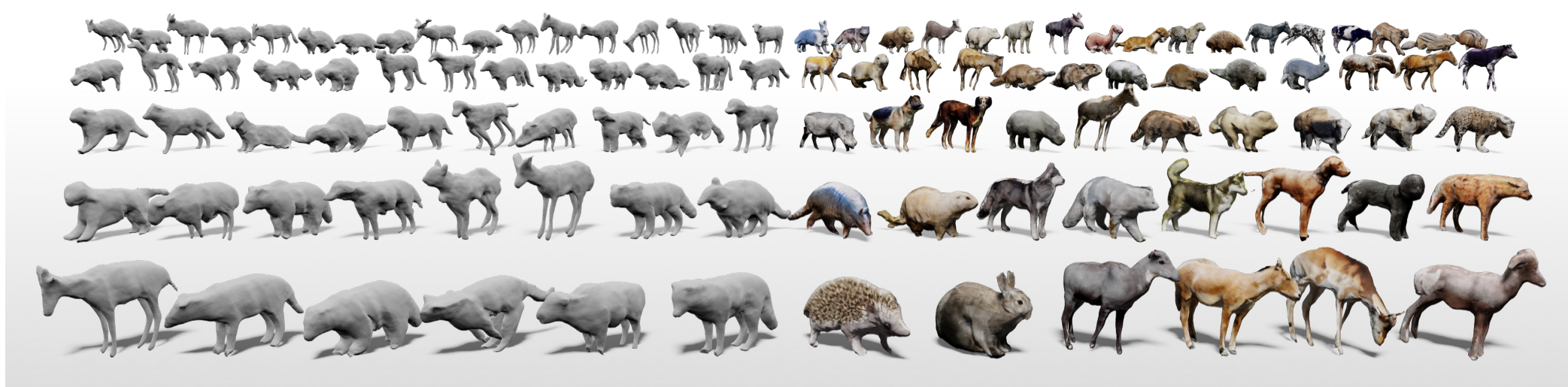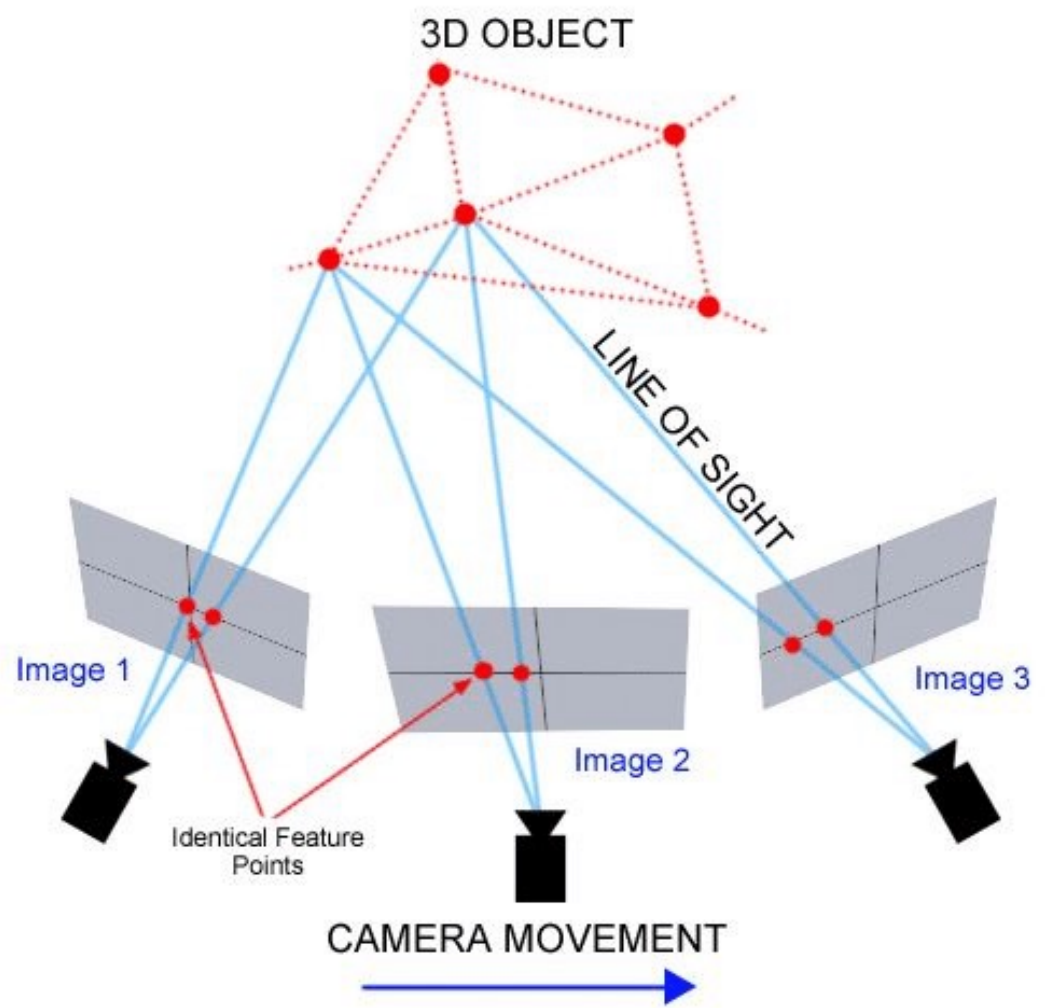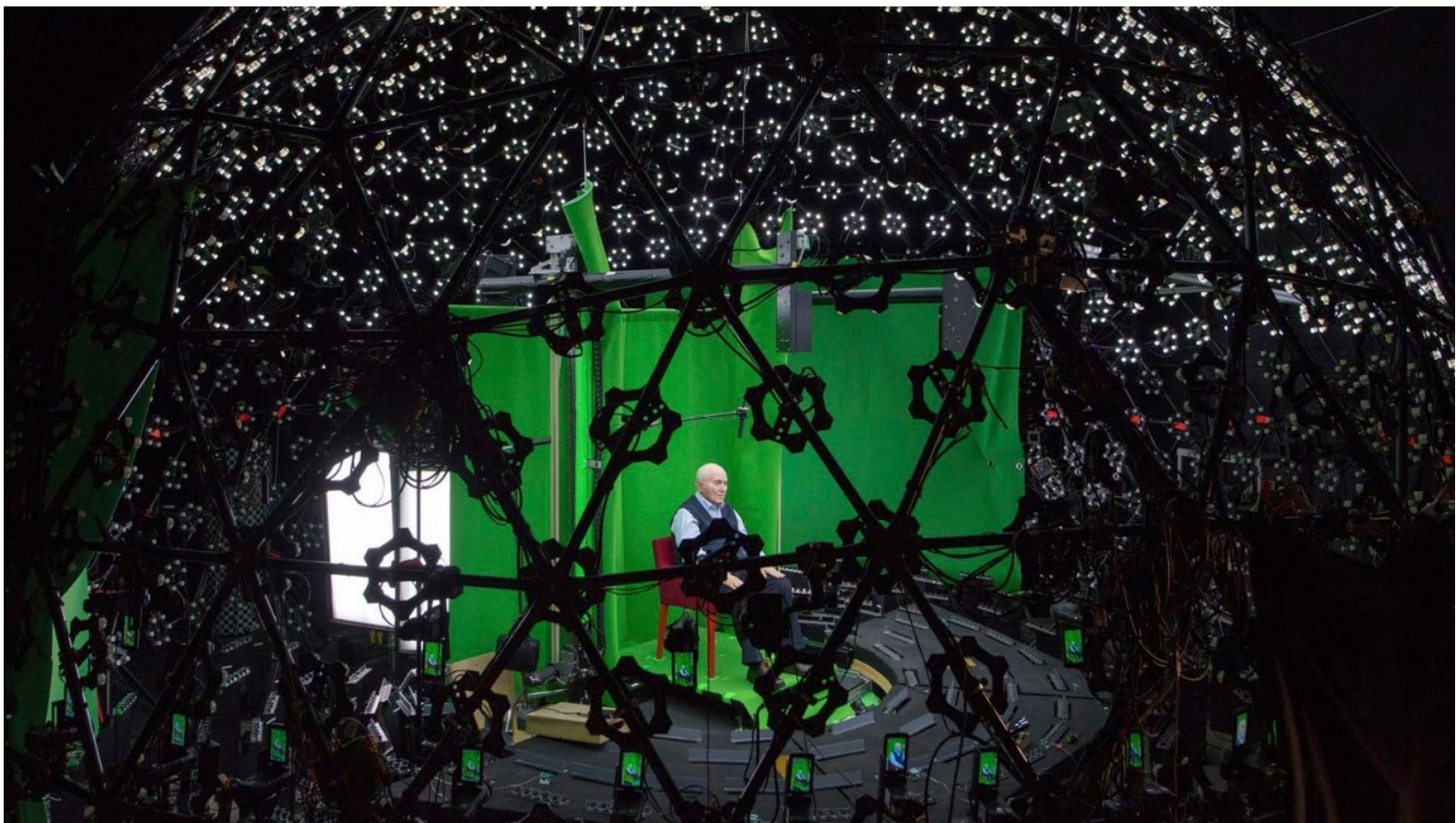# Learning Articulated 3D Animals from Internet Images

Tomas Jakab, University of Oxford, VGG

3D OBJECT

LINE OF SIGHT

Image 1

Image 3

Image 2

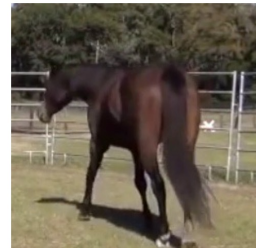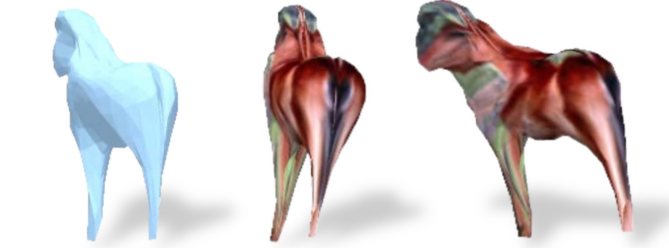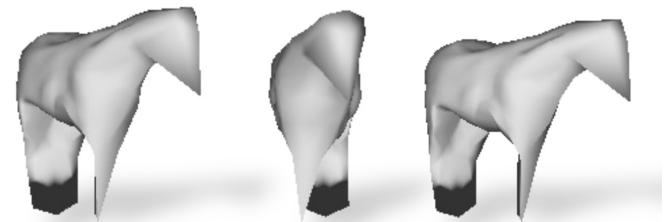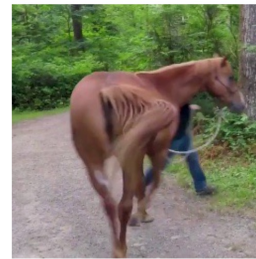Identical Feature Points

CAMERA MOVEMENT

# Training Data – Single View Images

# Prior work



| Input | Input view | Other views | Input | Input view | Other views |

UMR [1]

DOVE [2]

[1] Self-supervised single-view 3d reconstruction via semantic consistency.
Li et. al. ECCV 2020.

[2] DOVE: Learning deformable 3d objects by watching videos.
Wu et. al. IJCV, 2023.

# Training Data – Single View Images

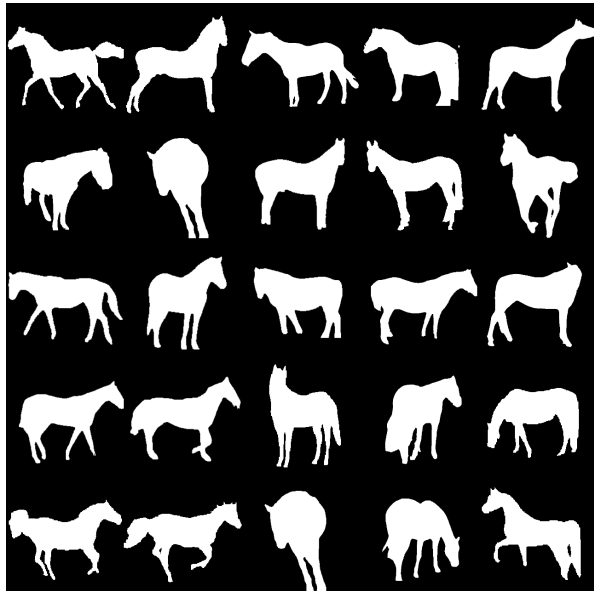# Training Data – Single View Images



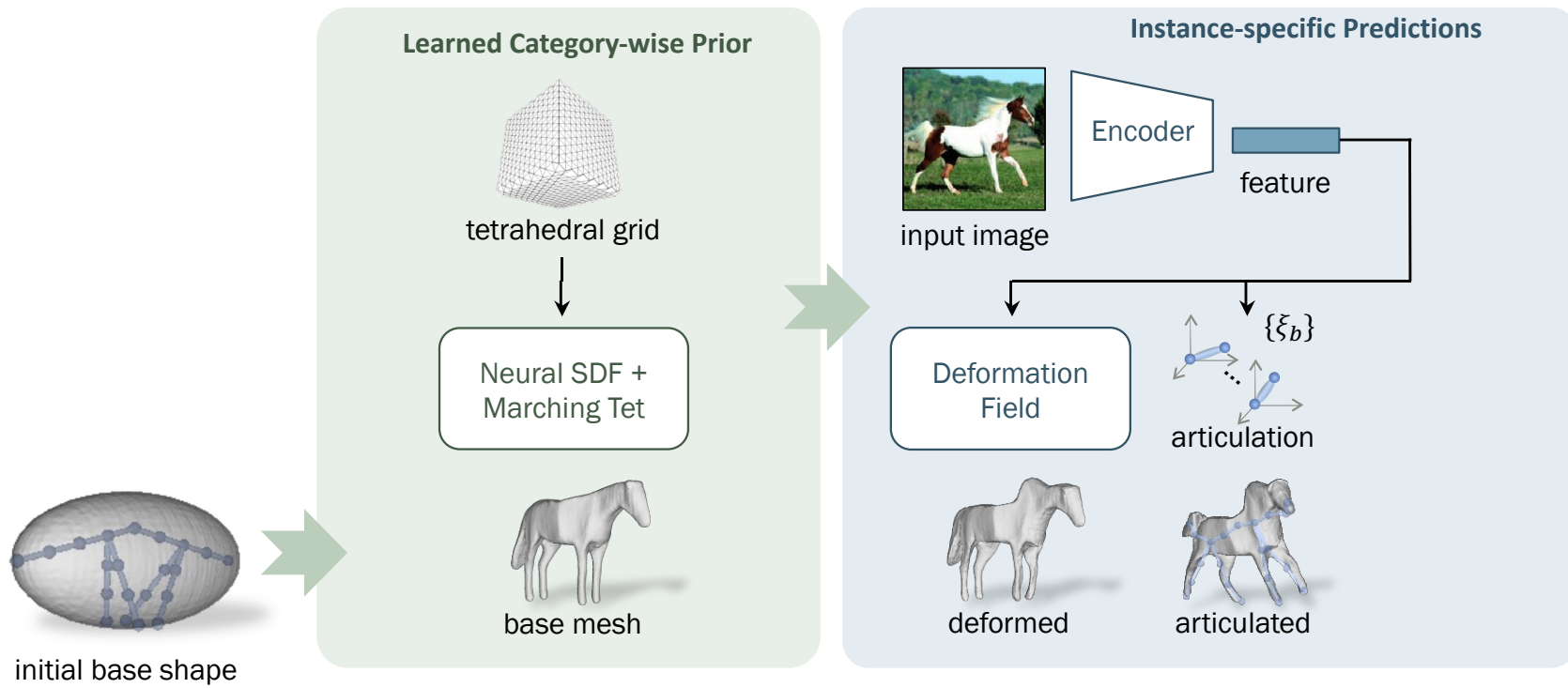## What is different?

Background

Shape

Appearance

Instance Masks

# Modeling shape
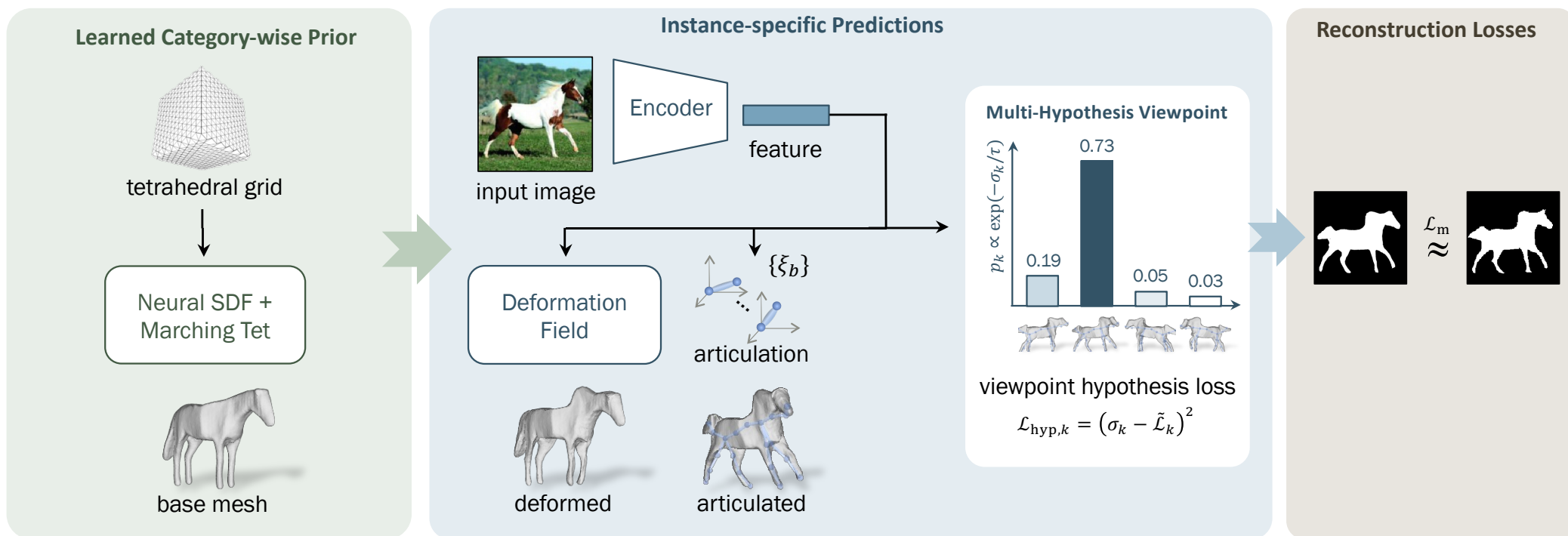
# Hierarchical Shape Prediction

# Camera Prediction



**Learned Category-wise Prior**

tetrahedral grid

Neural SDF + Marching Tet

base mesh

**Instance-specific Predictions**

input image

Encoder

feature

Deformation Field

$\{\xi_b\}$

articulation

deformed

articulated

**Multi-Hypothesis Viewpoint**

$p_k \propto \exp(-\sigma_k/\tau)$

0.19  0.73  0.05  0.03

viewpoint hypothesis loss

$\mathcal{L}_{\mathrm{hyp},k} = (\sigma_k - \tilde{\mathcal{L}}_k)^2$

**Reconstruction Losses**

$\mathcal{L}_{\mathrm{m}}$
$\approx$

# Category Appearance



Off-the-shelf
DINO-ViT [1]

Self-supervised Image Features

[1] Emerging Properties in Self-supervised Vision Transformers. Caron et. al. ICCV 2021.

# Category Appearance



**Learned Category-wise Prior**

tetrahedral grid

Feature Field

Neural SDF + Marching Tet

DINO feature

base mesh

**Instance-specific Predictions**

input image

Encoder

feature

Deformation Field

$\{\xi_b\}$

articulation

deformed

articulated

**Multi-Hypothesis Viewpoint**

$p_k \propto \exp(-\sigma_k/\tau)$

0.19   0.73   0.05   0.03

viewpoint hypothesis loss

$\mathcal{L}_{\mathrm{hyp},k} = \left(\sigma_k - \tilde{\mathcal{L}}_k\right)^2$

**Reconstruction Losses**

$\mathcal{L}_{\mathrm{m}}$ $\approx$

$\mathcal{L}_{\mathrm{feat}}$ $\approx$

# Instance Appearance

# Canonical Appearance



**Learned Category-wise Prior**

tetrahedral grid

Feature Field → DINO feature

Neural SDF + Marching Tet → base mesh

**Instance-specific Predictions**

input image → Encoder → feature

Deformation Field → deformed

articulation $\{\xi_b\}$ → articulated

Albedo Field → albedo

light → shading

**Multi-Hypothesis Viewpoint**

$p_k \propto \exp(-\sigma_k/\tau)$

0.19    0.73    0.05    0.03

viewpoint hypothesis loss

$$\mathcal{L}_{\mathrm{hyp},k} = (\sigma_k - \tilde{\mathcal{L}}_k)^2$$

**Reconstruction Losses**

$\mathcal{L}_{\mathrm{m}} \approx$

$\mathcal{L}_{\mathrm{feat}} \approx$

$\mathcal{L}_{\mathrm{im}} \approx$

$$\mathcal{L} = \mathbb{E}_{p_k}[\mathcal{L}_{\mathrm{feat}} + \mathcal{L}_{\mathrm{im}} + \mathcal{L}_{\mathrm{m}} + \mathcal{R}_{\mathrm{Eik}} + \mathcal{R}_{\mathrm{def}} + \mathcal{R}_{\mathrm{art}}] + \mathcal{L}_{\mathrm{hyp}}$$

## Entire pipeline trained end-to-end with reconstruction losses

(except for frozen DINO-ViT image encoder, pre-trained via self-supervision)

no keypoints, no template shapes

MagicPony: Learning Articulated 3D Animals in the Wild. Shangzhe Wu*, Ruining Li*, Tomas Jakab*, Christian Rupprecht, Andrea Vedaldi. CVPR 2023

# Results

Off-the-shelf
DINO-ViT [1]

Self-supervised Image Features

[1] Emerging Properties in Self-supervised Vision Transformers. Caron et. al. ICCV 2021.

# Instance Appearance



**Learned Category-wise Prior**

tetrahedral grid

Feature Field

Neural SDF + Marching Tet

DINO feature    base mesh

**Instance-specific Predictions**

input image    DINO    Encoder    feature

Deformation Field    $\{\xi_b\}$ articulation    Albedo Field    light
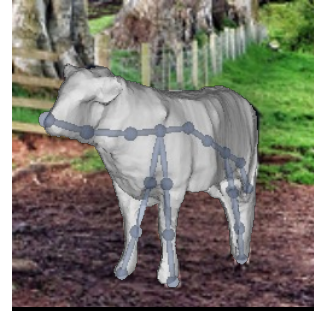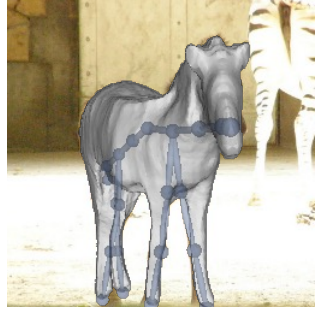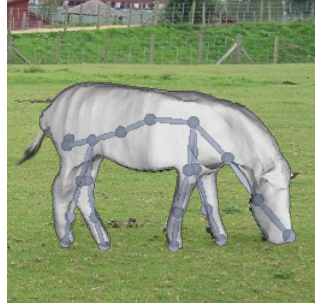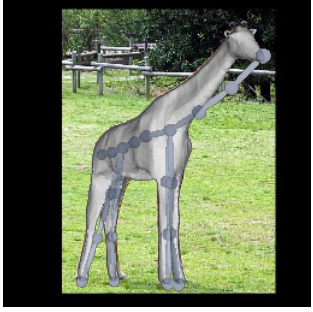
deformed    articulated    albedo    shading

**Multi-Hypothesis Viewpoint**

$p_k \propto \exp(-\sigma_k/\tau)$

0.19    0.73    0.05    0.03

viewpoint hypothesis loss

$$\mathcal{L}_{\mathrm{hyp},k} = \left(\sigma_k - \tilde{\mathcal{L}}_k\right)^2$$

**Reconstruction Losses**

$\mathcal{L}_{\mathrm{m}} \approx$

$\mathcal{L}_{\mathrm{feat}} \approx$

$\mathcal{L}_{\mathrm{im}} \approx$

# Frame-by-Frame Inference on Videos



Input Frames          Input View          360° Rotations

# Follow up works

# Real vs Diffusion Generated Images

Typical Unsuitable Real Images
from ImageNet



StableDiffusion Generated Images



*"Implicitly curated"*

# Synthetic Training Images

Real Training Images



Synthetic Training Images

**Generating Training Images**

Stable Diffusion

Prompt: *"cow"*

↓

synthetic training images $\mathcal{D}$

# Virtual Multi-view Supervision



Generating Training Images

Stable Diffusion

Prompt: *"cow"* ↓

synthetic training images $\mathcal{D}$

**Virtual Multi-view Supervision**

SDS gradient

random $\tilde{l}$    random $\tilde{v}$

virtual view $\tilde{I}$

Stable Diffusion

[1] DreamFusion: Text-to-3D using 2D Diffusion. Poole et. al. arXiv 2022.

# Training pipeline



Farm3D: Learning Articulated 3D Animals by Distilling 2D Diffusion. Tomas Jakab*, Ruining Li*, Shangzhe Wu, Christian Rupprecht, Andrea Vedaldi. 3DV 2024

# Comparison with MagicPony



Input View

Novel View

Ours

MagicPony

Input Image

MagicPony: Learning Articulated 3D Animals in the Wild
*Shangzhe Wu\*, Ruining Li\*, Tomas Jakab\*, Christian Rupprecht, Andrea Vedaldi*, CVPR 2023, *\*equal contribution*

# Comparison with MagicPony



Input Image

Ours

MagicPony

MagicPony: Learning Articulated 3D Animals in the Wild
*Shangzhe Wu*, Ruining Li*, Tomas Jakab*, Christian Rupprecht, Andrea Vedaldi*, CVPR 2023, *equal contribution*
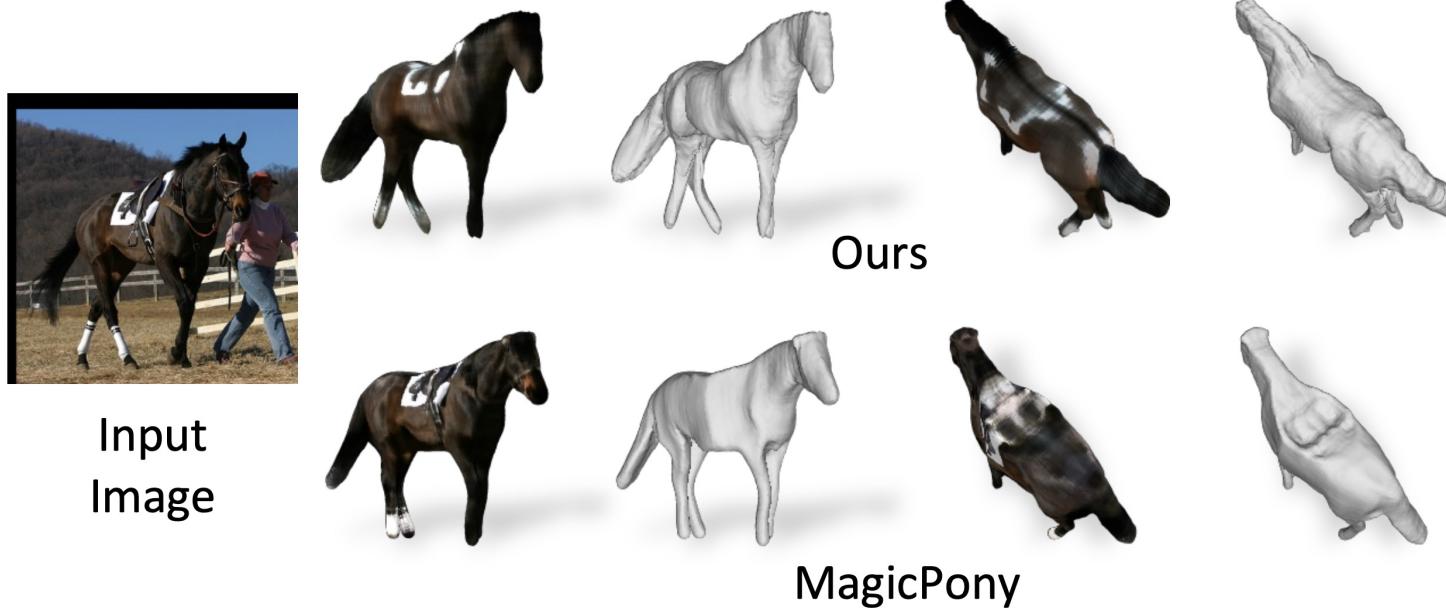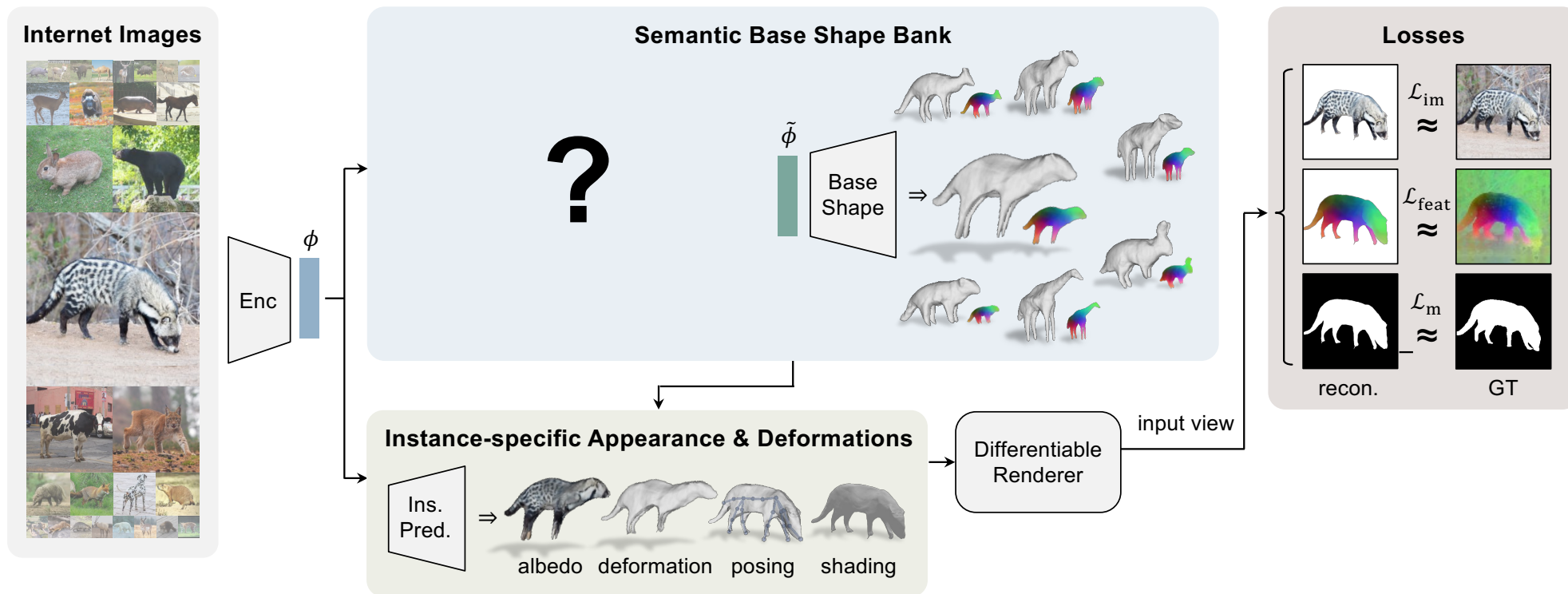
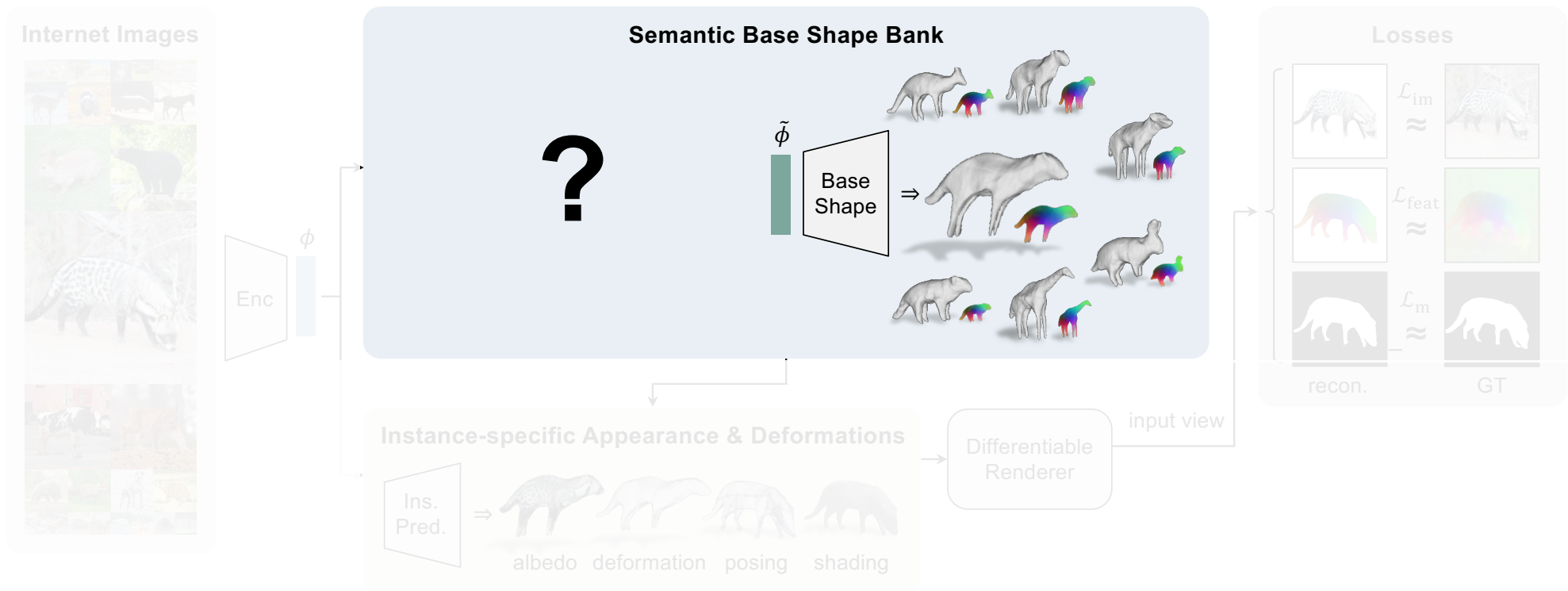# Towards Reconstructing the Animal Kingdom



Learning the 3D Fauna of the Web. Zizhang Li, Dor Litvak, Yunzhi Zhang, Ruining Li, Tomas Jakab, Christian Rupprecht, Shangzhe Wu, Andrea Vedaldi, Jiaiun Wu. arXiv:2401.02400
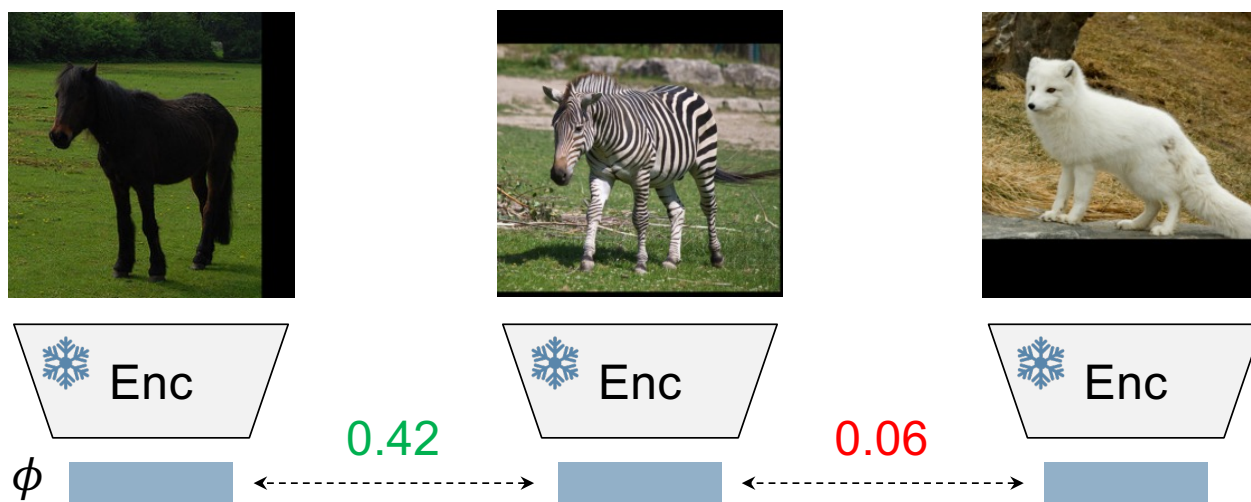
# Handling multiple categories
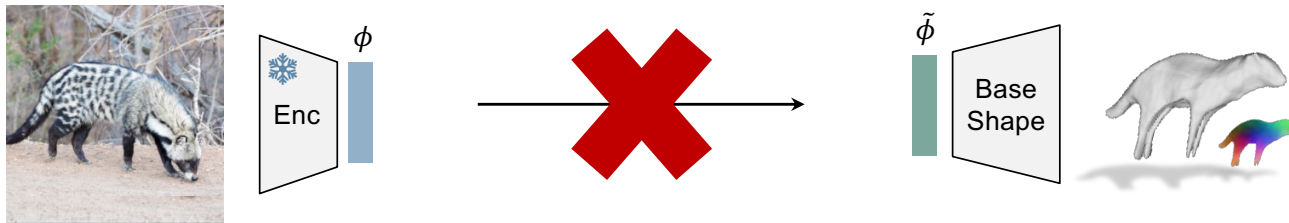
# Handling multiple categories

# Category embedding

- Leverage a pre-trained vision encoder - DINO

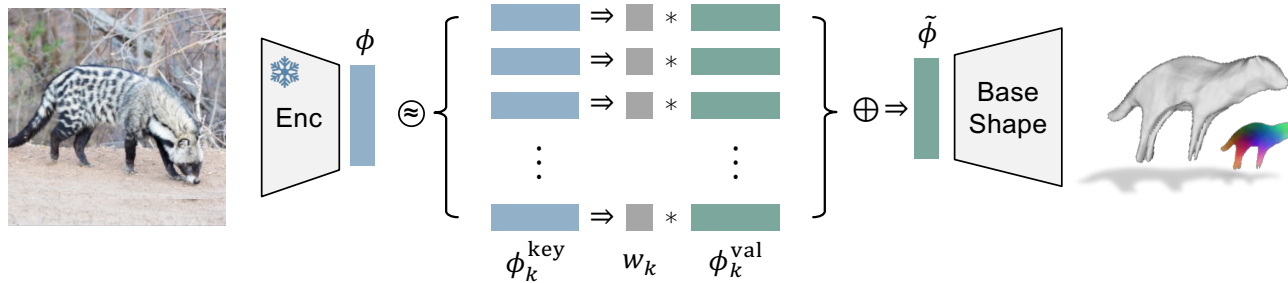  - Features serve as the soft definition of category

# Category embedding
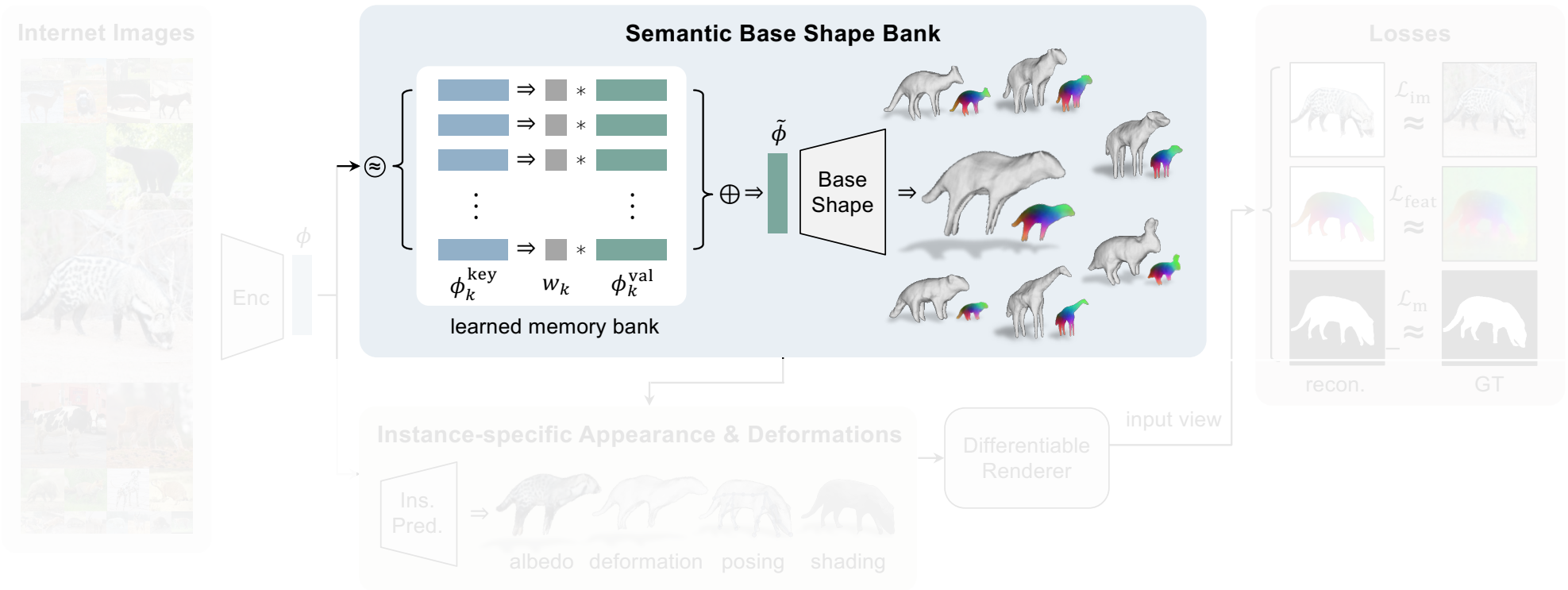
- Leverage a pre-trained vision encoder - DINO
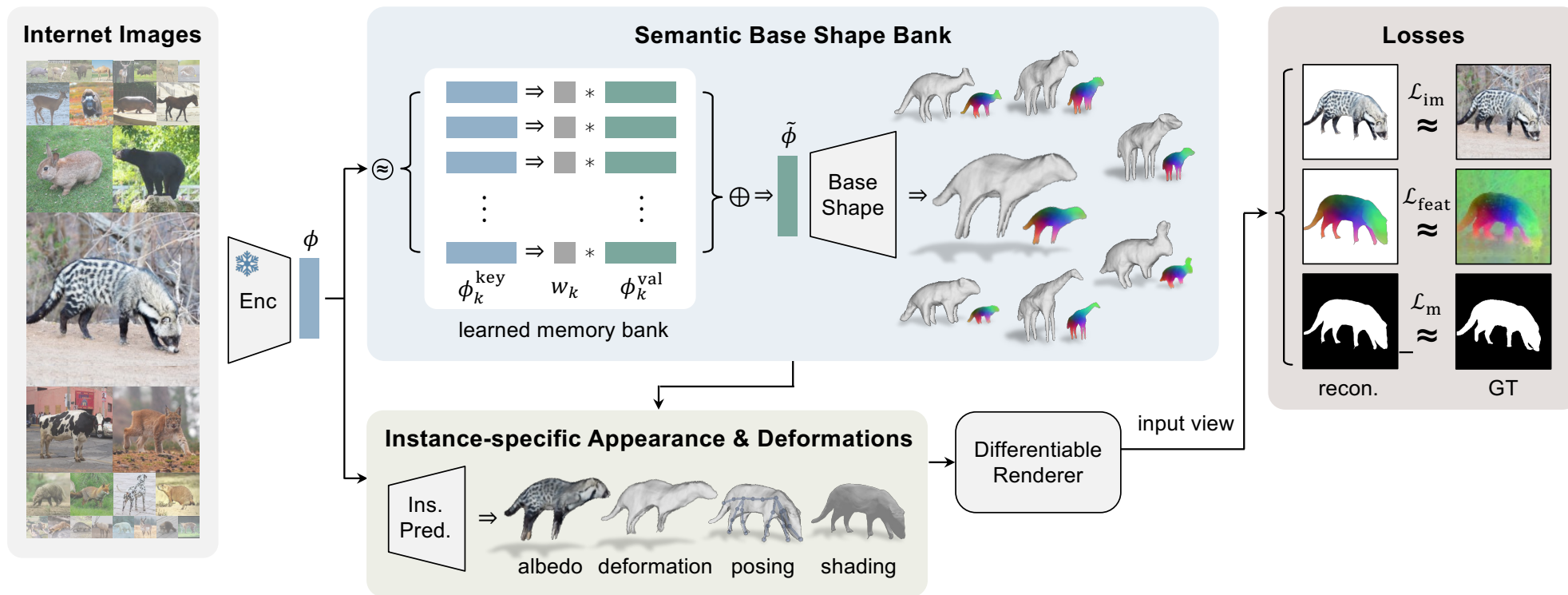
# Category embedding

- Leverage a pre-trained vision encoder - DINO

  - A memory bank to *distills* the category information and prevents overfitting
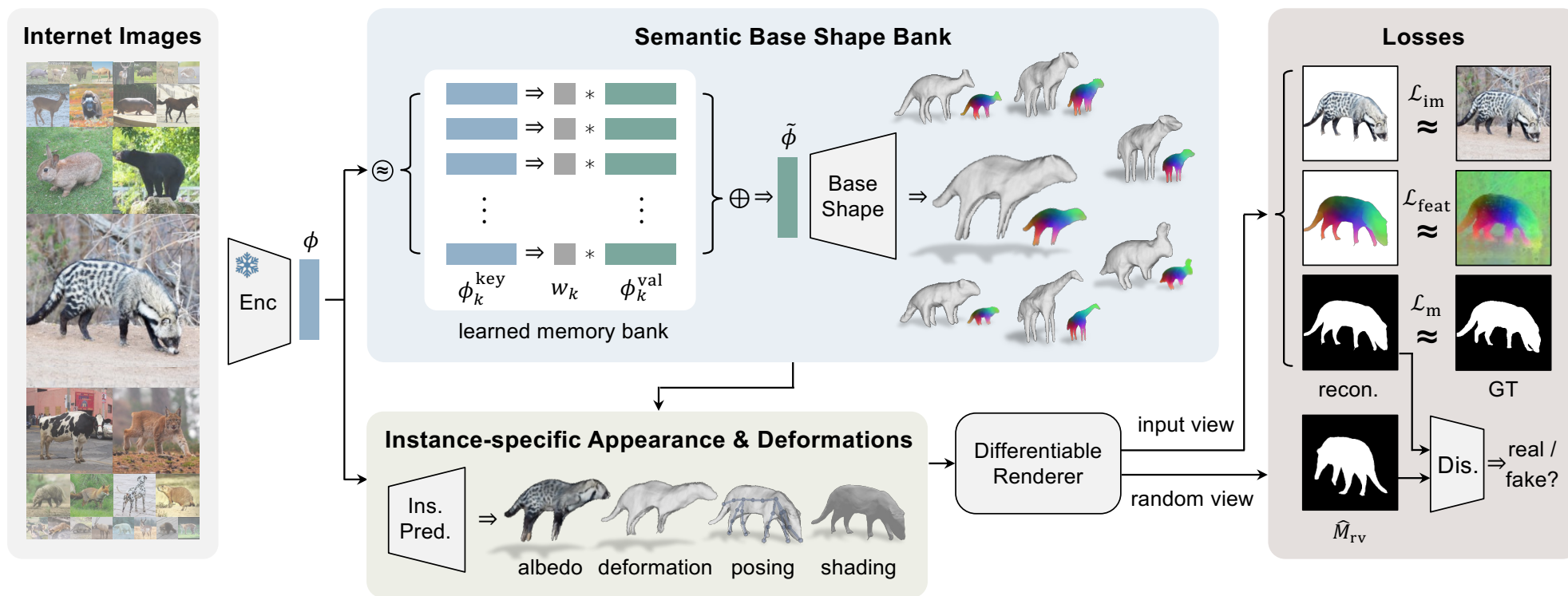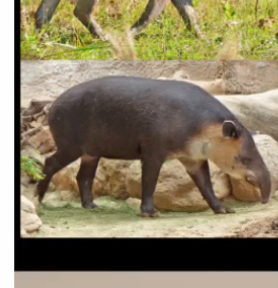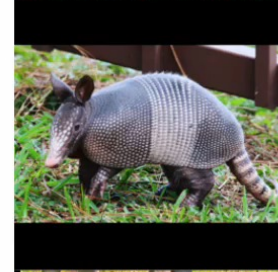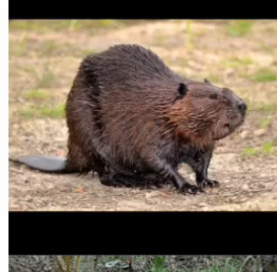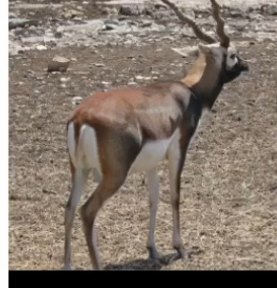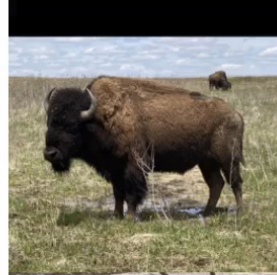
# Category embedding
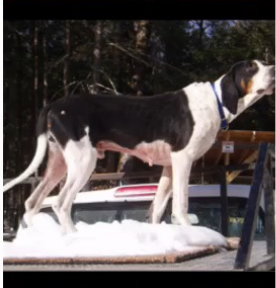
# Full pipeline

# Full pipeline
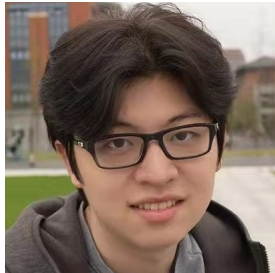
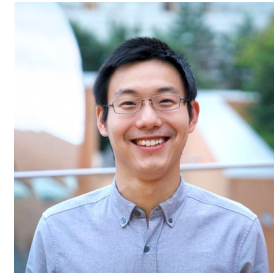Andrea Vedaldi     Shangzhe Wu     Christian Rupprecht     Ruining Li

Zizhang Li     Dor Litvak     Yunzhi Zhang     Jiajun Wu

# Learning Articulated 3D Animals from Internet Images

Tomas Jakab, University of Oxford, VGG