



# MagicPony: Learning Articulated 3D Animals in the Wild

Shangzhe Wu\* Ruining Li\* Tomas Jakab\* Christian Rupprecht Andrea Vedaldi

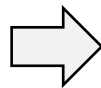
Visual Geometry Group, University of Oxford

(\* Equal Contribution)

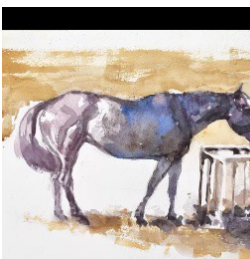
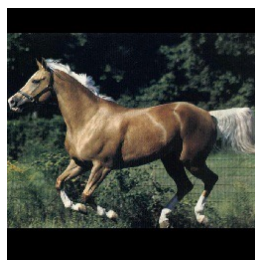
## Training



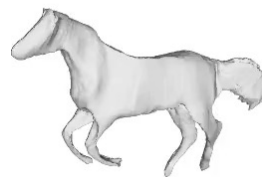
Single-view Images



## Single-Image Inference



Test Image



Articulated 3D Shape



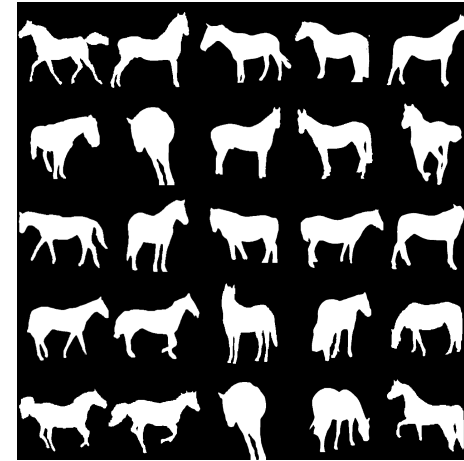
Animation

# Training Data



Single-view Images

No keypoint or viewpoint supervision,  
nor template shapes

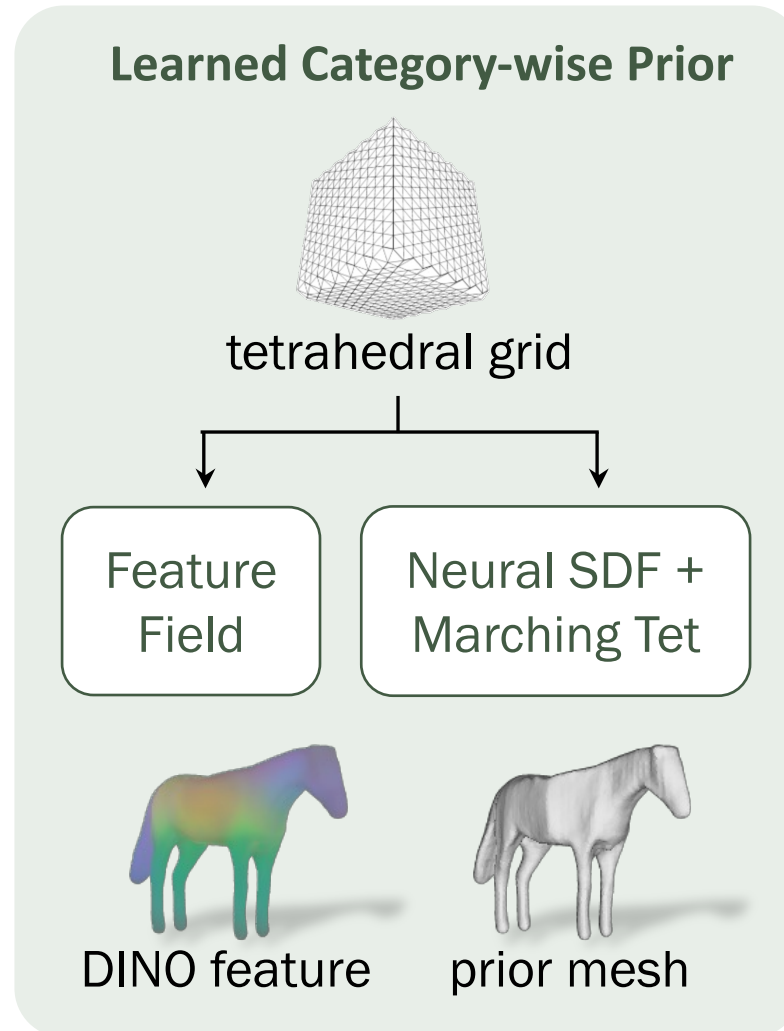


Instance Masks



Self-supervised Image Features

# Implicit-Explicit 3D Representation

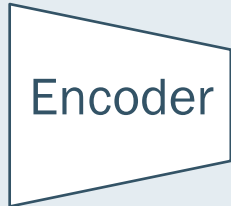


# Hierarchical Shape Prediction

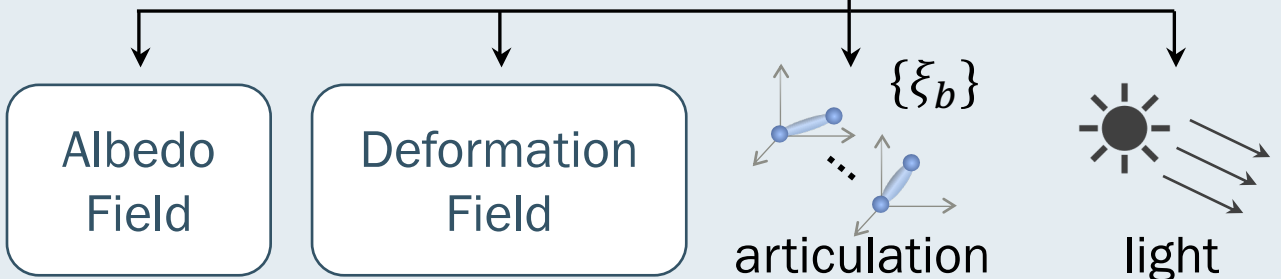
## Instance-specific Predictions



input image



feature

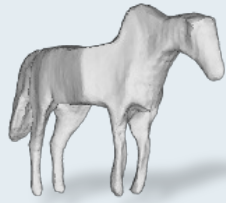


Albedo  
Field

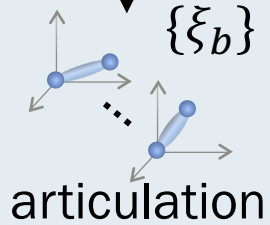


albedo

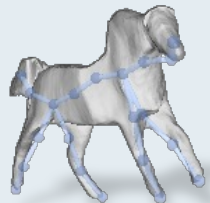
Deformation  
Field



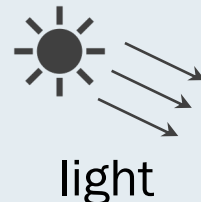
deformed



articulation



articulated



light



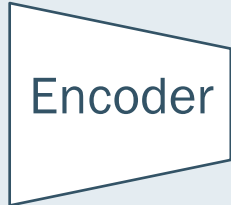
shading

# Hierarchical Shape Prediction

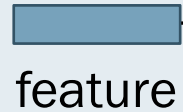
## Instance-specific Predictions



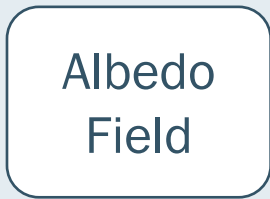
input image



Encoder



feature



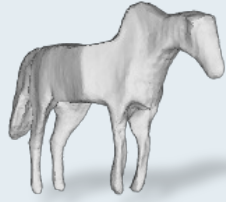
Albedo  
Field



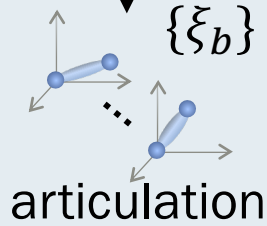
albedo



Deformation  
Field

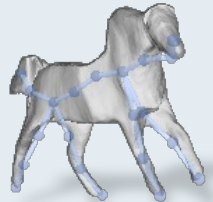


deformed

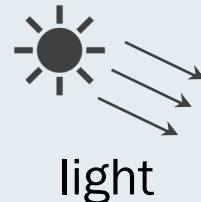


$\{\xi_b\}$

articulation



articulated

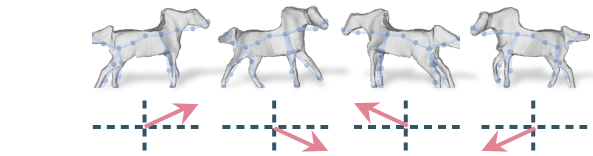
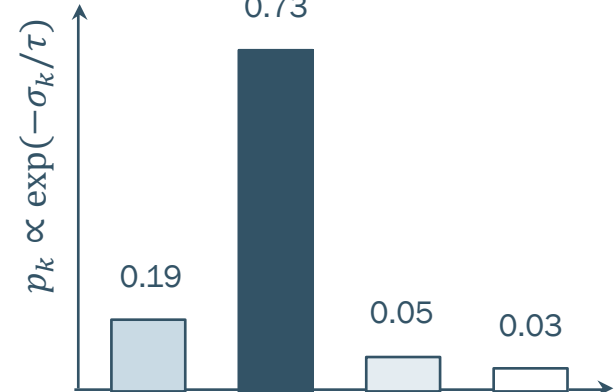


light



shading

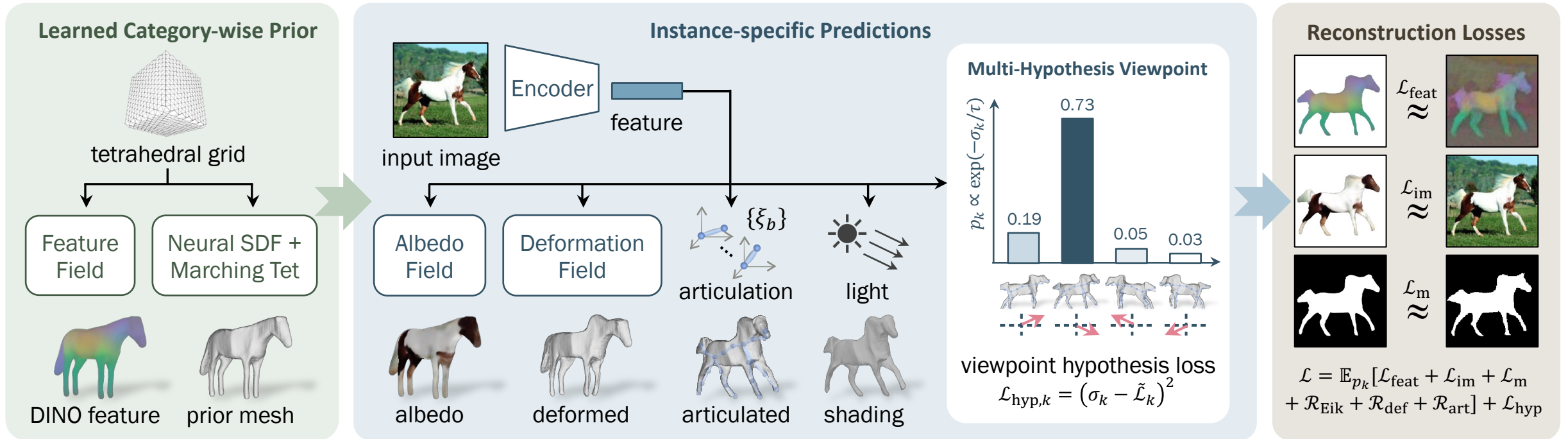
## Multi-Hypothesis Viewpoint



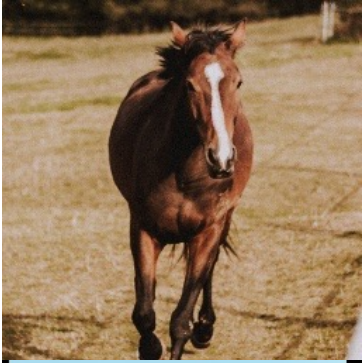
viewpoint hypothesis loss

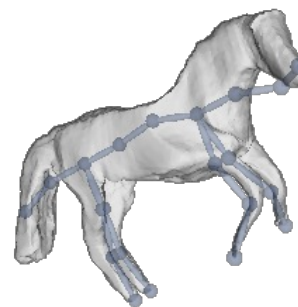
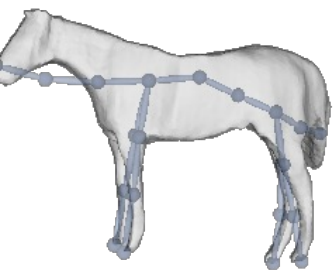
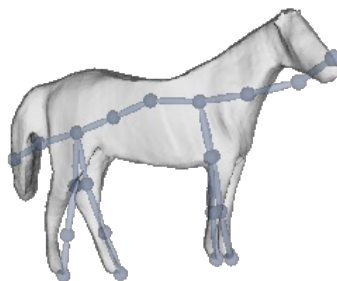
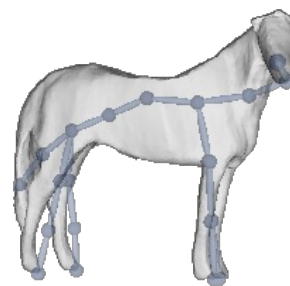
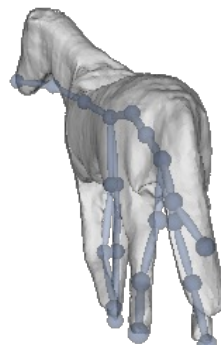
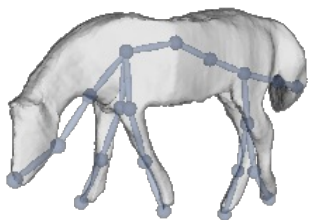
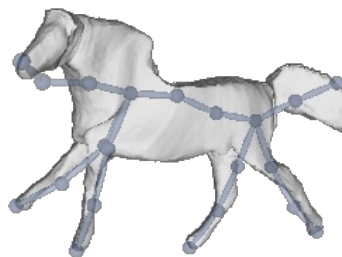
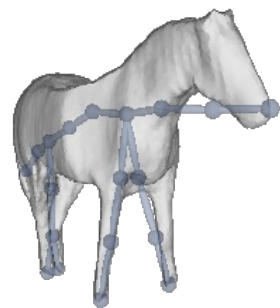
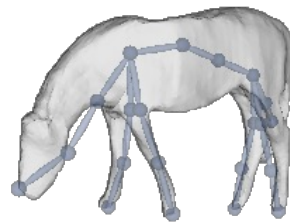
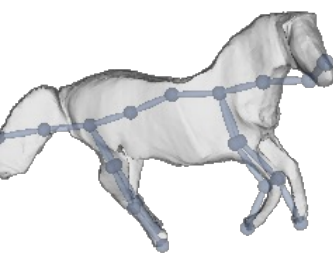
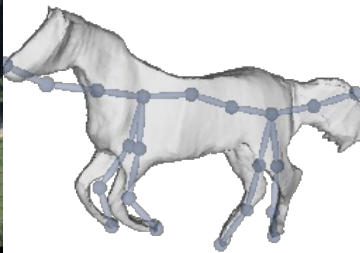
$$\mathcal{L}_{\text{hyp},k} = (\sigma_k - \tilde{\mathcal{L}}_k)^2$$

# Implicit-Explicit 3D Representation

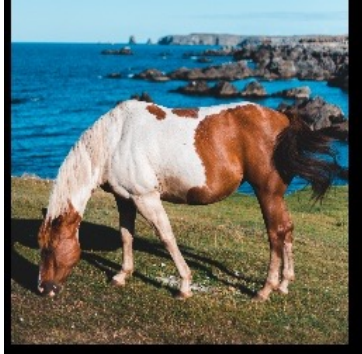
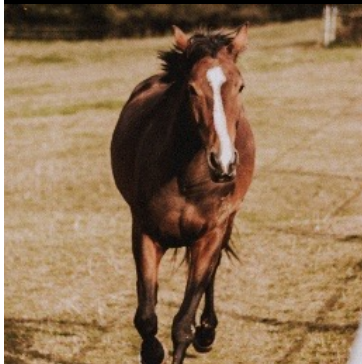


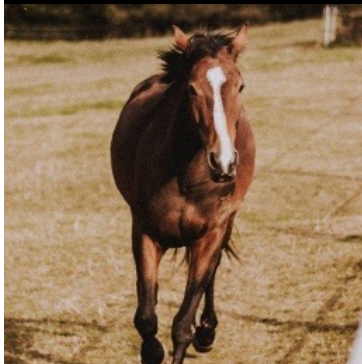
Entire pipeline trained end-to-end with reconstruction losses  
(except for frozen DINO-ViT [1] image encoder, pre-trained via self-supervision)

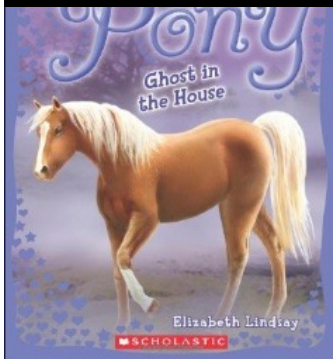
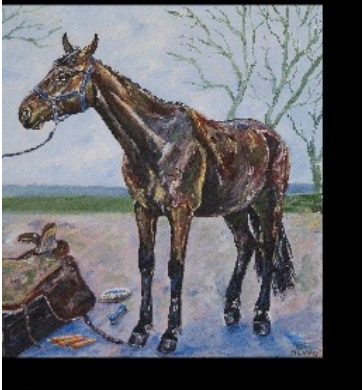
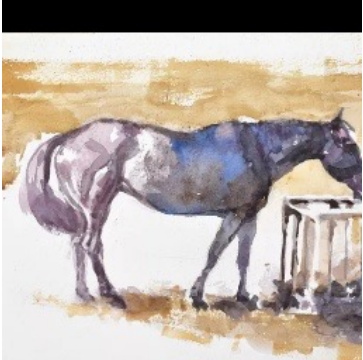


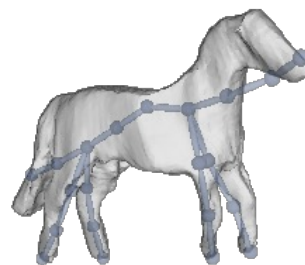
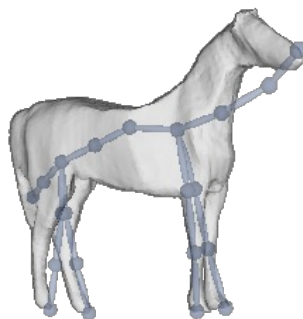
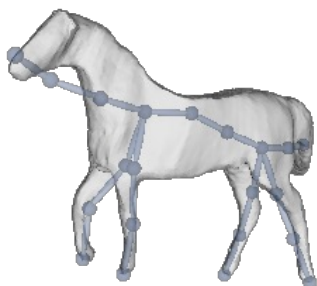
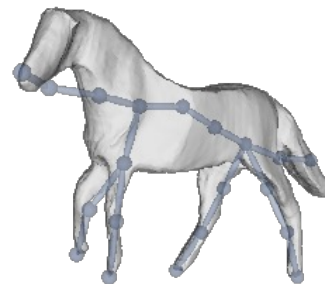
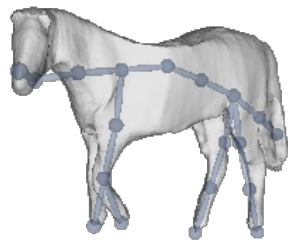
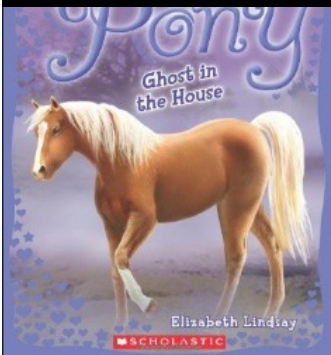
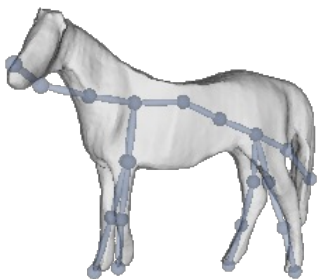
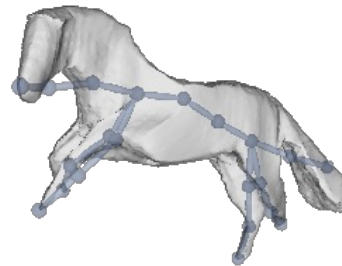
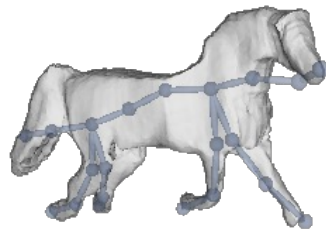
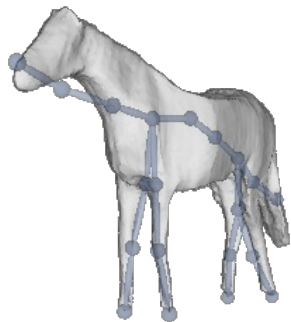
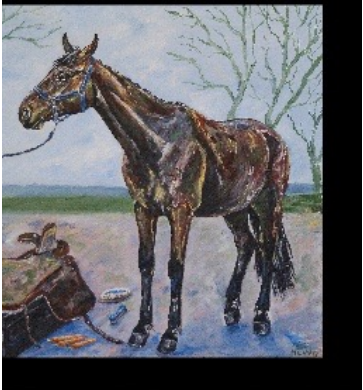
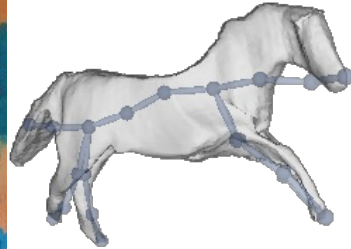
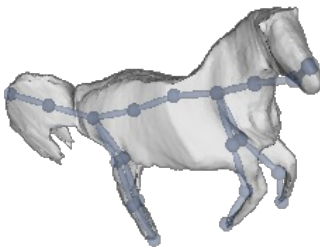
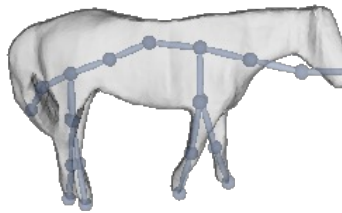
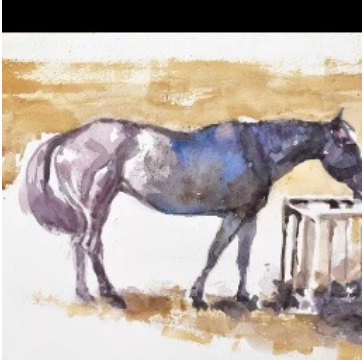


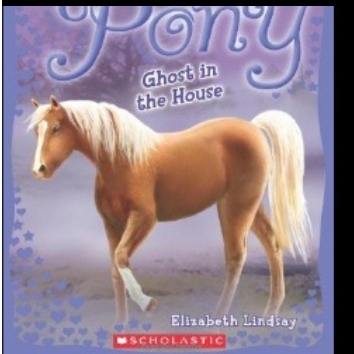
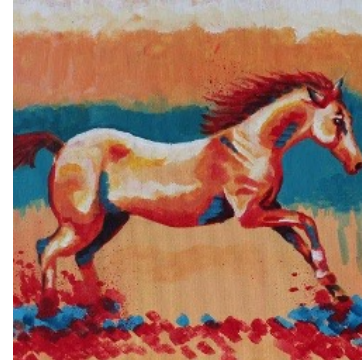


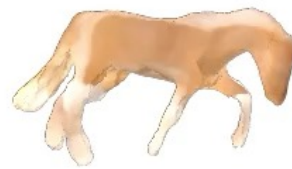
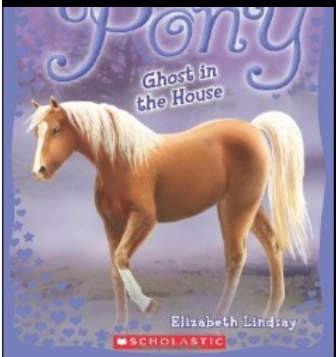


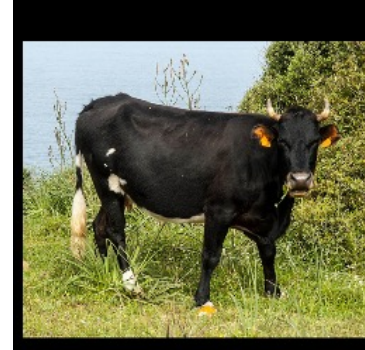
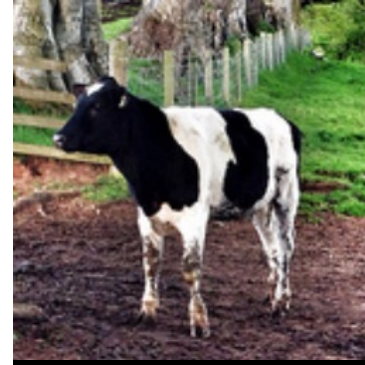
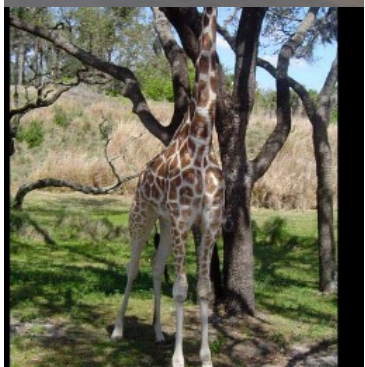
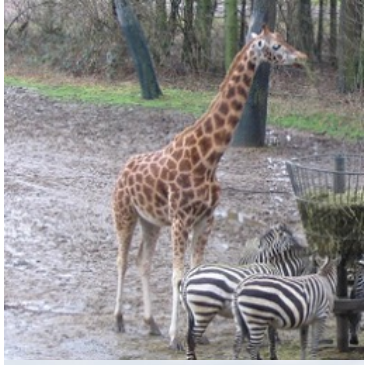


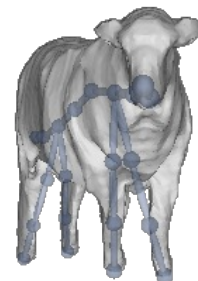
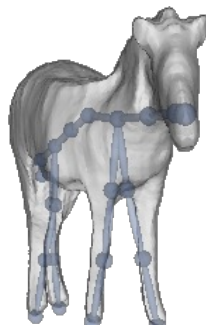
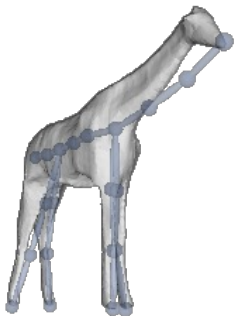
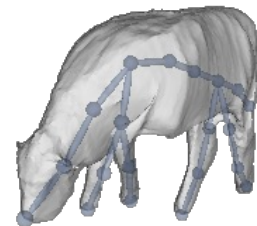
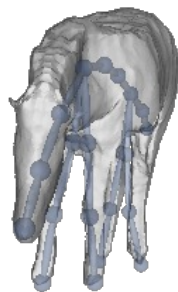
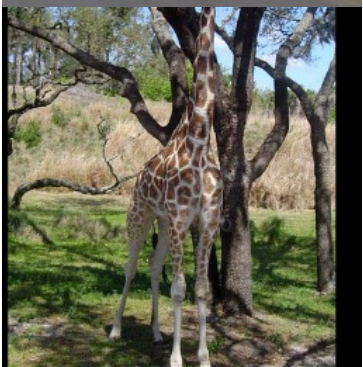
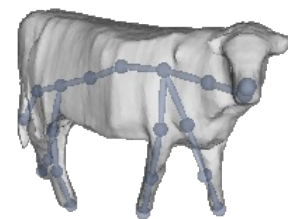
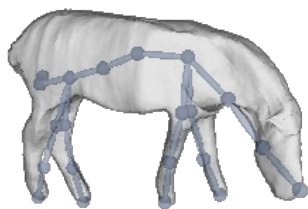
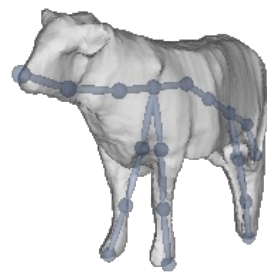
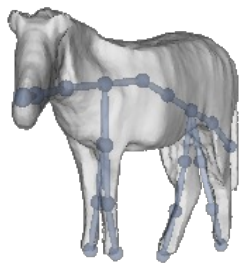
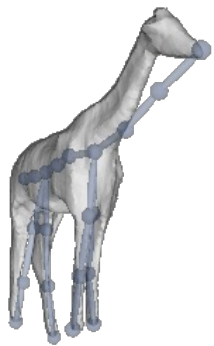




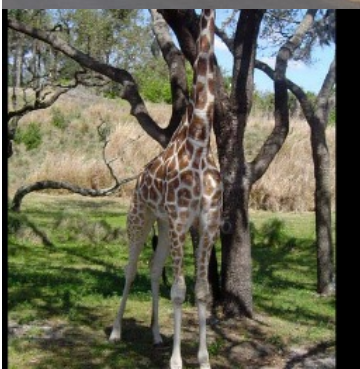


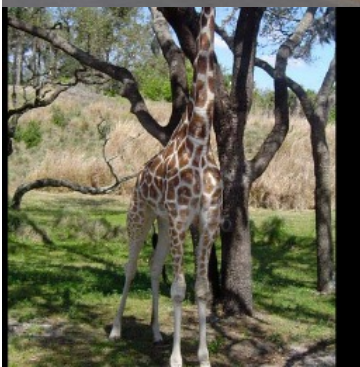








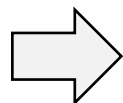
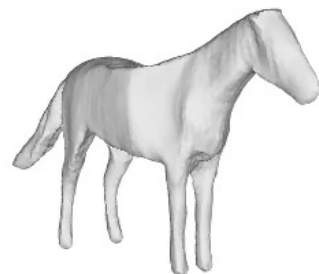
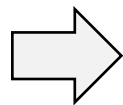
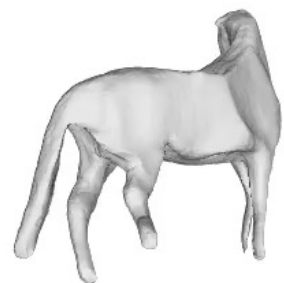
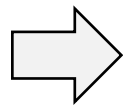




# Frame-by-Frame Inference on Videos

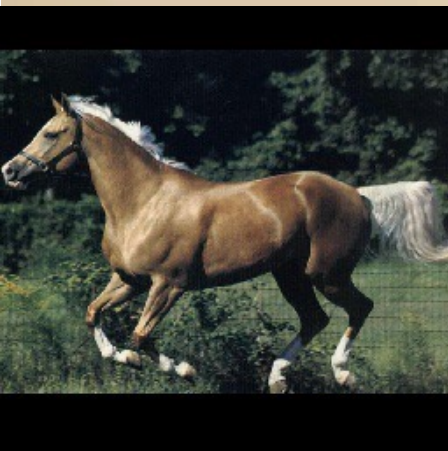


Input Frames



Input View

360° Rotations



3D Printed Horse Reconstruction