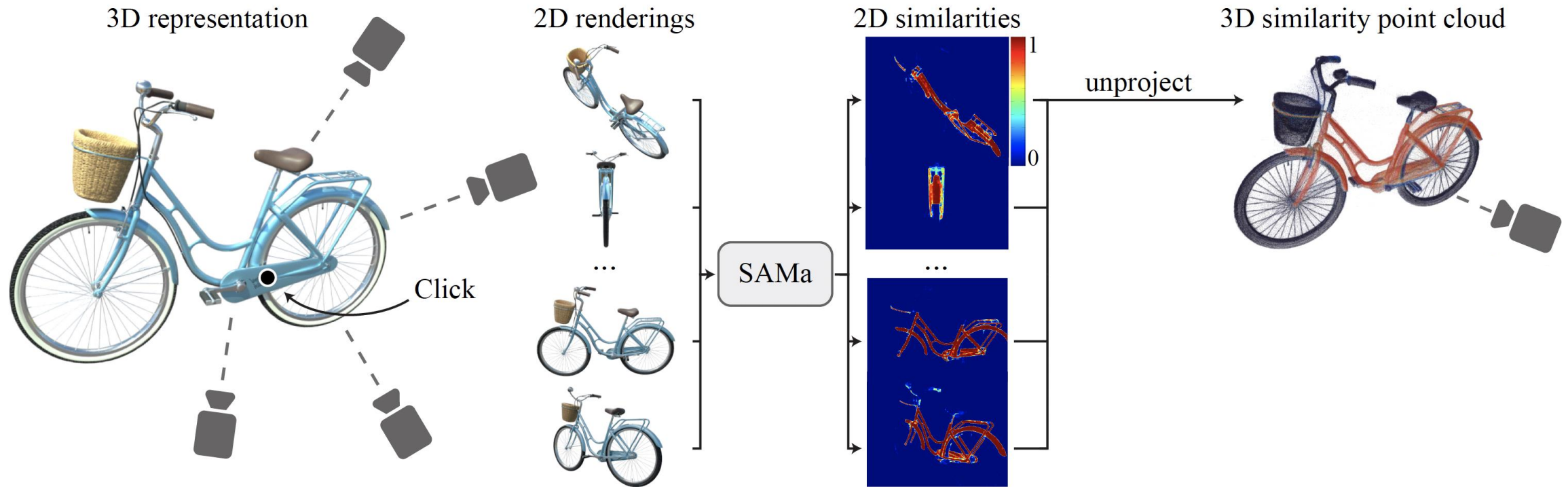
# Inferring 3D Segmentations

- Sample 25 views

# Inferring 3D Segmentations

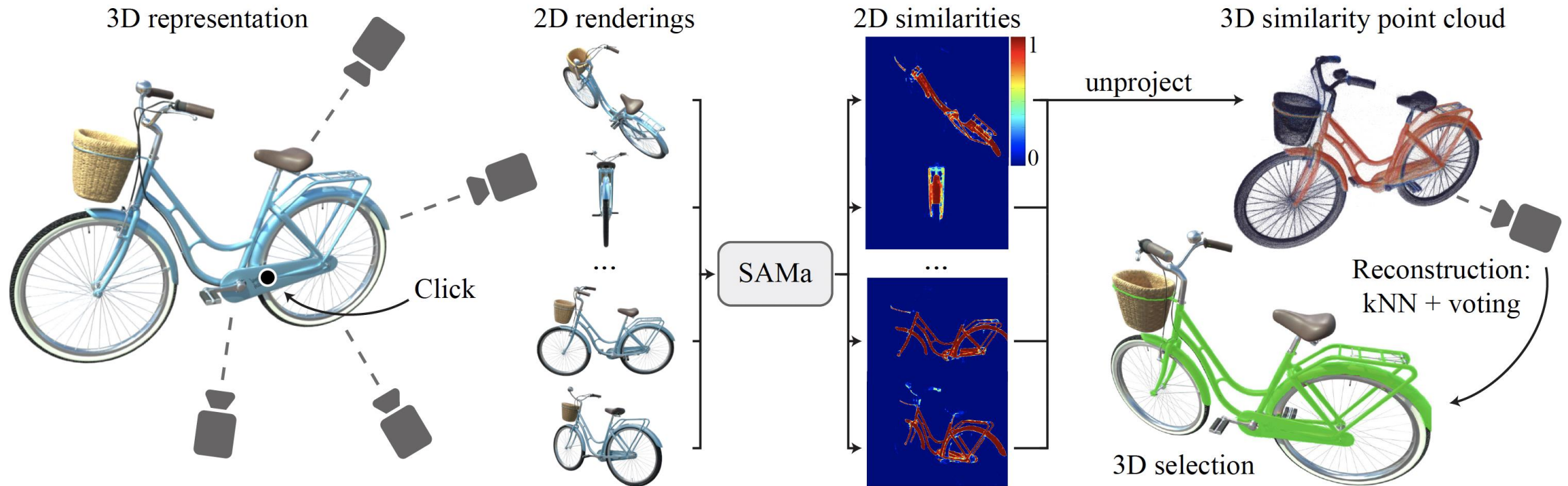- Sample 25 views => predict 2D features

# Inferring 3D Segmentations

- Sample 25 views => predict 2D features => store features in 3D point cloud



3D representation    2D renderings    2D similarities    3D similarity point cloud

Click

SAMa

unproject

# Inferring 3D Segmentations

▪ Sample 25 views => predict 2D features => store features in 3D point cloud

▪ Given user click => select via kNN voting

# Interactive Segmentation

- At runtime selection can be done in 2ms per click

# Representation-Agnostic

- Gaussian Splats

# Representation-Agnostic

- NeRFs

# Representation-Agnostic

- Meshes

# Automatic Segmentation
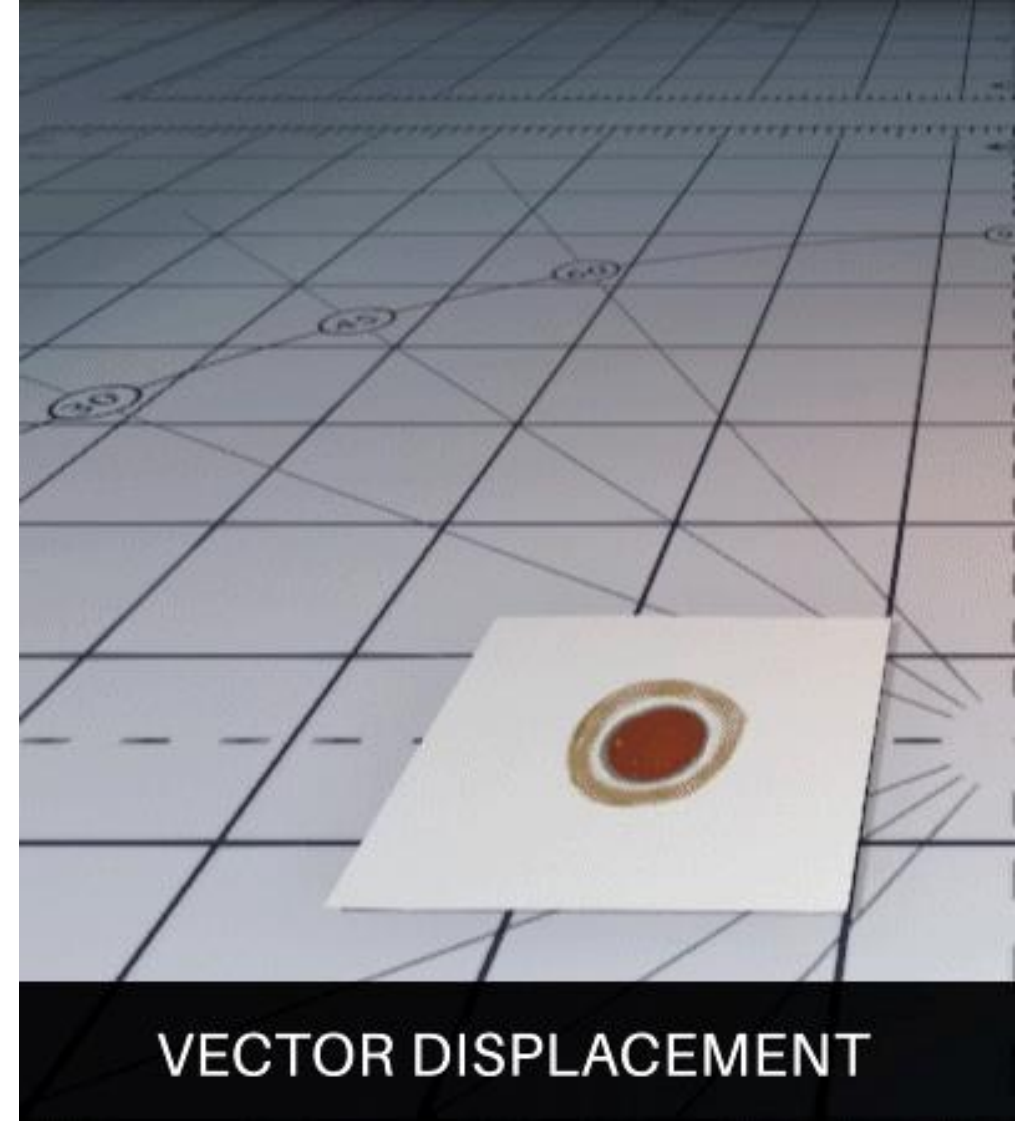
- Simulate many clicks => drop regions with high overlap

# Spatial Controlling the Details
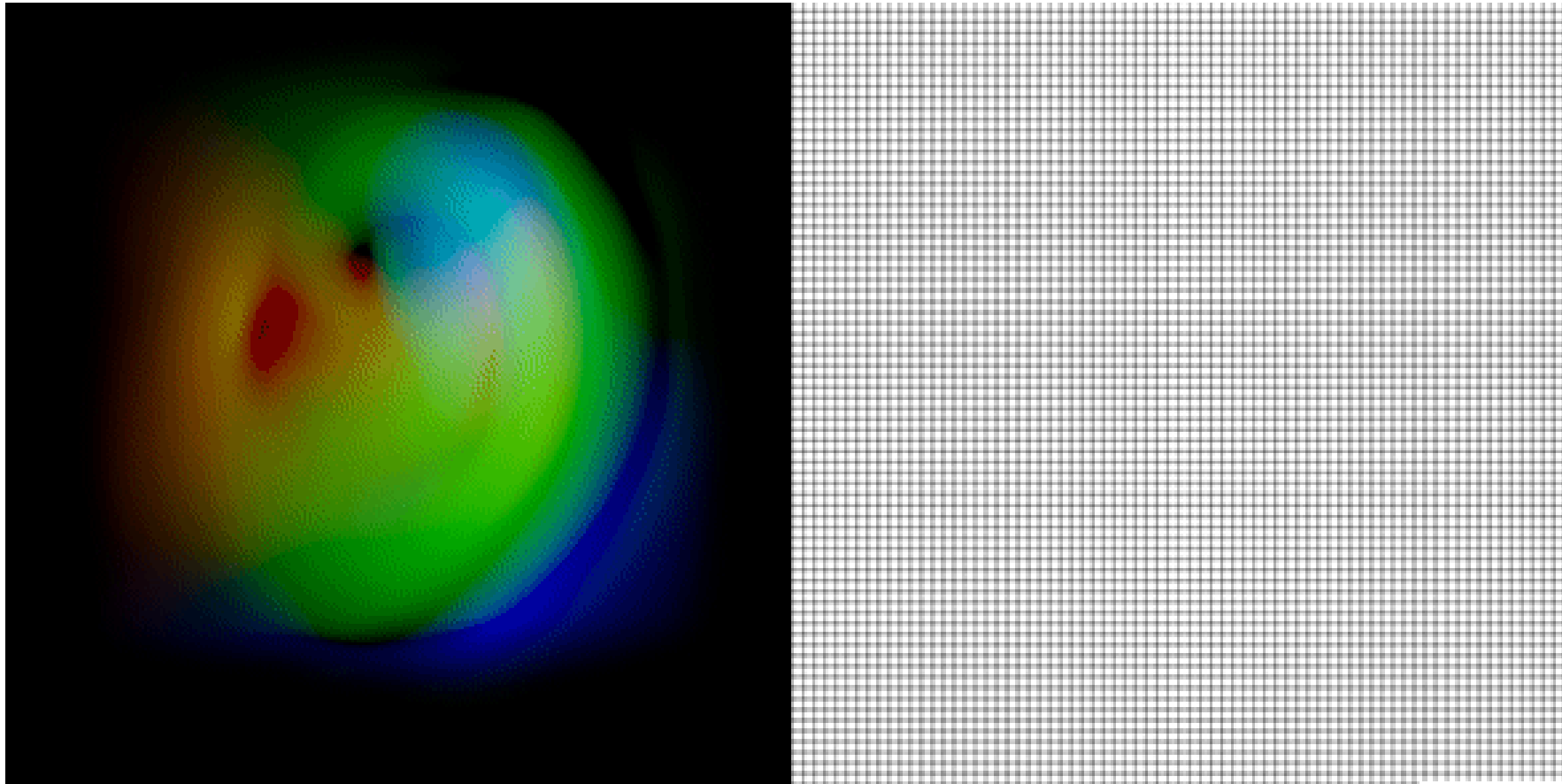
- Precisely position the ear on the head

# Vector Displacement Map (VDM)

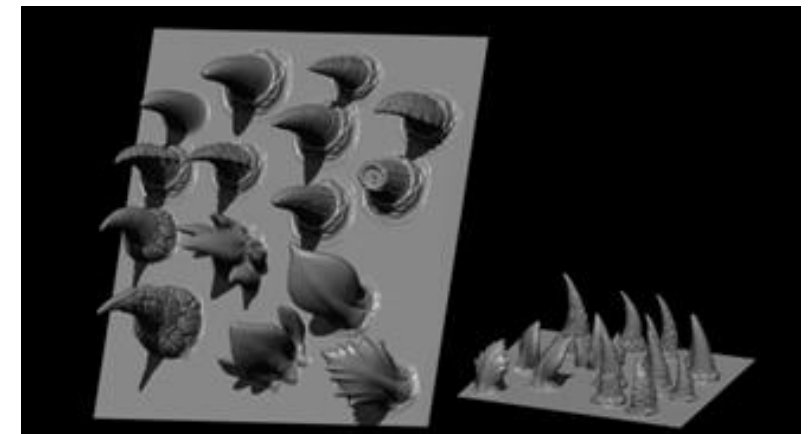- Map over 2D Plane: $f : [0,1]^2 \rightarrow \mathbb{R}^3$



VECTOR DISPLACEMENT

**Adobe**

# Vector Displacement Map (VDM)

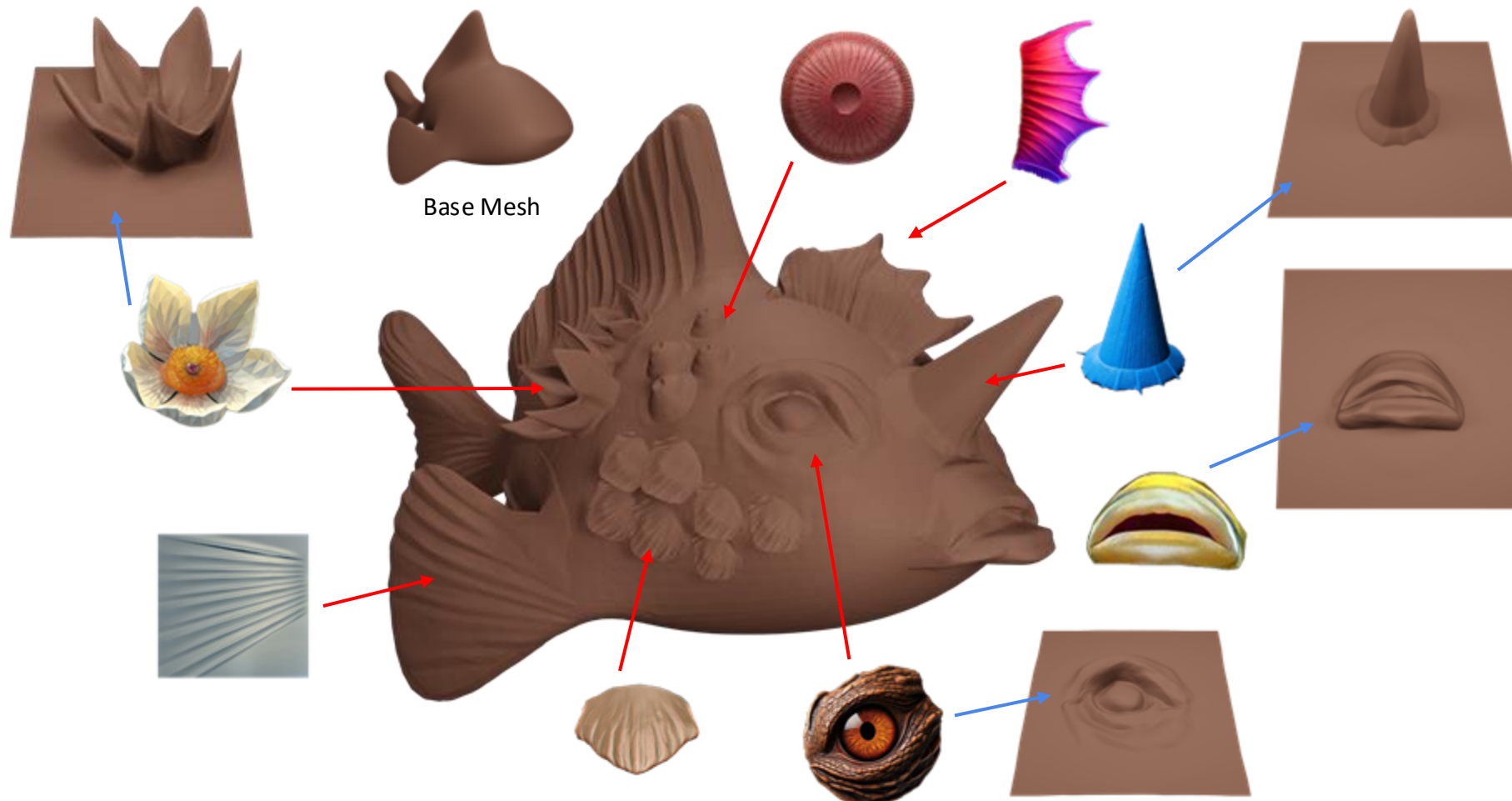- Map over 2D Plane: $f : [0, 1]^2 \rightarrow \mathbb{R}^3$

- Geometry Image

# Why do we care?

- Heavily used stock asset
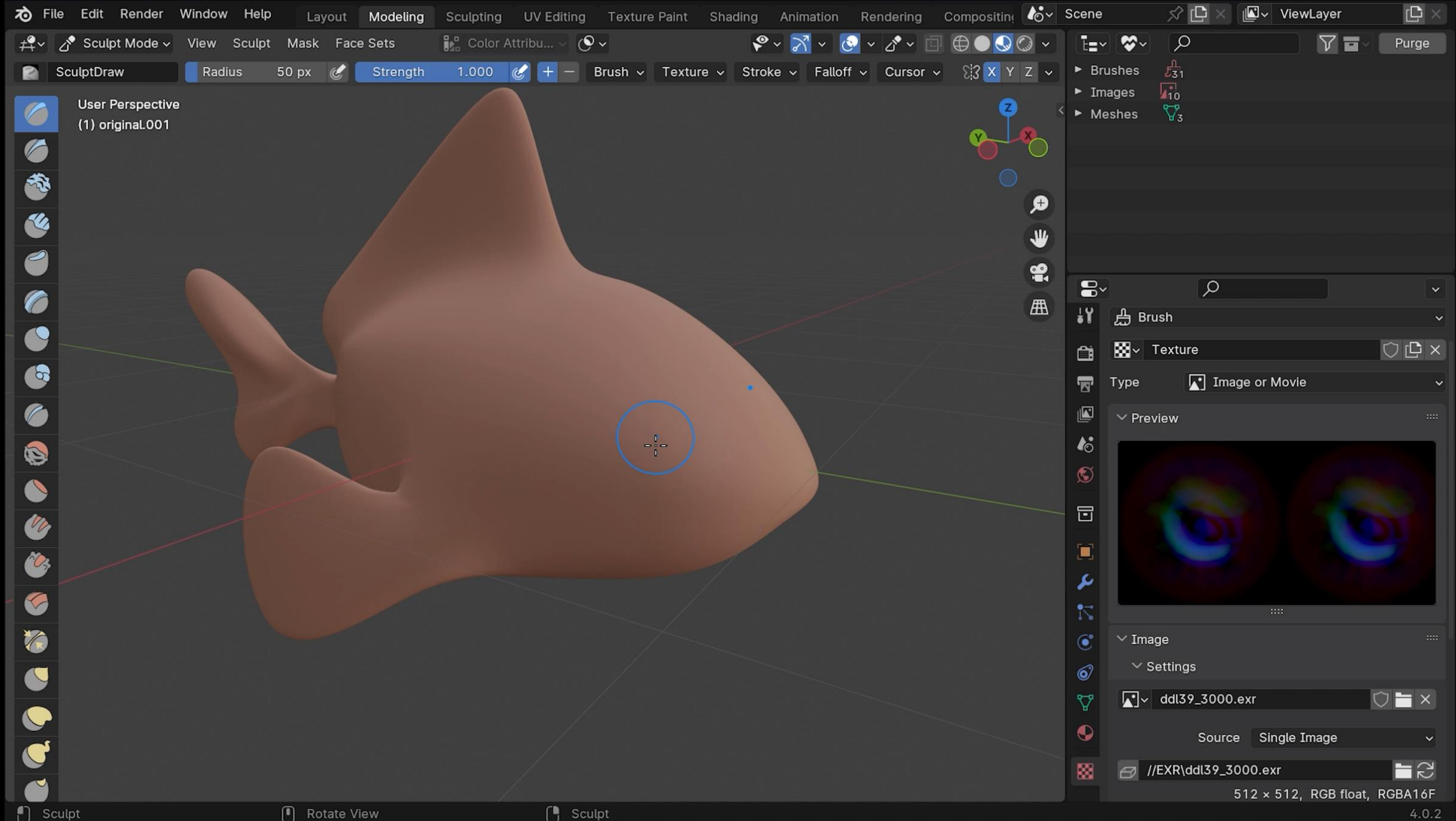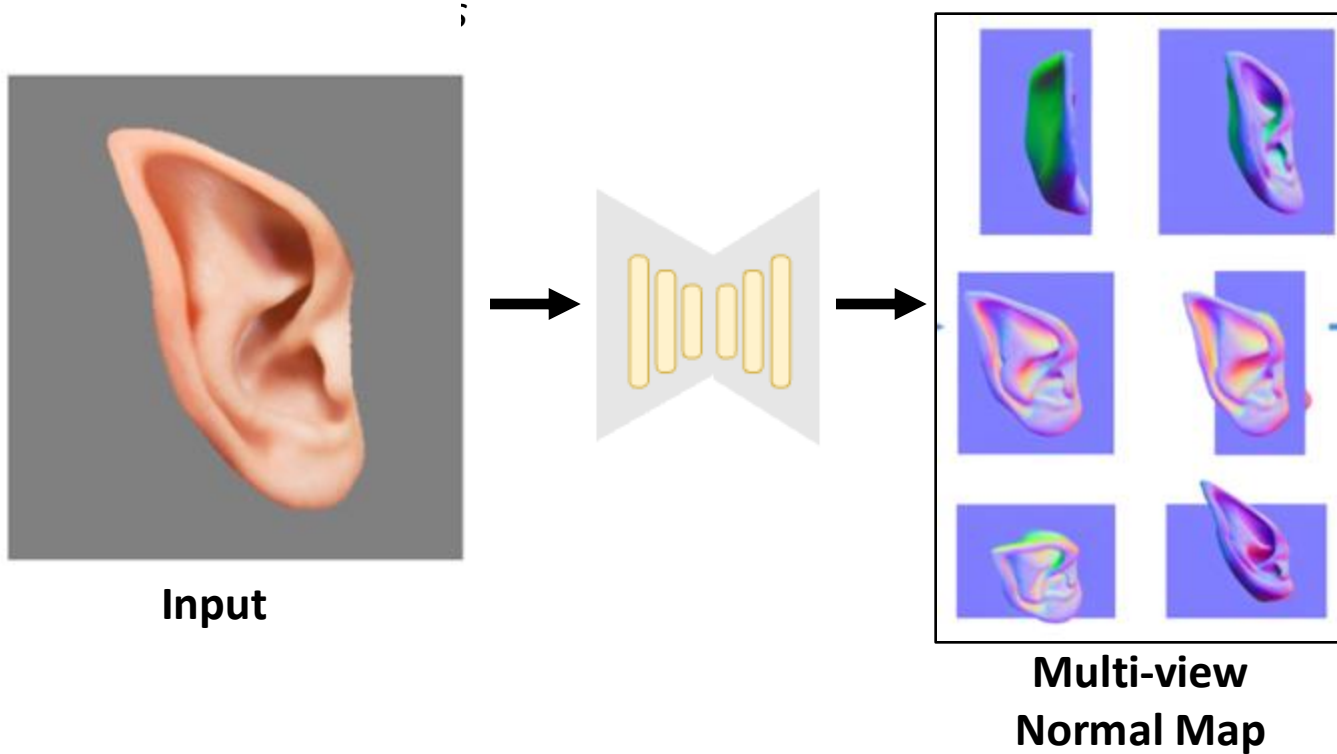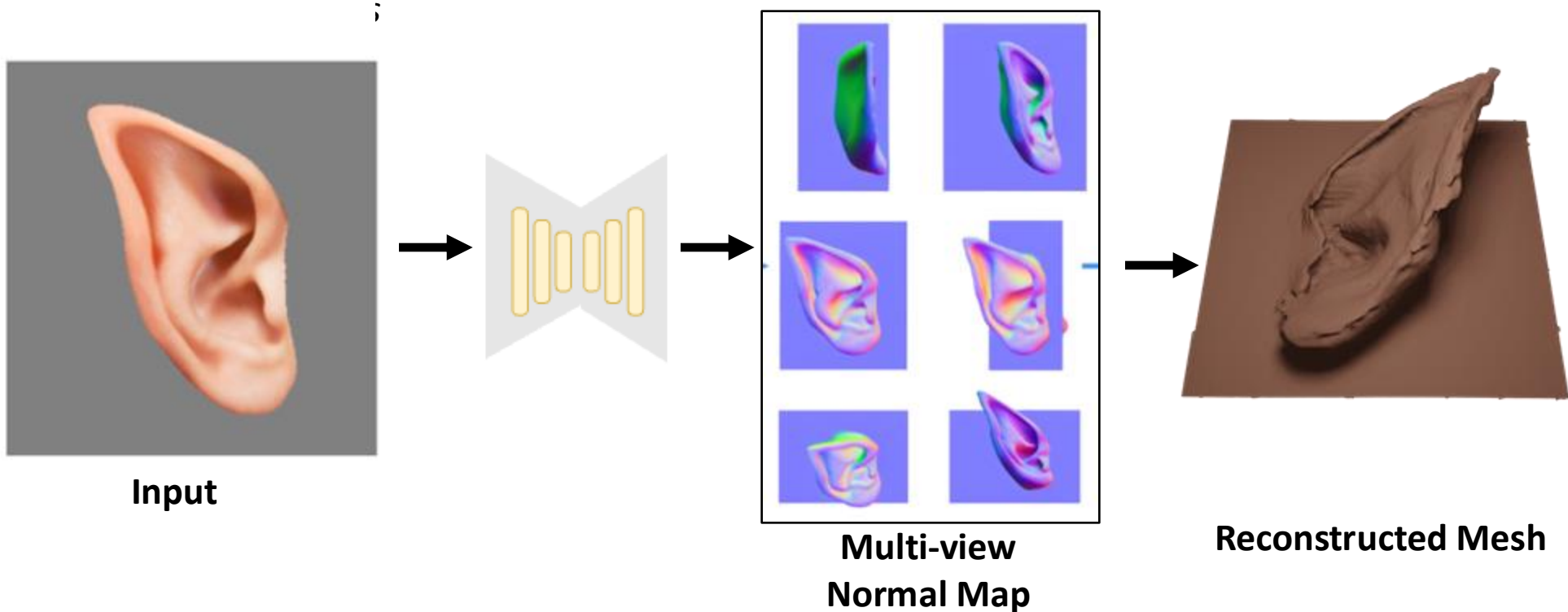
# Generative VDMs

- Create VDM from an image



Base Mesh

# Multi-view VDM Generation

- Start from an image



**Input**

**Multi-view Normal Map**

# Multi-view VDM Generation

- Start from an image, reconstruct



**Input**

**Multi-view Normal Map**

**Reconstructed Mesh**

# Multi-view VDM Generation

- Start from an image, reconstruct, optimize to get VDM



**Input**

**Multi-view Normal Map**

**Reconstructed Mesh**

**VDM Image**

# Multi-view VDM Generation

▪ Start from an image, reconstruct, optimize to get VDM



**Multi-View Stable Diffusion**

**Input**

**Multi-view Normal Map**

**Reconstructed Mesh**

**VDM Image**
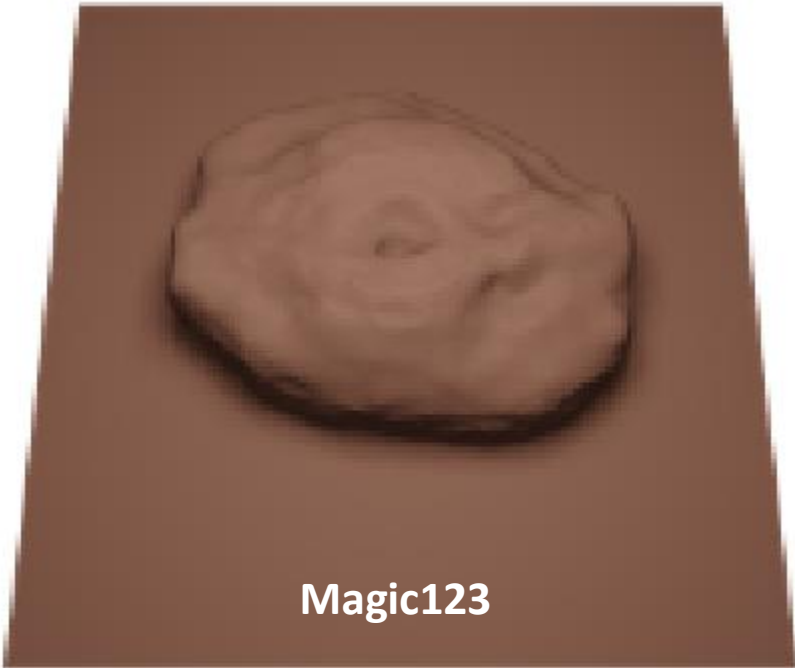
# Training Data

- Can we take a model pre-trained on full objects?



**Input**

# Training Data

- Can we take a model pre-trained on full objects? **NO!**



**Input**



**Magic123**

**LRM**

**Wonder3D**

# Training Data

- Can we take a model pre-trained on full objects? **NO!**



Input

Ours

Magic123

LRM

Wonder3D

# Training Data
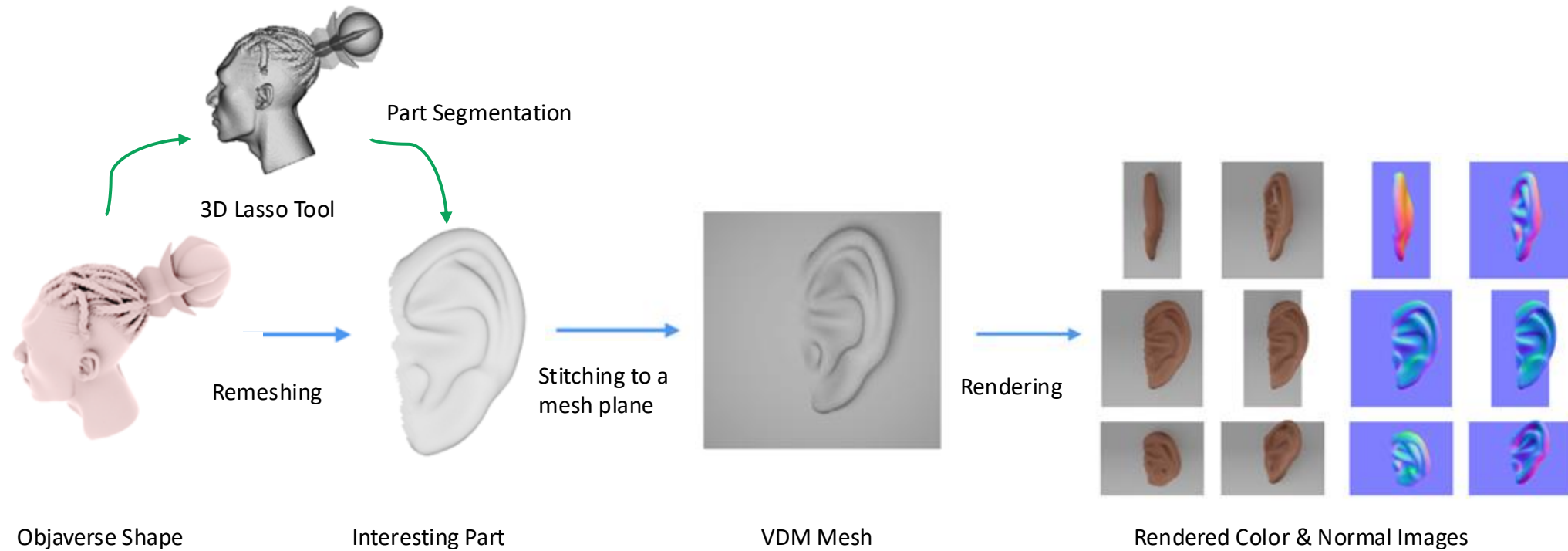
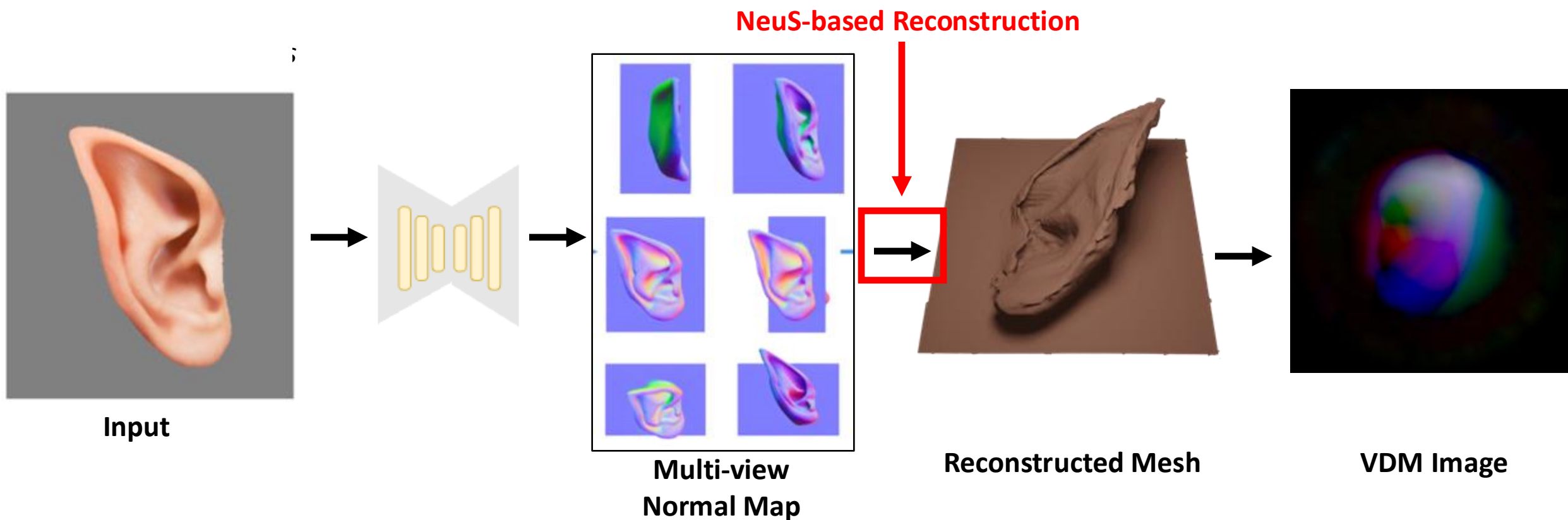- Can we take a model pre-trained on full objects? **NO!**



Input

Ours

Magic123

LRM

Wonder3D

# Training Data

- Can we take a model pre-trained on full objects? **NO!**

- VDM data pipeline using Objaverse



Objaverse Shape      Interesting Part      VDM Mesh      Rendered Color & Normal Images
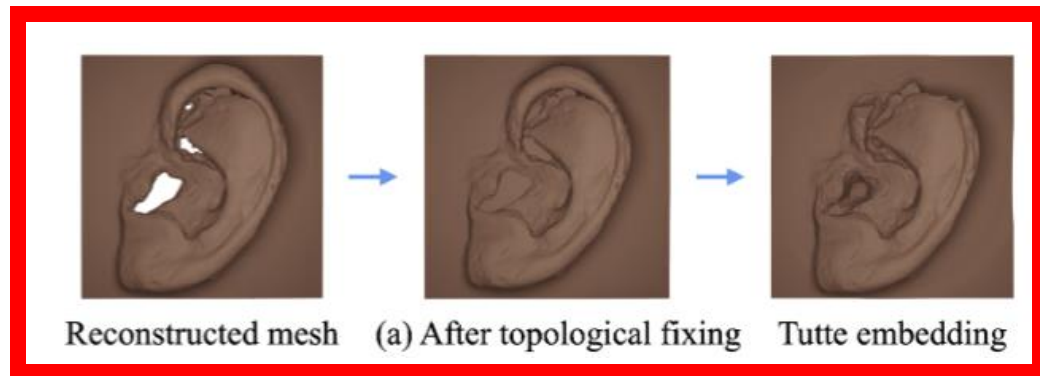
**We create a dataset of 1200 examples!**

# Multi-view VDM Generation

▪ Start from an image, reconstruct, optimize to get VDM



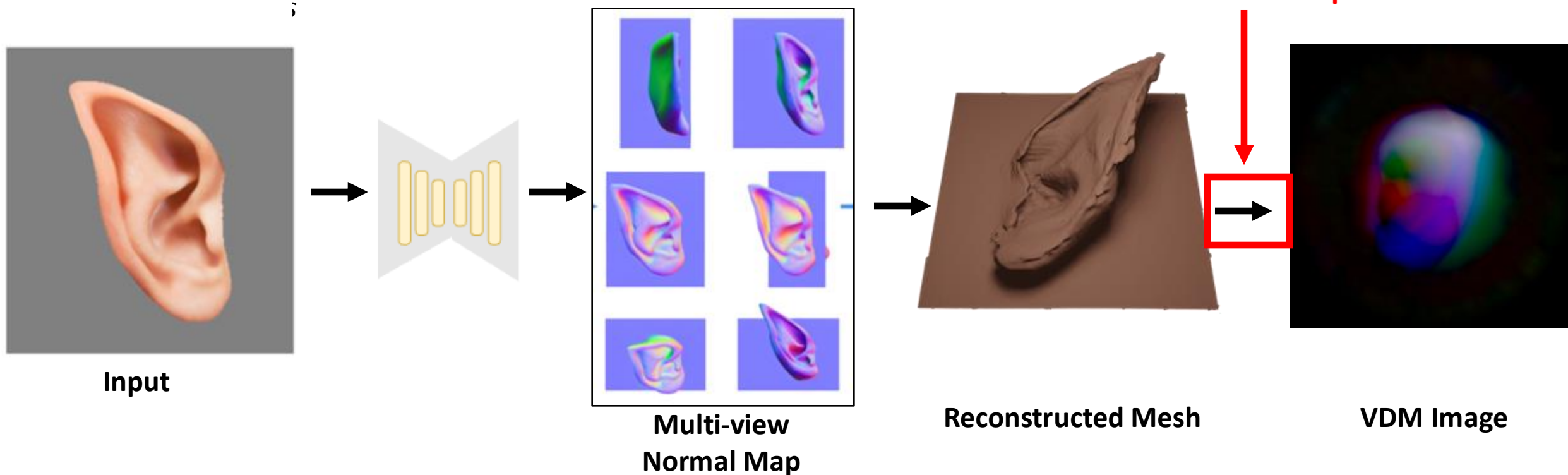**NeuS-based Reconstruction**

**Input**

**Multi-view Normal Map**

**Reconstructed Mesh**

**VDM Image**

**Adobe**

# Multi-view VDM Generation

▪ Start from an image, reconstruct, optimize to get VDM



Reconstructed mesh    (a) After topological fixing    Tutte embedding

**Naïve Baseline:**
**1. fill holes**
**2. Tutte embedding**

**Input**

**Multi-view
Normal Map**

**Reconstructed Mesh**

**VDM Image**

# Multi-view VDM Generation

- Start from an image, reconstruct, optimize to get VDM

$$f_\theta : [0,1]^2 \to \mathbb{R}^3$$

neural map

[0,1] x [0,1]

**Optimize VDM as a Neural Map**



**Input**

**Multi-view Normal Map**

**Reconstructed Mesh**

**VDM Image**
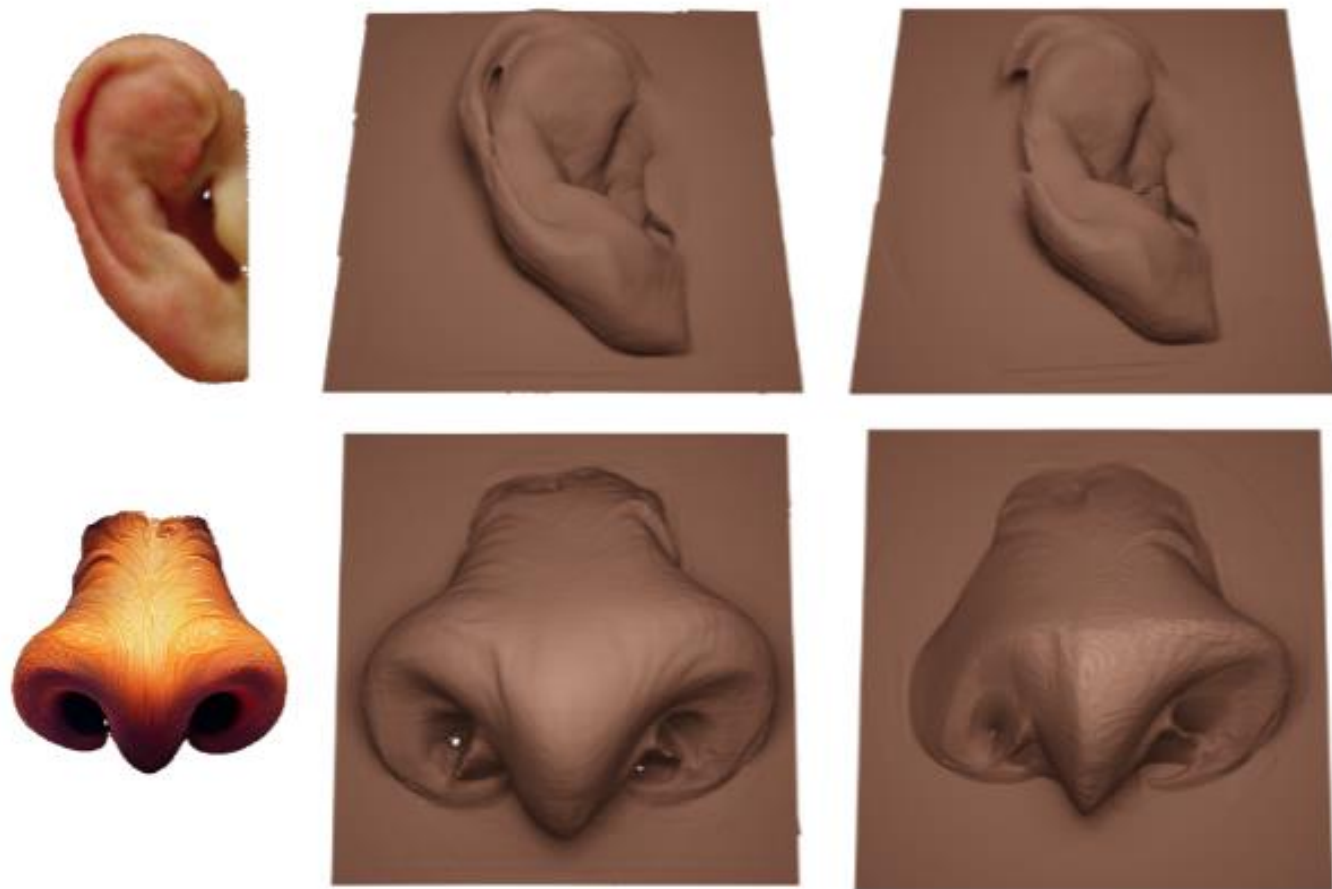
GenVDM, Yang et al., In Submission, 2025

# Ablation: Reconstruction Alternatives



Input images          Reconstructed meshes
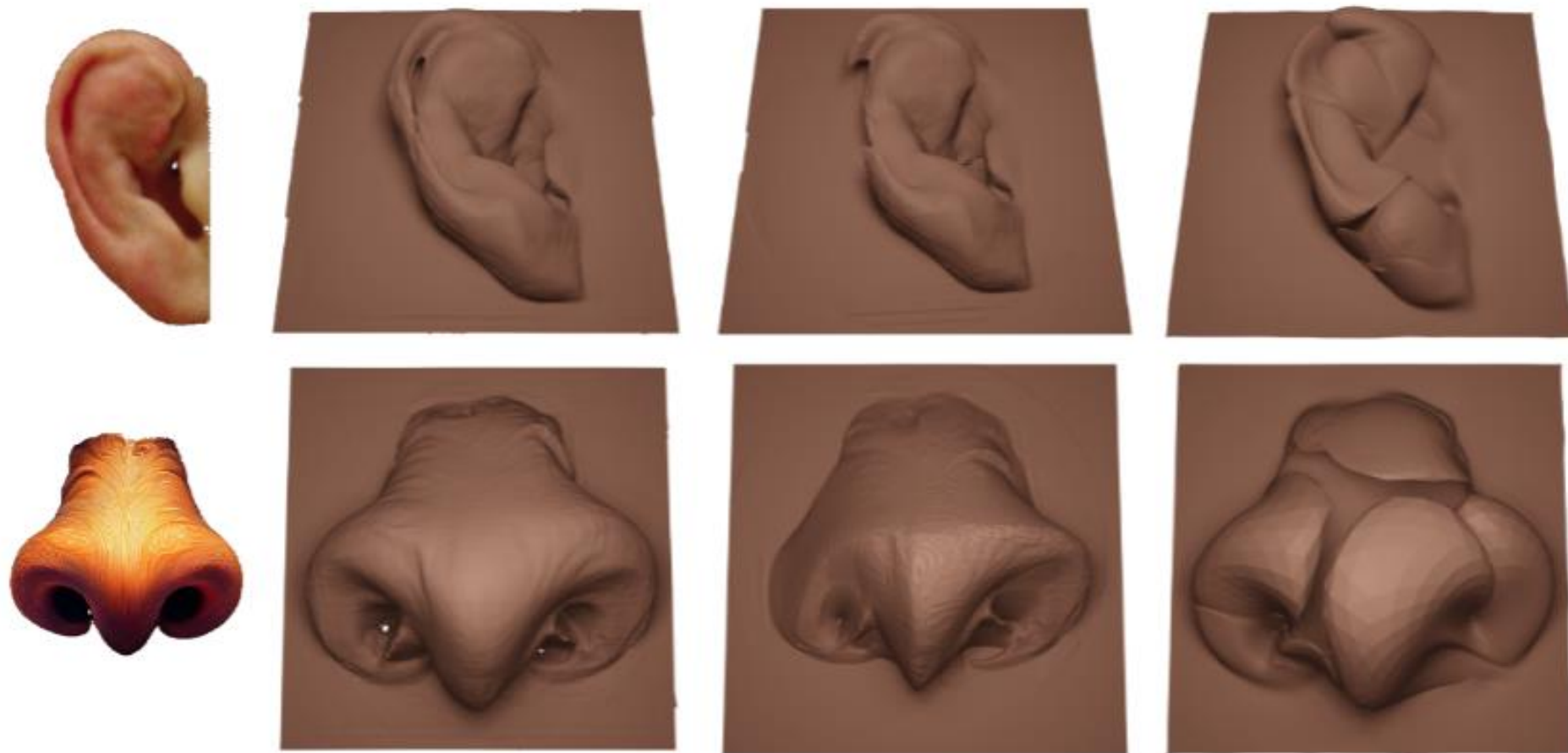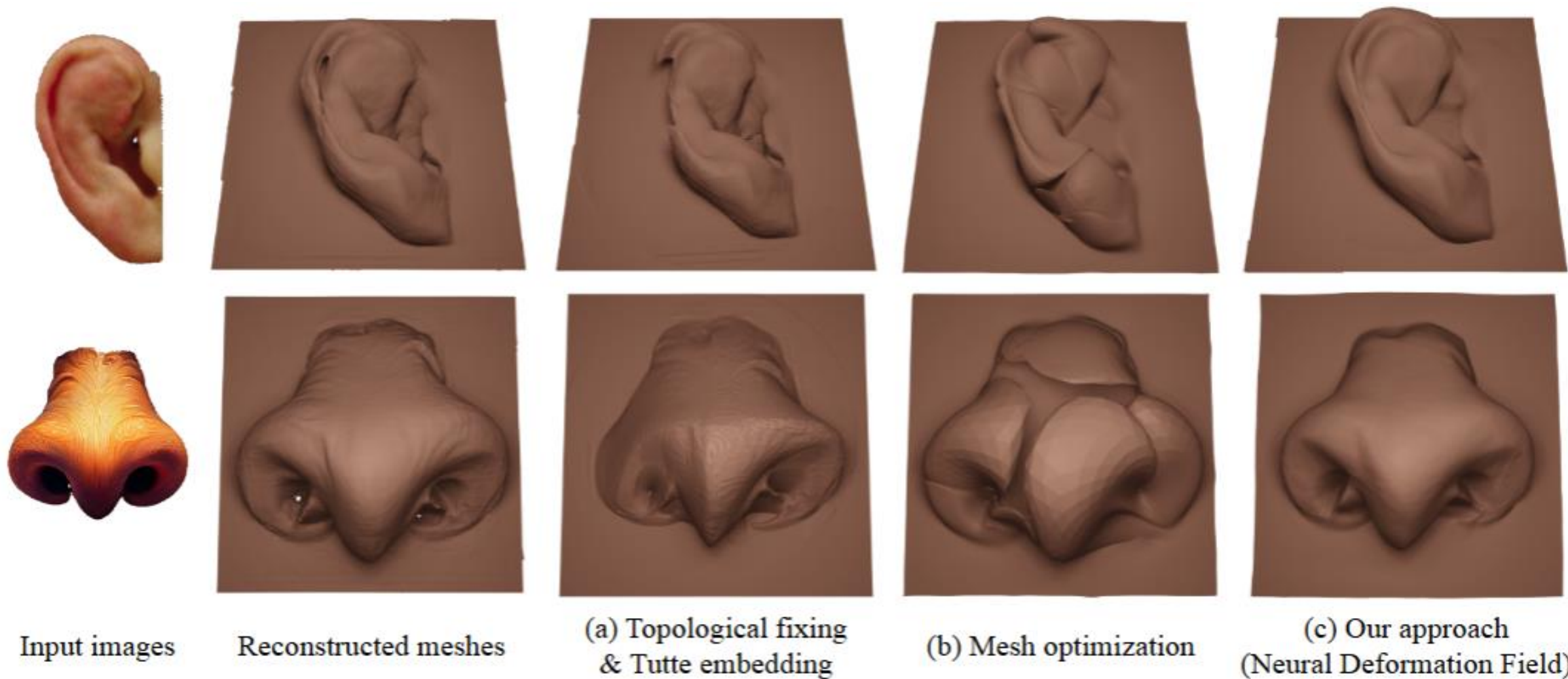
# Ablation: Reconstruction Alternatives



Input images    Reconstructed meshes    (a) Topological fixing
                                            & Tutte embedding

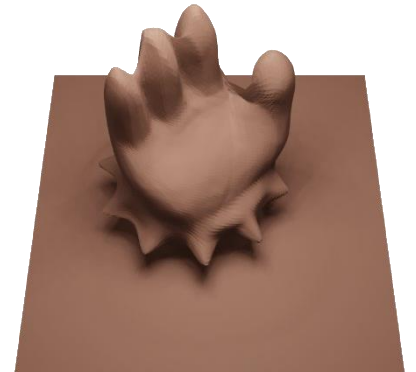# Ablation: Reconstruction Alternatives
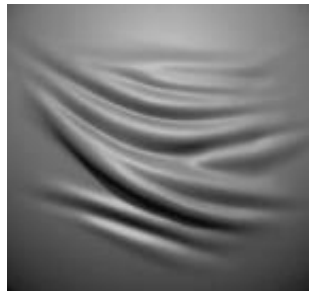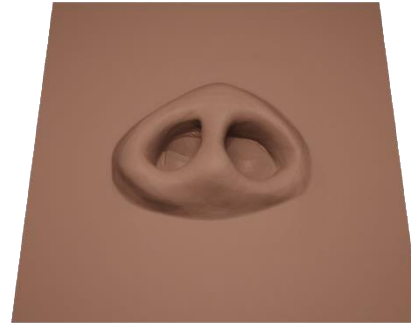


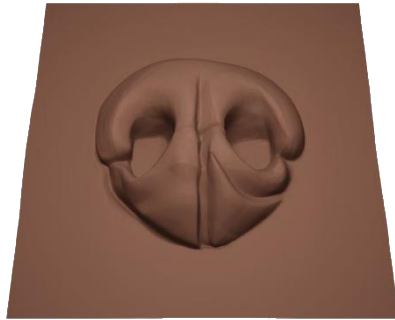Input images     Reconstructed meshes     (a) Topological fixing & Tutte embedding     (b) Mesh optimization

# Ablation: Reconstruction Alternatives



Input images | Reconstructed meshes | (a) Topological fixing & Tutte embedding | (b) Mesh optimization | (c) Our approach (Neural Deformation Field)

# Generated VDMs

# Generated VDMs

# Image Editing for VDM



Original                                    Edition 1                                    Edition 2
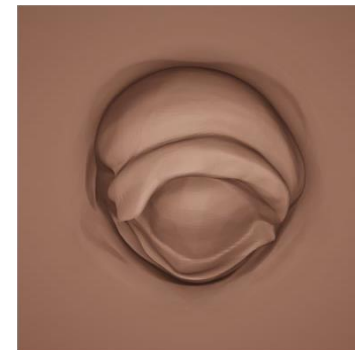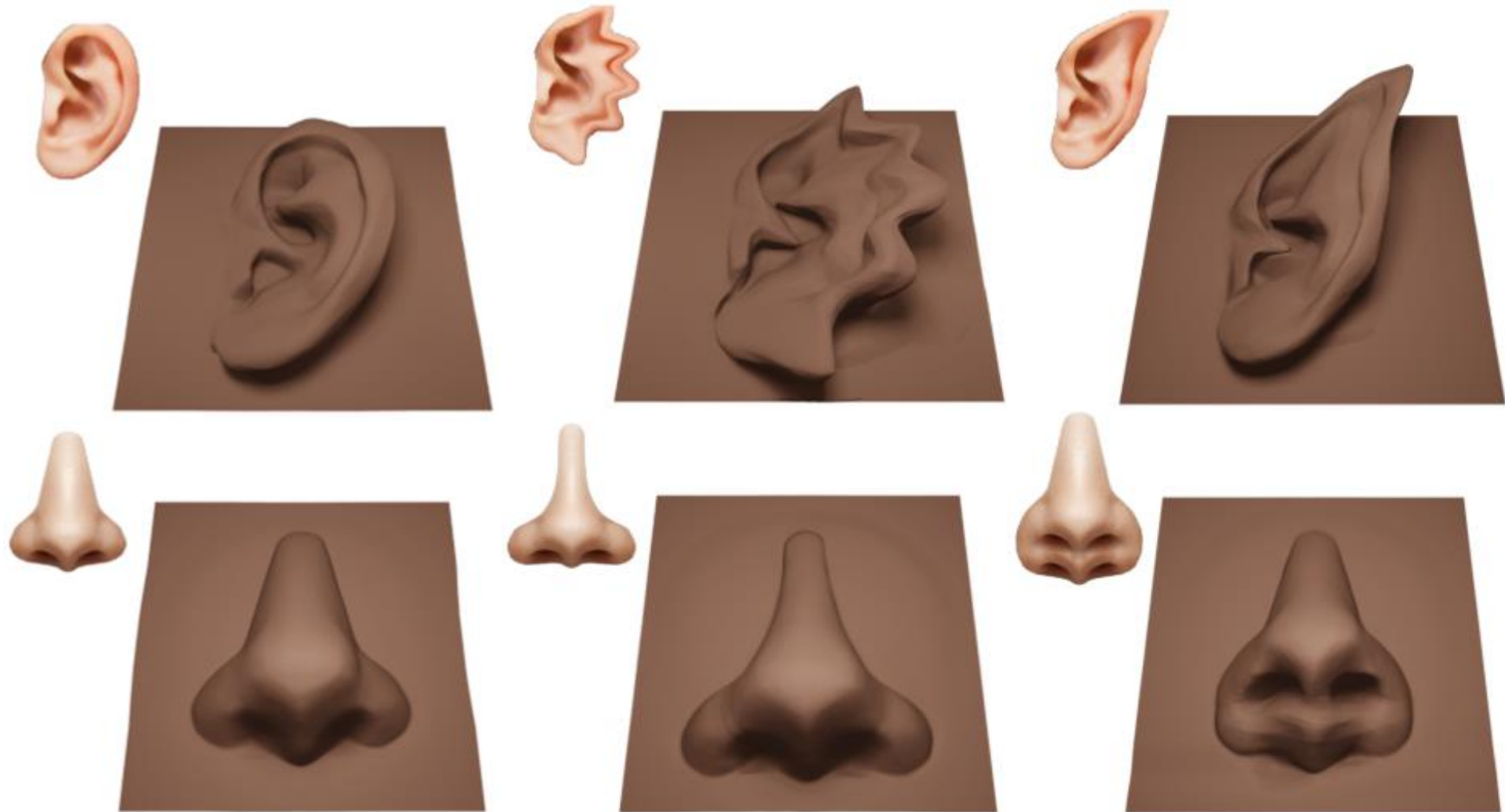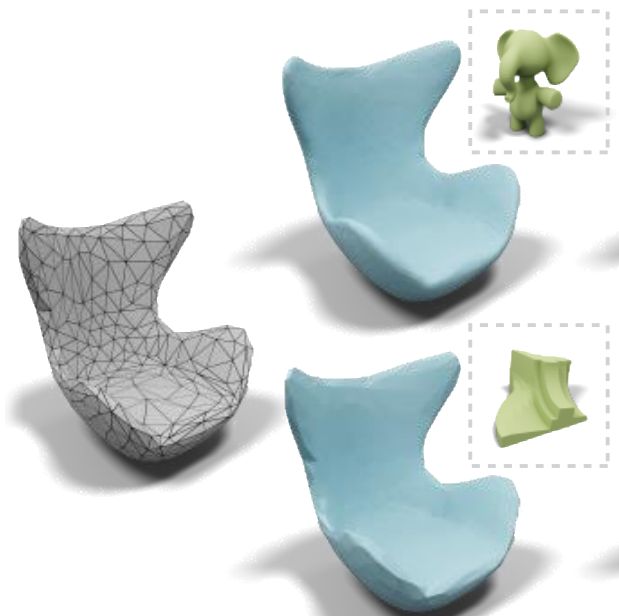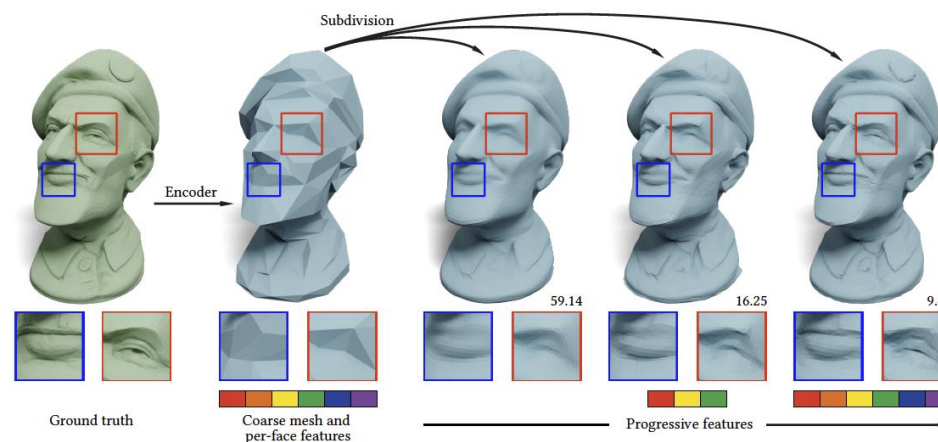
# Other Work on Detailization



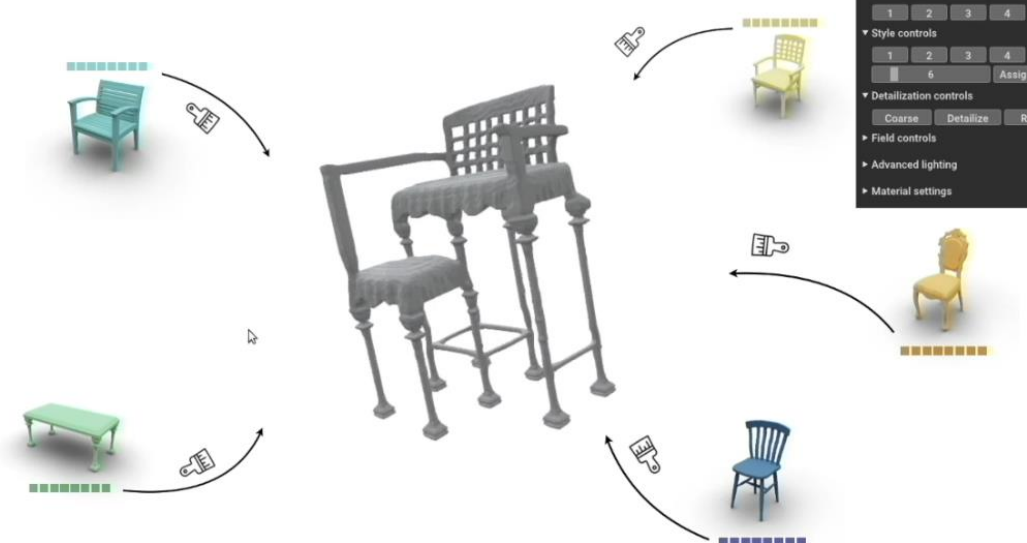"a classic furniture piece made of polished wood with subtle details"

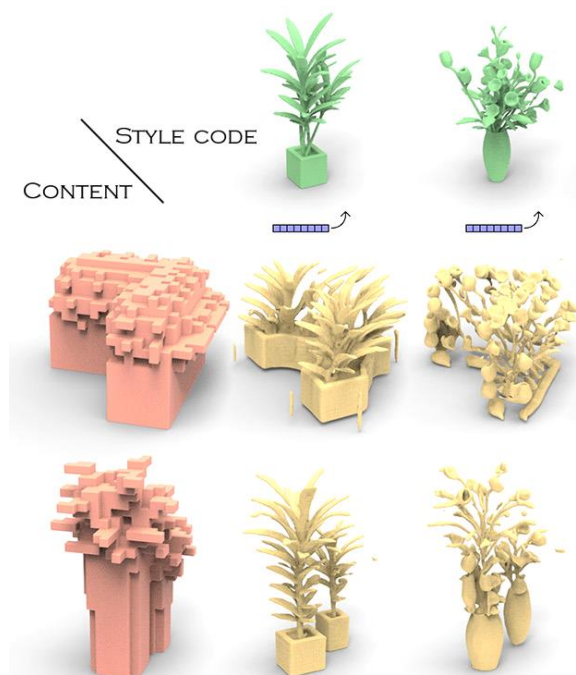ART-DECO (Qimin Chen et al.), under review, 2025

Neural Progressive Meshes (Hsueh-Ti Liu et al.), SIGGRAPH 2020

Neural Progressive Meshes (Yun-Chun Chen et al.), SIGGRAPH 2023

DECOLLAGE (Qimin Chen et al.), ECCV 2024
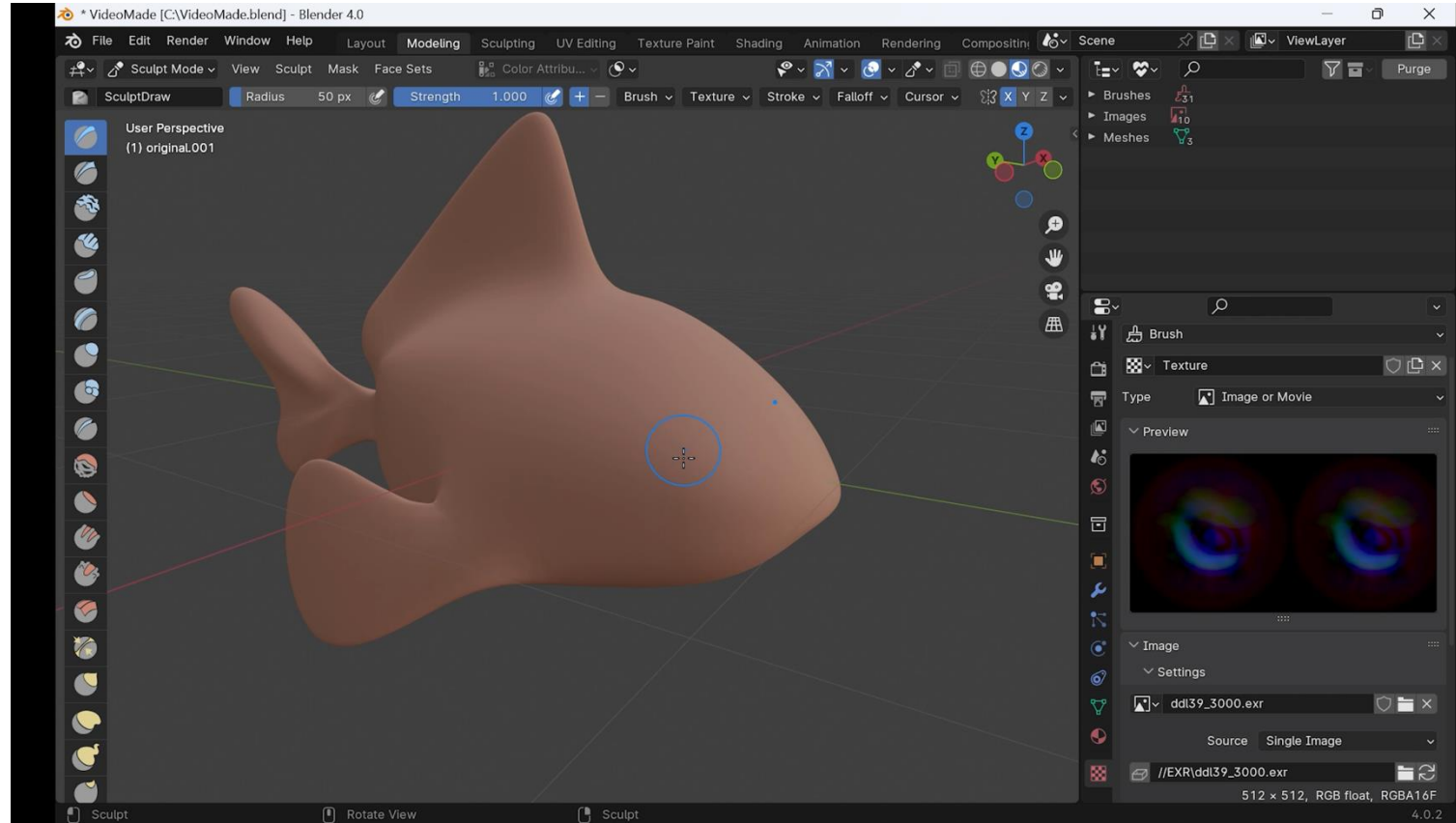
DecorGAN (Zhiqin Chen et al.), CVPR 2021

# Summary

- Workflows

- Representations
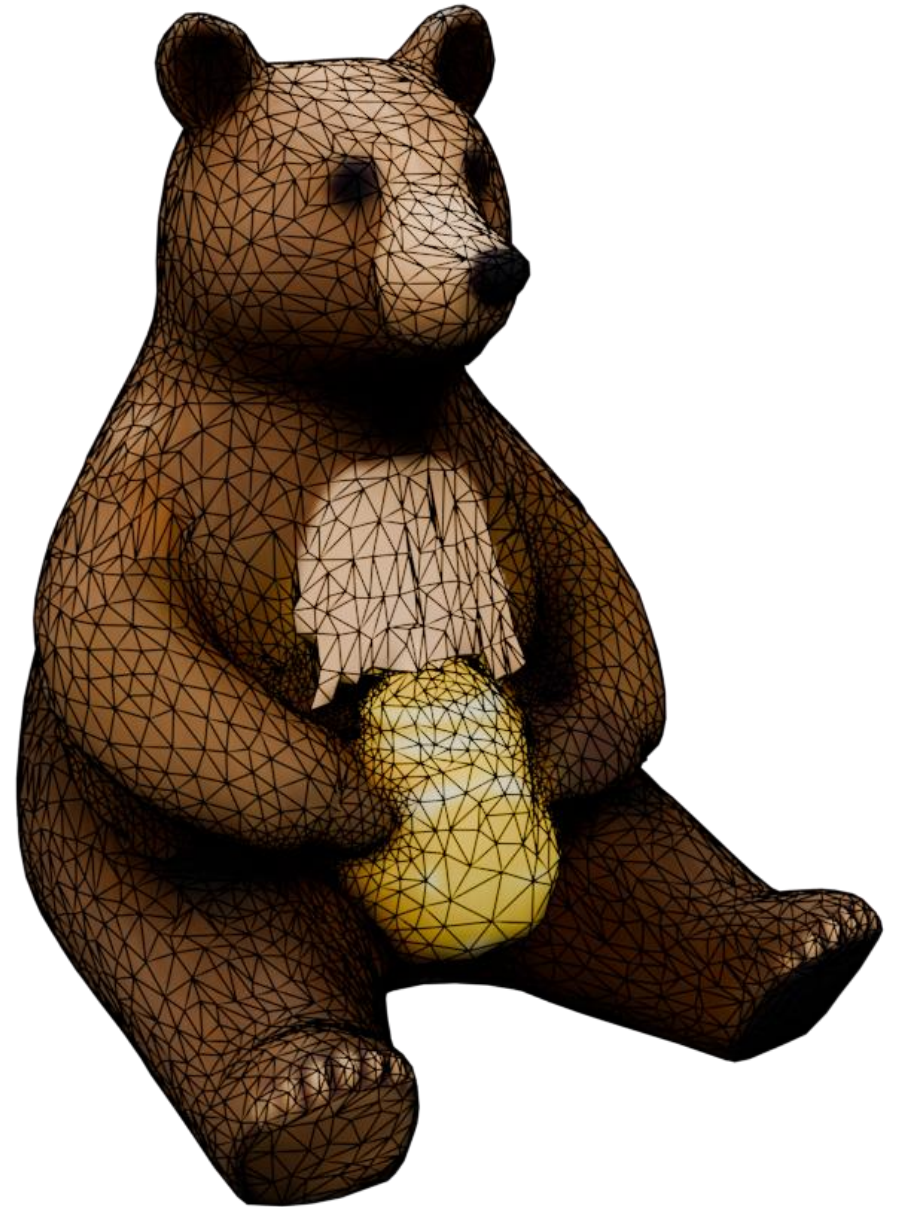
- Multi-view predictions

# Summary

- **Workflows**

  - Coarse deformation

  - Coarse sculpting

  - Detailize via normal maps

  - Detailize via VDMs

  - ...

- Representations

- Multi-view predictions

# Summary

- Workflows

- **Representations**

  - Meshes are essential for many existing workflows and pipelines

  - Hybrid representations allows to get the best of both worlds

  - Representation-agnostic methods via multi-view + LRM / Optimization

- Multi-view predictions

# Summary

- Workflows

- Representations

- Multi-view predictions
  - Leverage pre-trained 2D priors
  - Work with different LRMs
  - Often need to be fine-tuned



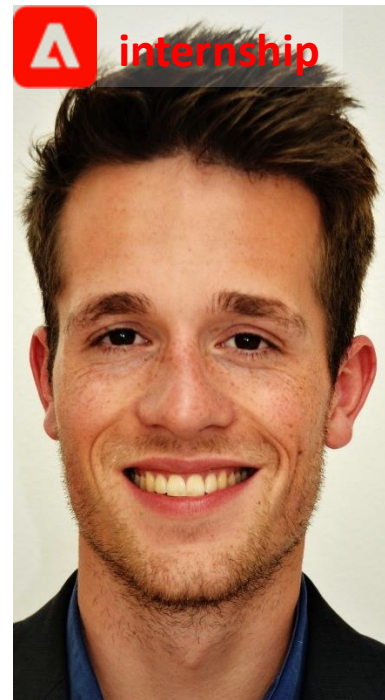**Adobe**

# Collaborators



**APAP**
Seungwoo Yoo
KAIST

**MeshUp**
Brian Kim

1. **MagicClay**
2. **Instant3dit**
Amir Barda

**Text-guided refinement**
Yun-Chun Chen

**SAMa**
Michael Fischer
UCL

**VDM**
Yuezhi Yang
THE UNIVERSITY OF TEXAS AT AUSTIN

**Adobe:** Matheus Gadelha, Thibault Groueix, Zhiqin Chen, Siddhartha Chaudhuri, Iliyan Giorgiev, Valentin Deschaintre, Alec Jacobson

**Academia:** Noam Aigerman, Amir Barda, Rana Hanocka, Qixing Huang, Kunho Kim, Itai Lang, Minhyuk Sung, Hao Zhang