



3D Vision and
Robotics Lab

[CVPR 2022] BANMo: Building Animatable 3D Neural Models from Many Casual Videos

Gengshan Yang^{2*} Minh Vo³ Natalia Neverova¹ Deva Ramanan² Andrea Vedaldi¹ Hanbyul Joo¹
¹Meta AI ²Carnegie Mellon University ³Meta Reality Labs

Lab Seminar

Gyeongsu Cho

@UNIST

2023.09.20. (Wed)

Contents

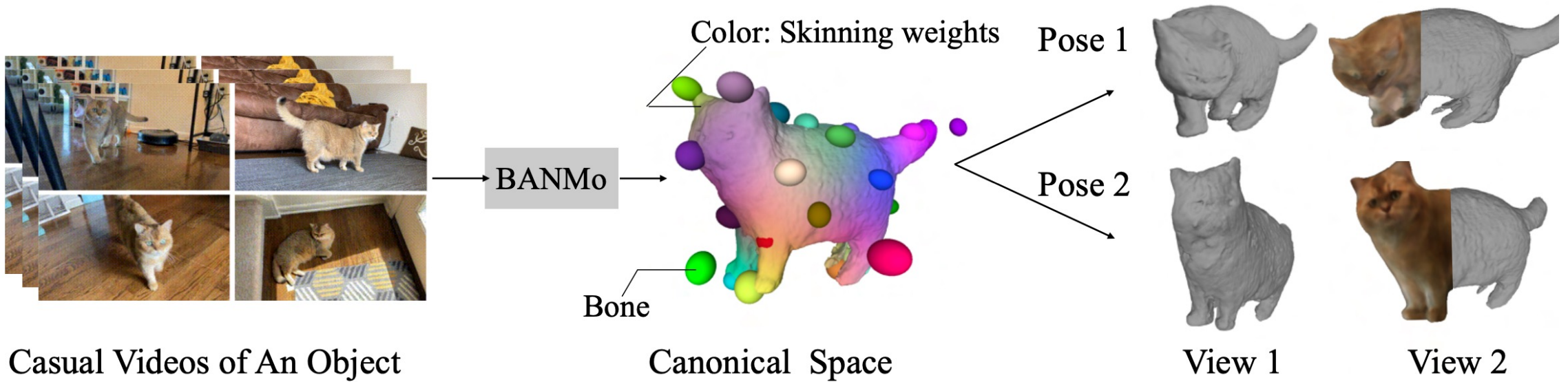
- Introduction
- Method
- Experiments
- Conclusion

Introduction

Goal: Animatable 3D Models from Causal Videos

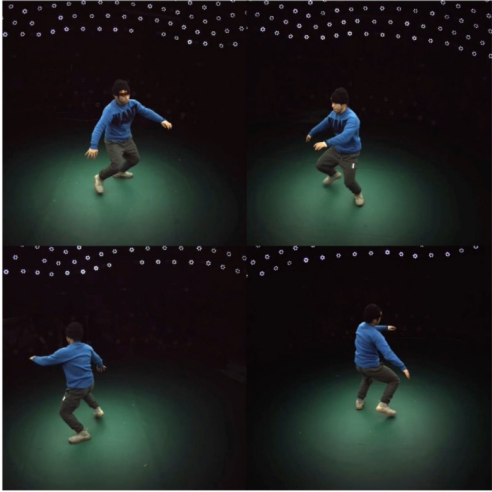
Introduction

Goal: Animatable 3D Models from Causal Videos

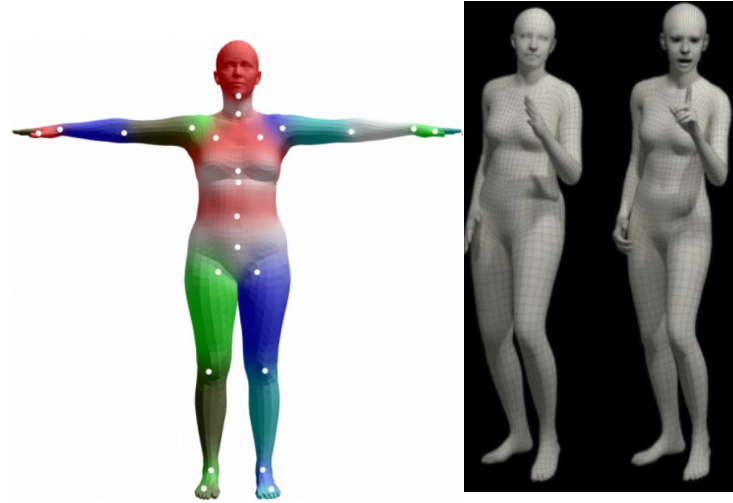


Introduction

Multi-view camera systems Category-specific body models



[CVPR 2021] Neural Body



[CVPR 2019] SMPL-X

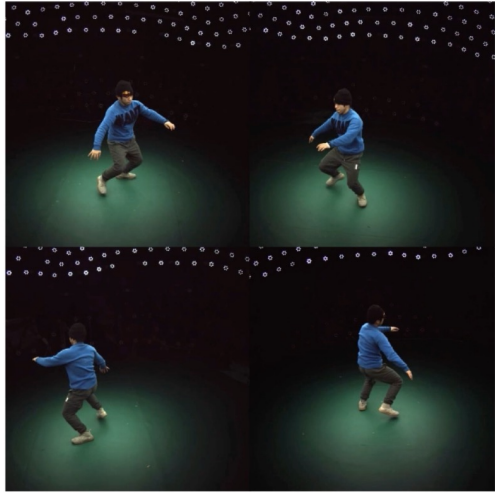
3D Reconstruction
with differentiable rendering



[ICCV 2021] Nerfies

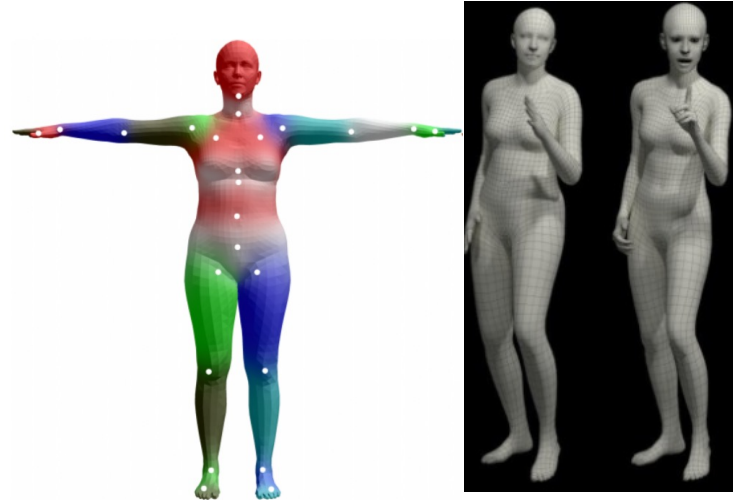
Introduction

Multi-view camera systems



[CVPR 2021] Neural Body

Category-specific body models



[CVPR 2019] SMPL-X

3D Reconstruction
with differentiable rendering



[ICCV 2021] Nerfies

Monocular casual videos

Template-free

Freely-moving targets

Introduction

Introduction

1. How to represent 3D geometry and appearance of the target in a canonical space?

Introduction

1. How to represent 3D geometry and appearance of the target in a canonical space?
2. How to deform 3D points between canonical space and individual time instances?

Introduction

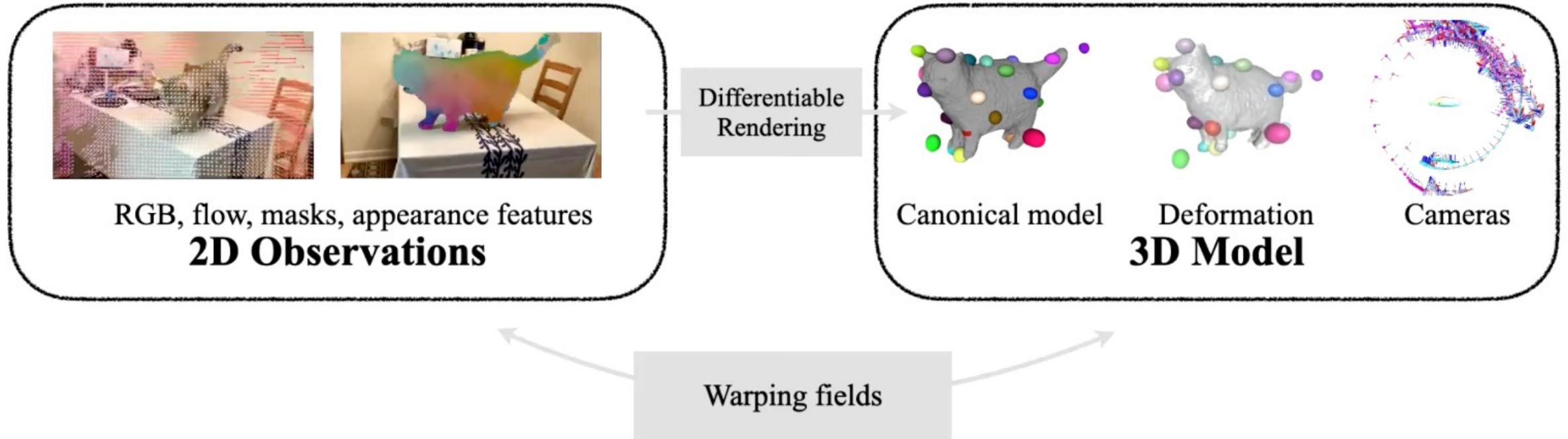
1. How to represent 3D geometry and appearance of the target in a canonical space?
2. How to deform 3D points between canonical space and individual time instances?
3. How to find pixel or part correspondences over videos given different viewpoint, lighting, background, and object deformations?

Introduction

1. How to represent 3D geometry and appearance of the target in a canonical space?
2. How to deform 3D points between canonical space and individual time instances?
3. How to find pixel or part correspondences over videos given different viewpoint, lighting, background, and object deformations?

Key Challenge : How to find the warping between the image and canonical 3D space?

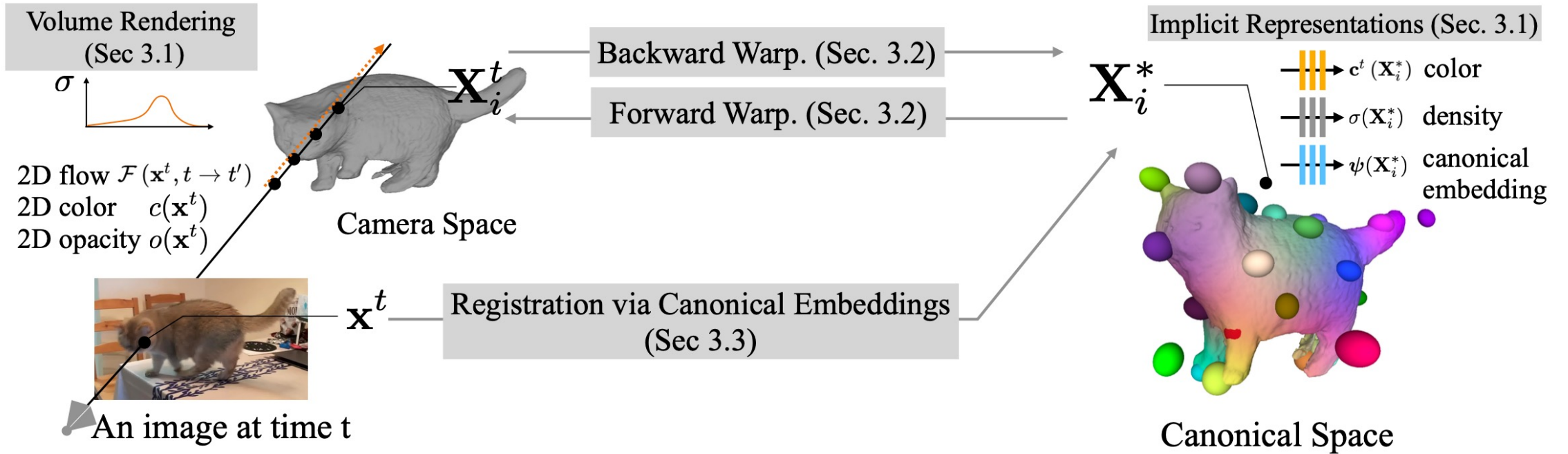
Introduction



Key Challenge : How to find the warping between the image and canonical 3D space?

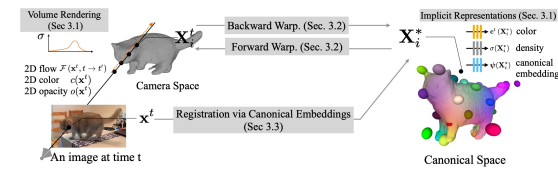
Method

Method



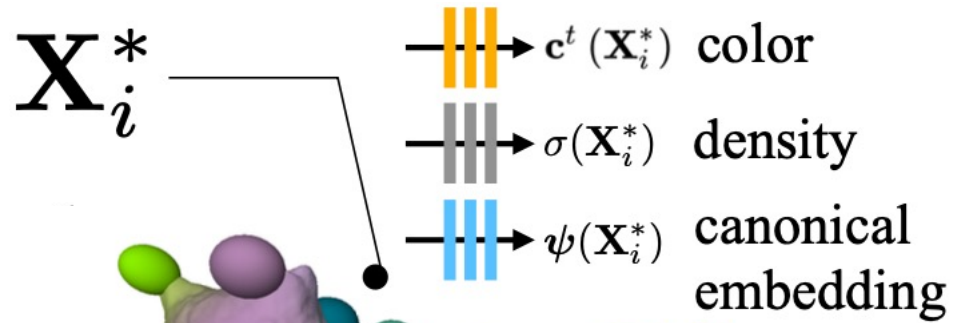
Method

1. Implicit 3D Representation
2. Differentiable Volume Rendering
3. Deformation Model via Neural Blend Skinning
4. Registration
5. Optimization

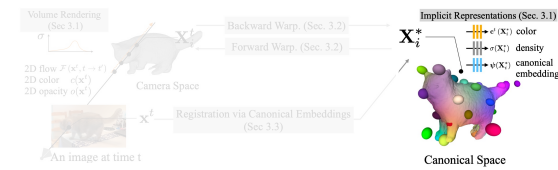


Method

1. Implicit 3D Representation

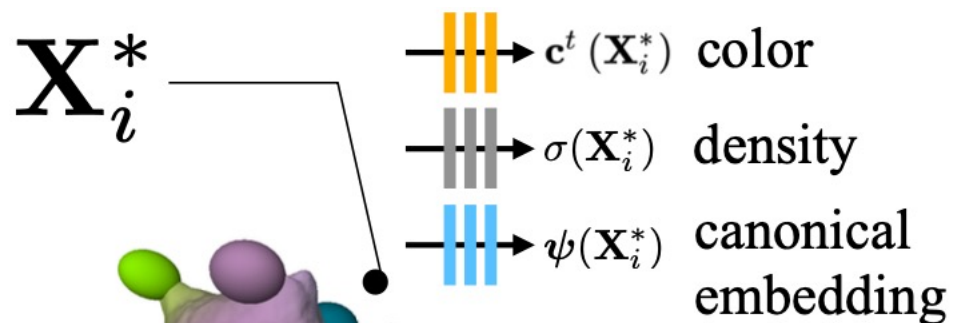


Canonical Space

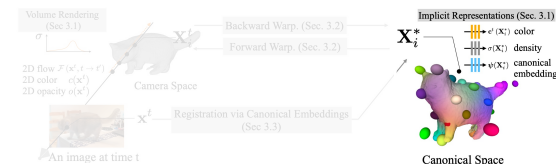


Method

1. Implicit 3D Representation



Canonical Space



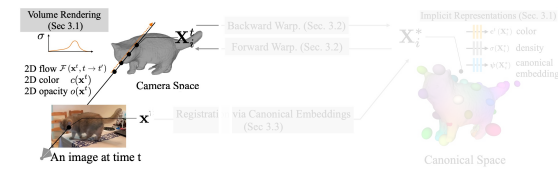
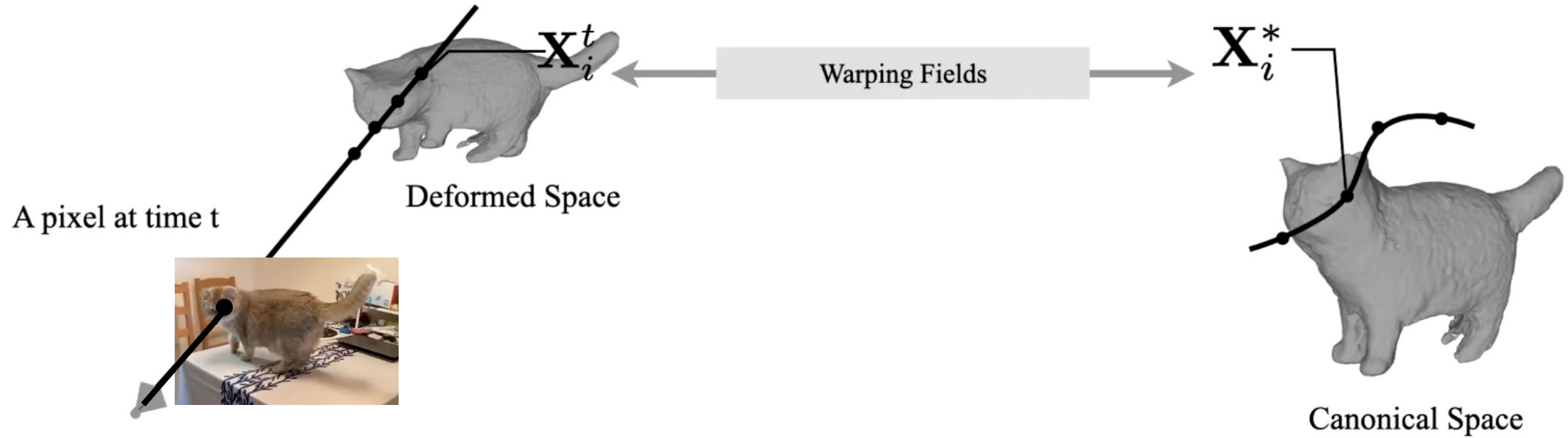
$$\mathbf{c}^t = \text{MLP}_{\mathbf{c}}(\mathbf{X}^*, \mathbf{v}^t, \omega_e^t), \quad (1)$$

$$\sigma = \Gamma_{\beta}(\text{MLP}_{\text{SDF}}(\mathbf{X}^*)), \quad (2)$$

$$\psi = \text{MLP}_{\psi}(\mathbf{X}^*). \quad (3)$$

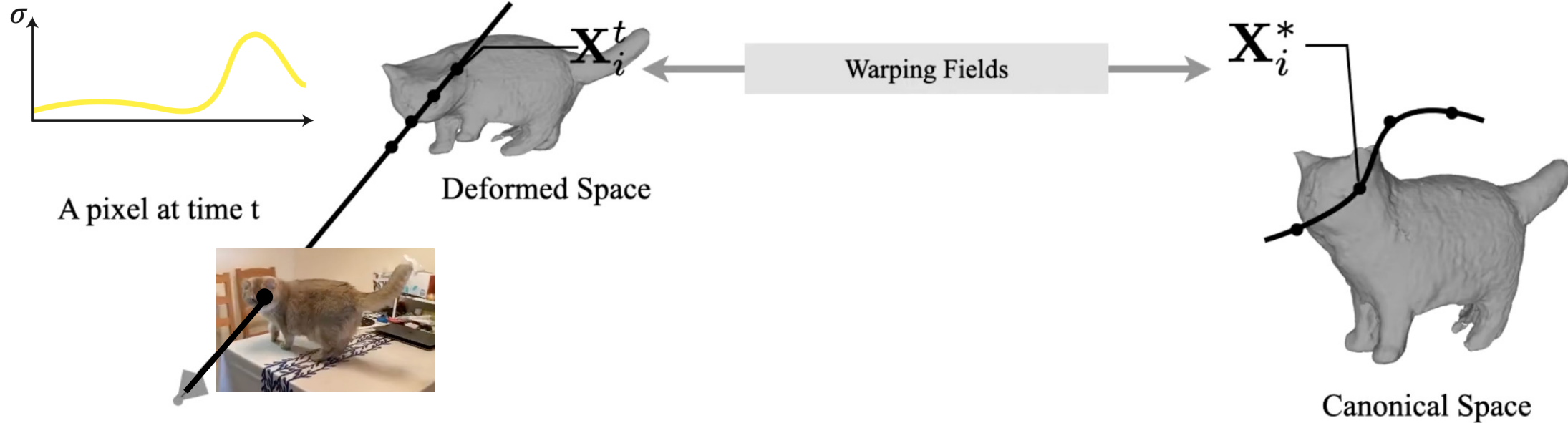
Method

2. Differentiable Volume Rendering



Method

2. Differentiable Volume Rendering

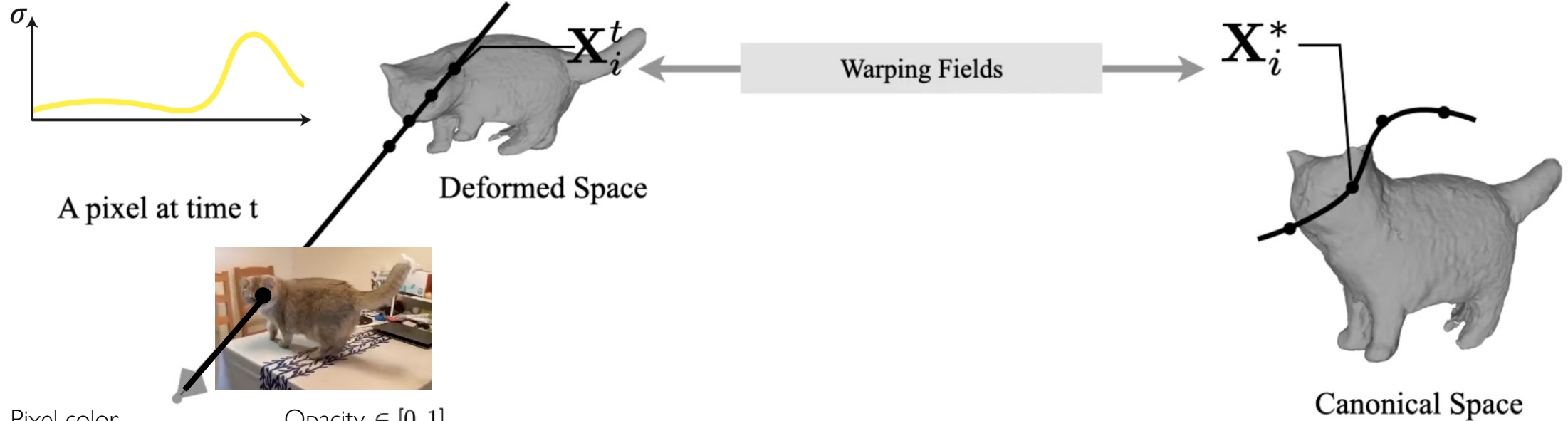
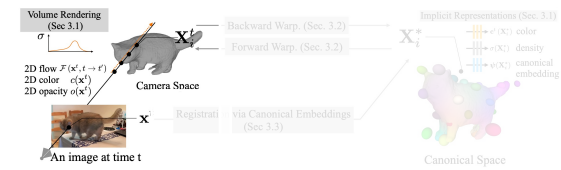


$$\mathbf{c}(\mathbf{x}^t) = \sum_{i=1}^N \tau_i \mathbf{c}_i^t, \quad o(\mathbf{x}^t) = \sum_{i=1}^N \tau_i, \quad (4)$$

$$\mathbf{X}^*(\mathbf{x}^t) = \sum_{i=1}^N \tau_i \mathbf{X}_i^*. \quad (5)$$

Method

2. Differentiable Volume Rendering



Pixel color Opacity $\in [0, 1]$

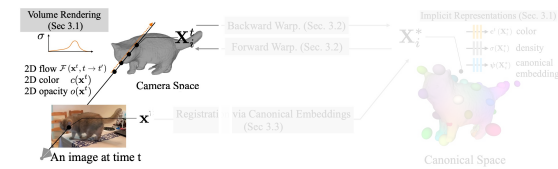
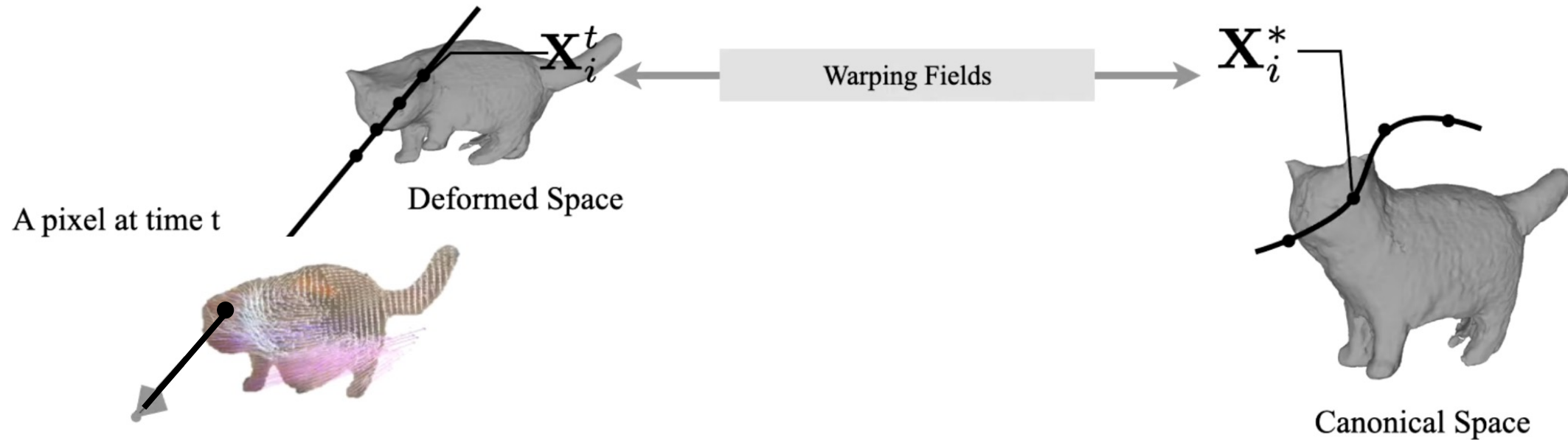
$$\mathbf{c}(\mathbf{x}^t) = \sum_{i=1}^N \tau_i \mathbf{c}_i^t, \quad o(\mathbf{x}^t) = \sum_{i=1}^N \tau_i, \quad (4)$$

$$\mathbf{X}^*(\mathbf{x}^t) = \sum_{i=1}^N \tau_i \mathbf{X}_i^*. \quad (5)$$

expected surface intersection

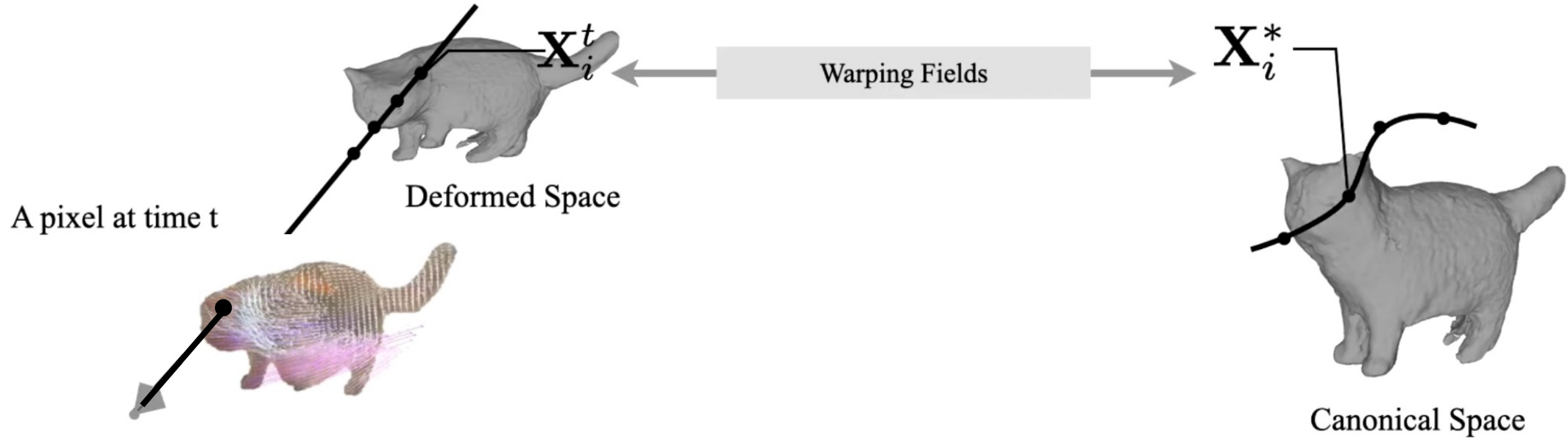
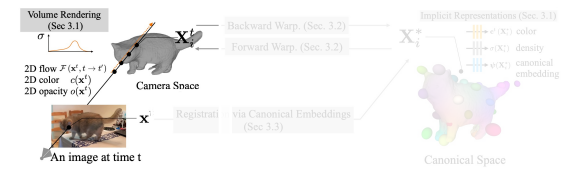
Method

2. Differentiable Volume Rendering



Method

2. Differentiable Volume Rendering



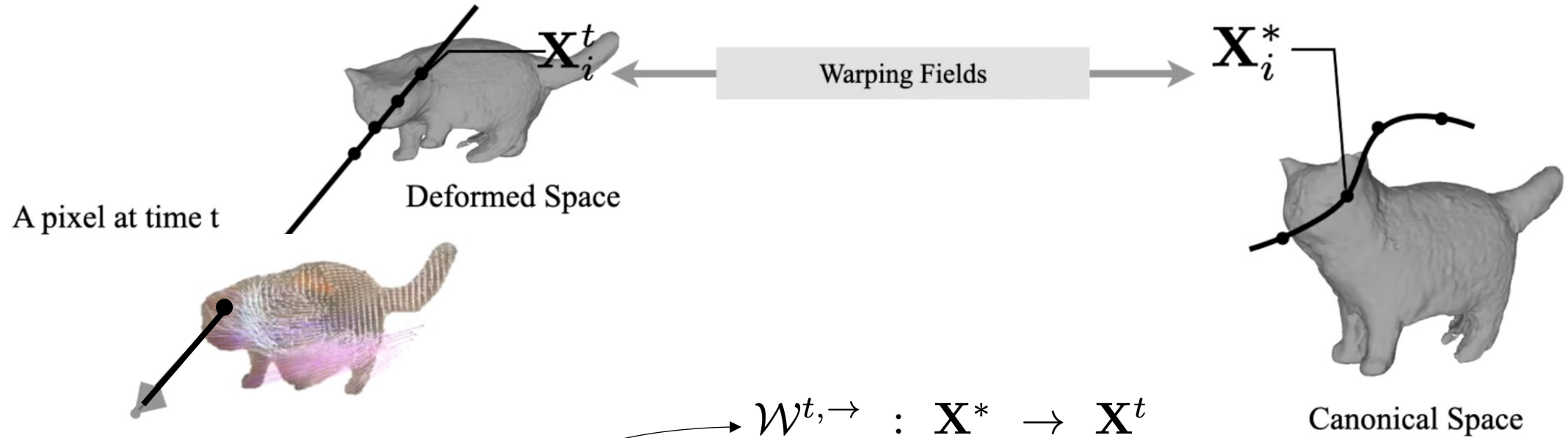
$$\mathbf{x}^{t'} = \sum_{i=1}^N \tau_i \Pi^{t'} \left(\mathcal{W}^{t', \rightarrow} (\mathbf{X}_i^*) \right), \quad (6)$$

$$\mathcal{F}(\mathbf{x}^t, t \rightarrow t') = \mathbf{x}^{t'} - \mathbf{x}^t. \quad (7)$$

for optical flow loss

Method

2. Differentiable Volume Rendering



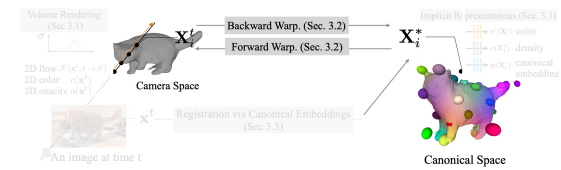
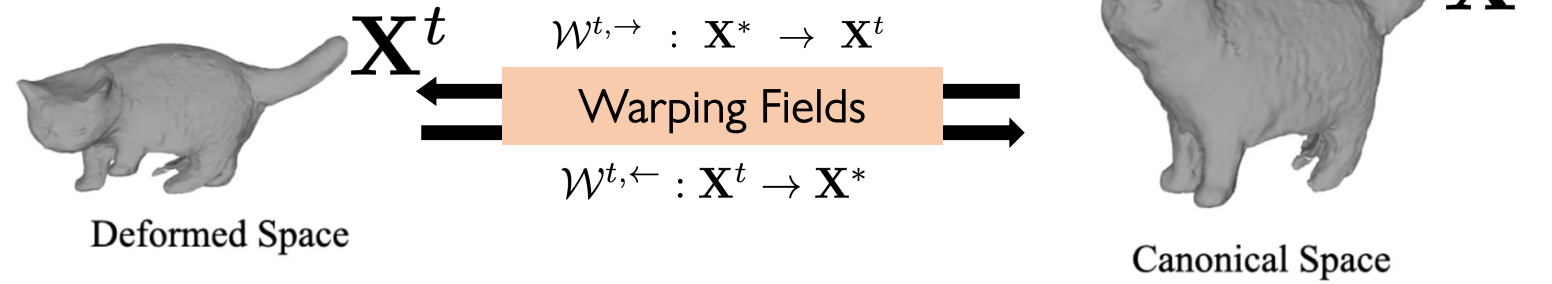
expected 2D re-projection

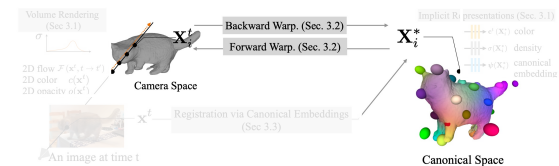
$$\mathbf{x}^{t'} = \sum_{i=1}^N \tau_i \Pi^{t'} \left(\mathcal{W}^{t', \rightarrow} (\mathbf{X}_i^*) \right), \quad (6)$$

$$\mathcal{F}(\mathbf{x}^t, t \rightarrow t') = \mathbf{x}^{t'} - \mathbf{x}^t. \quad (7) \quad \text{for optical flow loss}$$

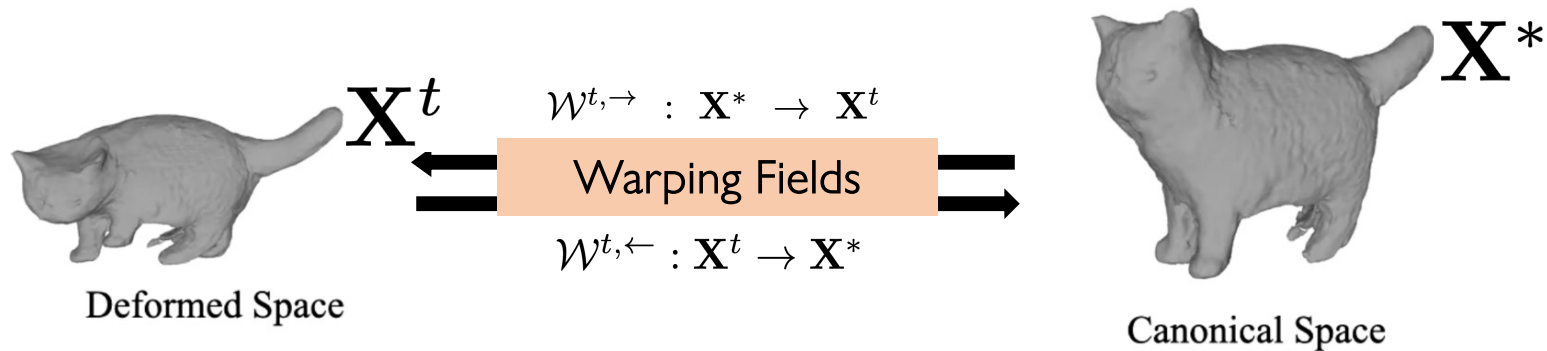
2D flow

3. Deformation Model via Neural Blend Skinning





3. Deformation Model via Neural Blend Skinning



$$\mathbf{X}^t = \mathcal{W}^{t, \rightarrow}(\mathbf{X}^*) = \mathbf{G}^t \mathbf{J}^{t, \rightarrow} \mathbf{X}^*, \quad (8)$$

$$\mathbf{W}_\sigma = (\mathbf{X} - \mathbf{C}_b)^T \mathbf{Q}_b (\mathbf{X} - \mathbf{C}_b), \quad (13)$$

$$\mathbf{X}^* = \mathcal{W}^{t, \leftarrow}(\mathbf{X}^t) = \mathbf{J}^{t, \leftarrow} (\mathbf{G}^t)^{-1} \mathbf{X}^t, \quad (9)$$

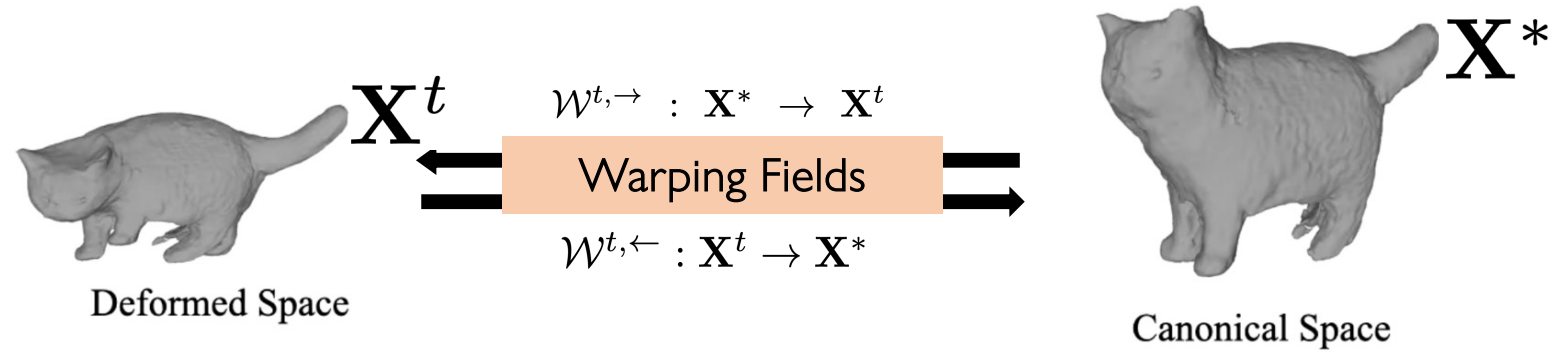
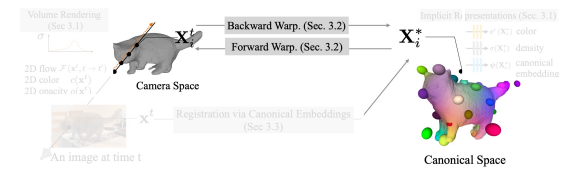
$$\mathbf{W} = \mathcal{S}(\mathbf{X}, \omega_b) = \sigma_{\text{softmax}}(\mathbf{W}_\sigma + \mathbf{W}_\Delta) \quad (14)$$

$$\mathbf{J}^{t, \rightarrow} = \sum_{b=1}^B \mathbf{W}_b^{t, \rightarrow} \mathbf{J}_b^t, \quad \mathbf{J}^{t, \leftarrow} = \sum_{b=1}^B \mathbf{W}_b^{t, \leftarrow} (\mathbf{J}_b^t)^{-1} \quad (10)$$

$$\mathbf{G}^t = \text{MLP}_{\mathbf{G}}(\omega_r^t), \quad \mathbf{J}_b^t = \text{MLP}_{\mathbf{J}}(\omega_b^t) \quad (11)$$

$$\omega_t^b = \mathbf{A}_i \mathcal{F}(t) \quad (12)$$

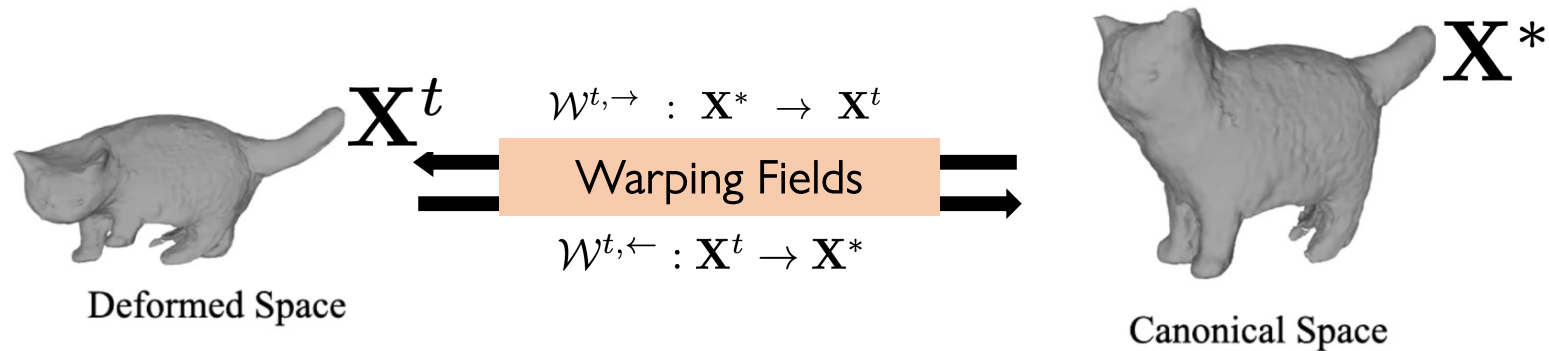
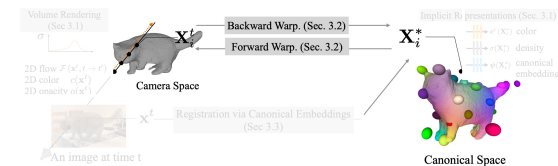
3. Deformation Model via Neural Blend Skinning



$$\mathbf{X}^t = \mathcal{W}^{t, \rightarrow}(\mathbf{X}^*) = \mathbf{G}^t \mathbf{J}^{t, \rightarrow} \mathbf{X}^*, \quad (8)$$

$$\mathbf{X}^* = \mathcal{W}^{t, \leftarrow}(\mathbf{X}^t) = \mathbf{J}^{t, \leftarrow} (\mathbf{G}^t)^{-1} \mathbf{X}^t, \quad (9)$$

3. Deformation Model via Neural Blend Skinning

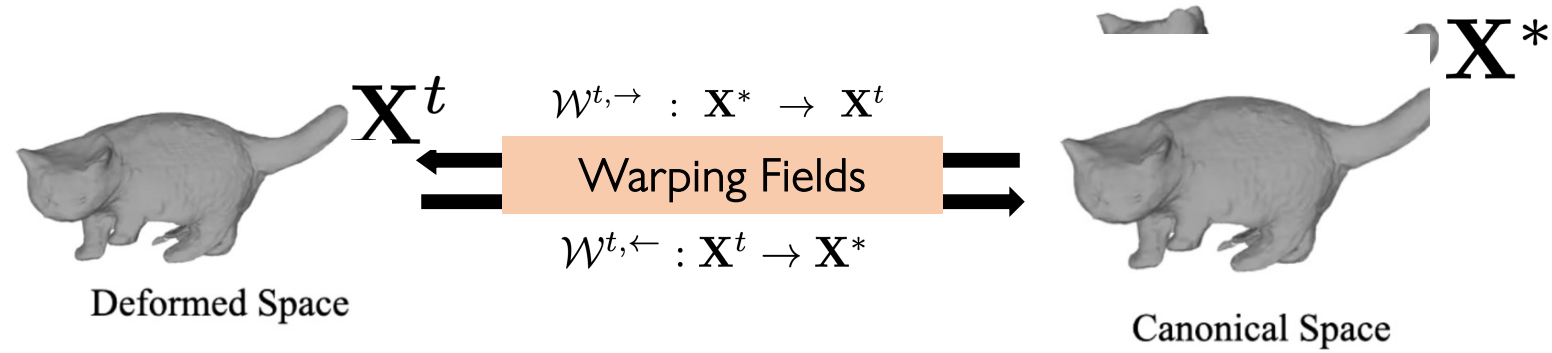
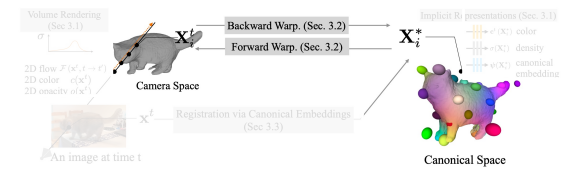


root body pose Bone(joint) pose

$$\mathbf{X}^t = \mathcal{W}^{t, \rightarrow}(\mathbf{X}^*) = \mathbf{G}^t \mathbf{J}^{t, \rightarrow} \mathbf{X}^*, \quad (8)$$

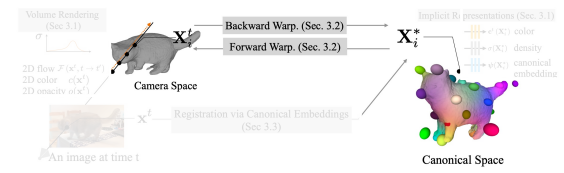
$$\mathbf{X}^* = \mathcal{W}^{t, \leftarrow}(\mathbf{X}^t) = \mathbf{J}^{t, \leftarrow} (\mathbf{G}^t)^{-1} \mathbf{X}^t, \quad (9)$$

3. Deformation Model via Neural Blend Skinning

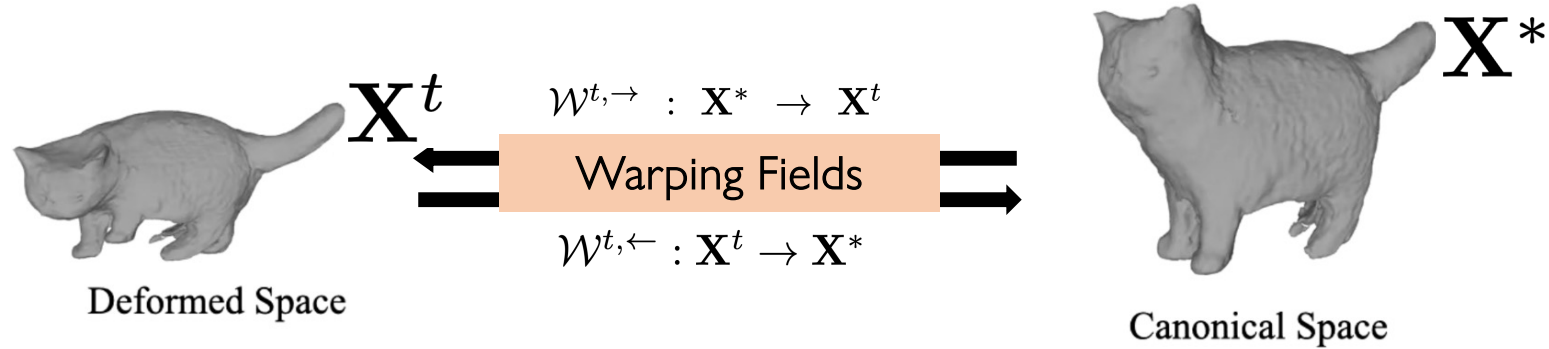


$$\mathbf{X}^t = \mathcal{W}^{t, \rightarrow}(\mathbf{X}^*) = \mathbf{G}^t \mathbf{J}^{t, \rightarrow} \mathbf{X}^*, \quad (8)$$

$$\mathbf{X}^* = \mathcal{W}^{t, \leftarrow}(\mathbf{X}^t) = \mathbf{J}^{t, \leftarrow} (\mathbf{G}^t)^{-1} \mathbf{X}^t, \quad (9)$$



3. Deformation Model via Neural Blend Skinning



root body pose Bone(joint) pose

$$\mathbf{X}^t = \mathcal{W}^{t, \rightarrow}(\mathbf{X}^*) = \mathbf{G}^t \mathbf{J}^{t, \rightarrow} \mathbf{X}^*, \quad (8) \quad \mathbf{W}_\sigma = (\mathbf{X} - \mathbf{C}_b)^T \mathbf{Q}_b (\mathbf{X} - \mathbf{C}_b), \quad (13)$$

$$\mathbf{X}^* = \mathcal{W}^{t, \leftarrow}(\mathbf{X}^t) = \mathbf{J}^{t, \leftarrow} (\mathbf{G}^t)^{-1} \mathbf{X}^t, \quad (9) \quad \mathbf{W} = \mathcal{S}(\mathbf{X}, \omega_b) = \sigma_{\text{softmax}}(\mathbf{W}_\sigma + \mathbf{W}_\Delta) \quad (14)$$

$$\mathbf{J}^{t, \rightarrow} = \sum_{b=1}^B \mathbf{W}_b^{t, \rightarrow} \mathbf{J}_b^t, \quad \mathbf{J}^{t, \leftarrow} = \sum_{b=1}^B \mathbf{W}_b^{t, \leftarrow} (\mathbf{J}_b^t)^{-1} \quad (10)$$

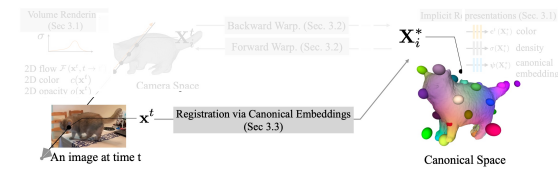
$$\mathbf{G}^t = \text{MLP}_{\mathbf{G}}(\omega_r^t), \quad \mathbf{J}_b^t = \text{MLP}_{\mathbf{J}}(\omega_b^t) \quad (11)$$

$$\omega_t^b = \mathbf{A}_i \mathcal{F}(t) \quad (12)$$

latent code

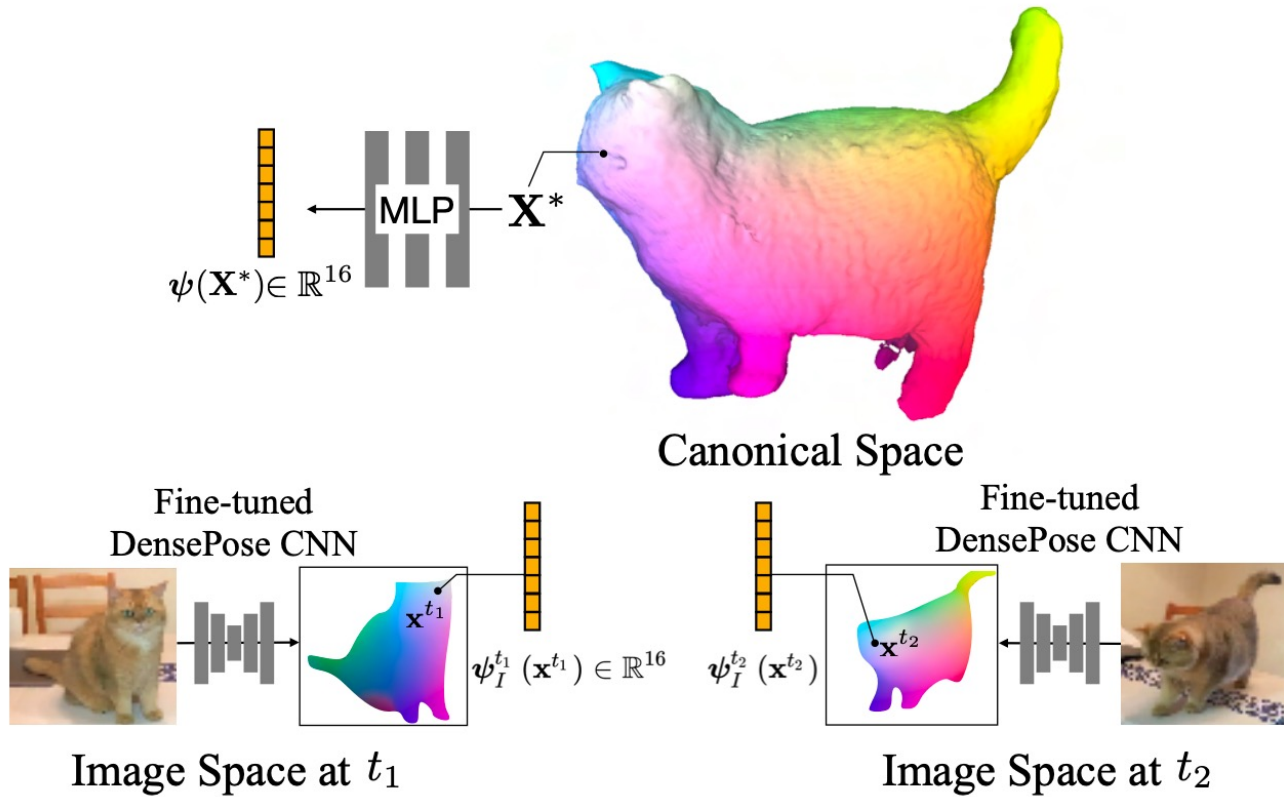
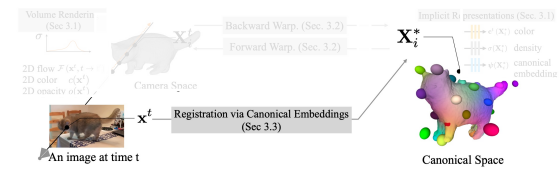
Method

4. Registration



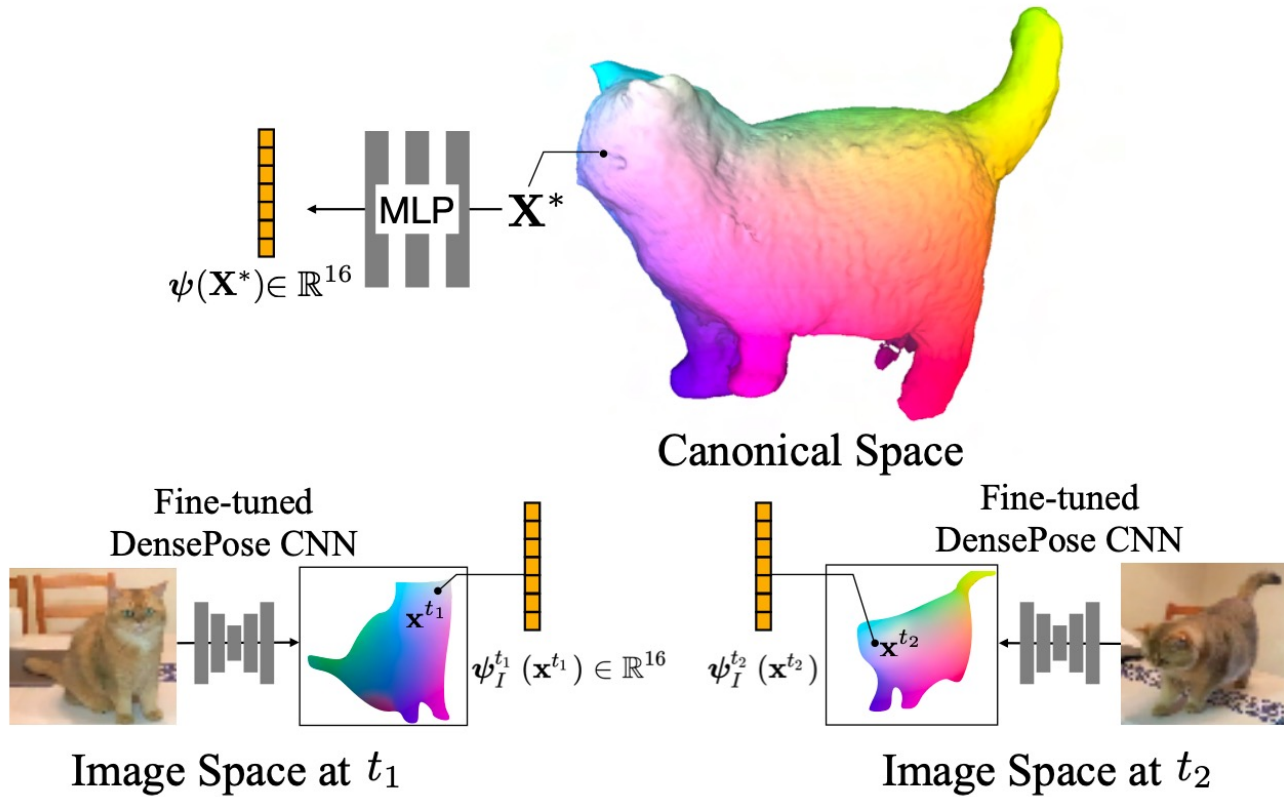
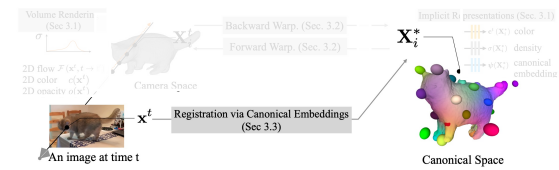
Method

4. Registration



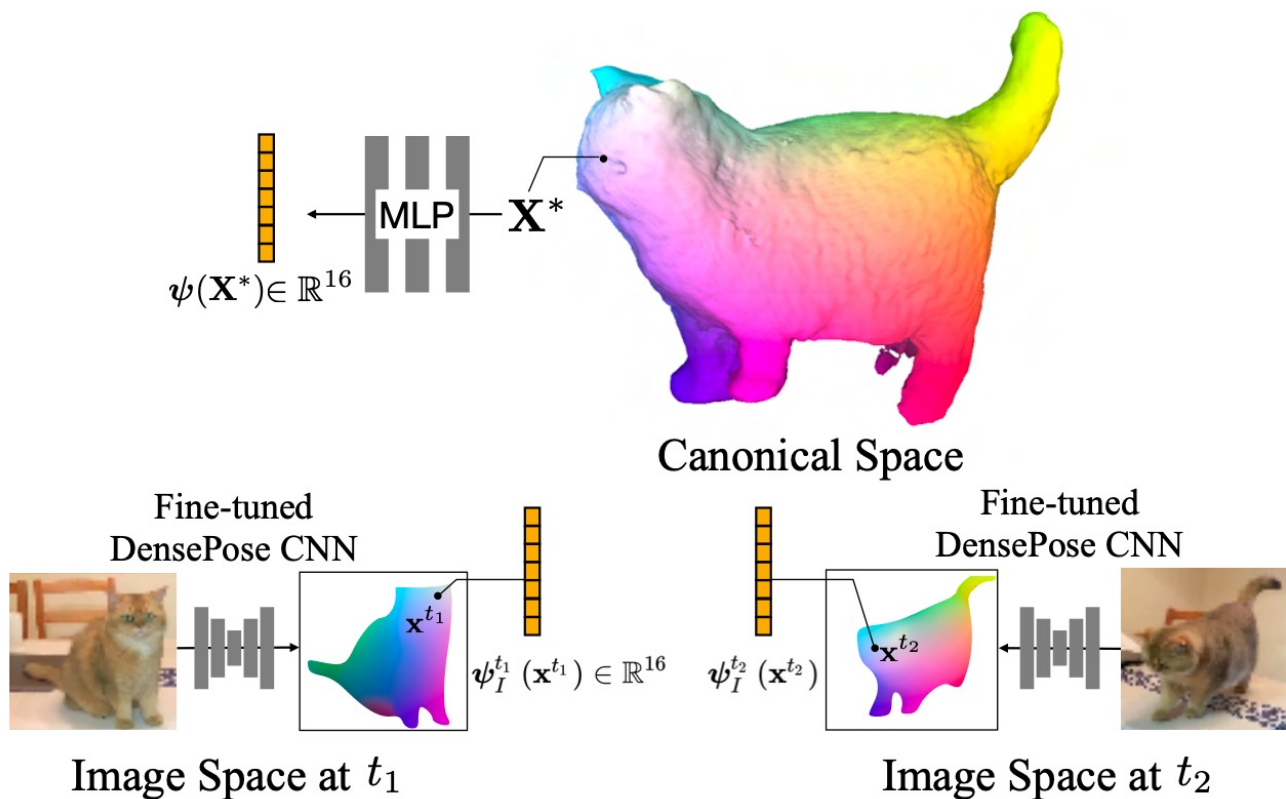
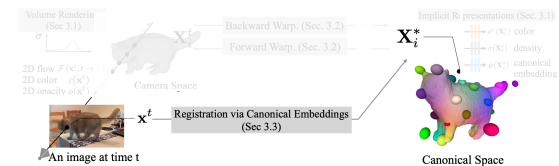
Method

4. Registration



$$\hat{\mathbf{X}}^*(\mathbf{x}^t) = \sum_{\mathbf{X} \in \mathbf{V}^*} \tilde{\mathbf{s}}^t(\mathbf{x}^t) \mathbf{X}, \quad (15)$$

4. Registration



$$\hat{\mathbf{X}}^*(\mathbf{x}^t) = \sum_{\mathbf{X} \in \mathbf{V}^*} \tilde{\mathbf{s}}^t(\mathbf{x}^t) \mathbf{X}, \quad (15)$$

$$\tilde{\mathbf{s}}^t(\mathbf{x}^t) = \sigma_{\text{softmax}} \left(\alpha_s \langle \psi_I^t(\mathbf{x}^t), \psi(\mathbf{X}) \rangle \right)$$

5. Optimization

1 $\{\mathbf{MLP}_c, \mathbf{MLP}_{\text{SDF}}, \mathbf{MLP}_\psi, \mathbf{MLP}_G, \mathbf{MLP}_J, \mathbf{MLP}_\Delta\}$

2 $\{\omega_e^t, \omega_r^t, \omega_b^t, \omega_b^*\}$

3 pixel embeddings ψ_I

5. Optimization

1 $\{\mathbf{MLP}_{\mathbf{c}}, \mathbf{MLP}_{\text{SDF}}, \mathbf{MLP}_{\psi}, \mathbf{MLP}_{\mathbf{G}}, \mathbf{MLP}_{\mathbf{J}}, \mathbf{MLP}_{\Delta}\}$

2 $\{\omega_e^t, \omega_r^t, \omega_b^t, \omega_b^*\}$

3 pixel embeddings ψ_I

$$\text{Total Loss } \mathcal{L} = \underbrace{\left(\mathcal{L}_{\text{sil}} + \mathcal{L}_{\text{rgb}} + \mathcal{L}_{\text{OF}} \right)}_{\text{reconstruction losses}} + \underbrace{\left(\mathcal{L}_{\text{match}} + \mathcal{L}_{\text{2D-cyc}} \right)}_{\text{feature registration losses}} + \mathcal{L}_{\text{3D-cyc}}.$$

5. Optimization

$$\text{Total Loss } \mathcal{L} = \underbrace{\left(\mathcal{L}_{\text{sil}} + \mathcal{L}_{\text{rgb}} + \mathcal{L}_{\text{OF}} \right)}_{\text{reconstruction losses}} + \underbrace{\left(\mathcal{L}_{\text{match}} + \mathcal{L}_{\text{2D-cyc}} \right)}_{\text{feature registration losses}} + \mathcal{L}_{\text{3D-cyc}}.$$

5. Optimization

$$\text{Total Loss } \mathcal{L} = \underbrace{\left(\mathcal{L}_{\text{sil}} + \mathcal{L}_{\text{rgb}} + \mathcal{L}_{\text{OF}} \right)}_{\text{reconstruction losses}} + \underbrace{\left(\mathcal{L}_{\text{match}} + \mathcal{L}_{\text{2D-cyc}} \right)}_{\text{feature registration losses}} + \mathcal{L}_{\text{3D-cyc}}.$$

$$\mathcal{L}_{\text{rgb}} = \sum_{\mathbf{x}^t} \left\| \mathbf{c}(\mathbf{x}^t) - \hat{\mathbf{c}}(\mathbf{x}^t) \right\|^2, \quad \mathcal{L}_{\text{sil}} = \sum_{\mathbf{x}^t} \left\| \mathbf{o}(\mathbf{x}^t) - \hat{\mathbf{s}}(\mathbf{x}^t) \right\|^2$$

$$\mathcal{L}_{\text{OF}} = \sum_{\mathbf{x}^t, (t, t')} \left\| \mathcal{F}(\mathbf{x}^t, t \rightarrow t') - \hat{\mathcal{F}}(\mathbf{x}^t, t \rightarrow t') \right\|^2, \quad (16)$$

5. Optimization

$$\text{Total Loss } \mathcal{L} = \underbrace{\left(\mathcal{L}_{\text{sil}} + \mathcal{L}_{\text{rgb}} + \mathcal{L}_{\text{OF}} \right)}_{\text{reconstruction losses}} + \underbrace{\left(\mathcal{L}_{\text{match}} + \mathcal{L}_{\text{2D-cyc}} \right)}_{\text{feature registration losses}} + \mathcal{L}_{\text{3D-cyc}}.$$

$$\mathcal{L}_{\text{match}} = \sum_{\mathbf{x}^t} \left\| \hat{\mathbf{X}}^*(\mathbf{x}^t) - \mathbf{X}^*(\mathbf{x}^t) \right\|_2^2, \quad (17)$$

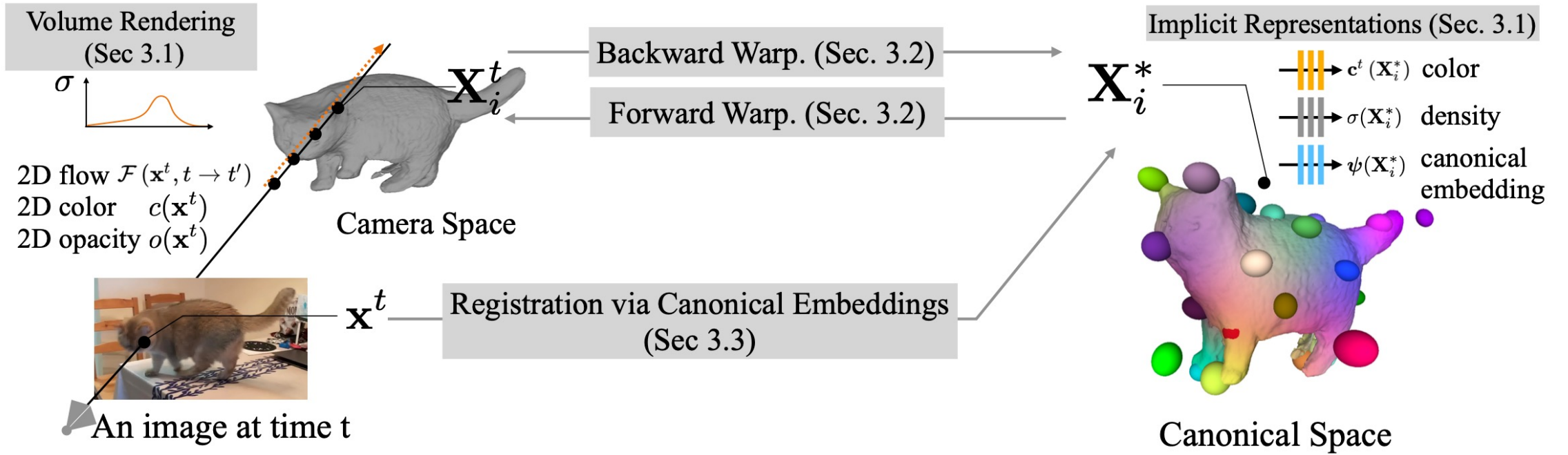
$$\mathcal{L}_{\text{2D-cyc}} = \sum_{\mathbf{x}^t} \left\| \Pi^t \left(\mathcal{W}^{t, \rightarrow}(\hat{\mathbf{X}}^*(\mathbf{x}^t)) \right) - \mathbf{x}^t \right\|_2^2. \quad (18)$$

5. Optimization

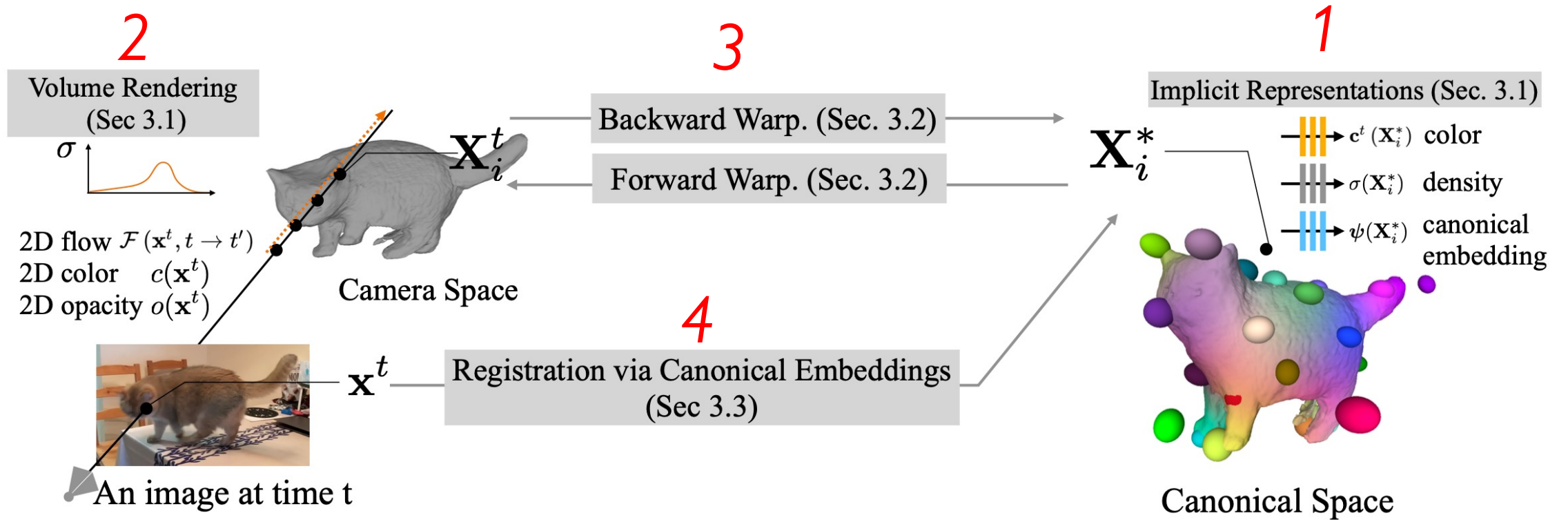
$$\text{Total Loss } \mathcal{L} = \underbrace{\left(\mathcal{L}_{\text{sil}} + \mathcal{L}_{\text{rgb}} + \mathcal{L}_{\text{OF}} \right)}_{\text{reconstruction losses}} + \underbrace{\left(\mathcal{L}_{\text{match}} + \mathcal{L}_{\text{2D-cyc}} \right)}_{\text{feature registration losses}} + \mathcal{L}_{\text{3D-cyc}}.$$

$$\mathcal{L}_{\text{3D-cyc}} = \sum_i \tau_i \left\| \mathcal{W}^{t, \rightarrow} \left(\mathcal{W}^{t, \leftarrow} (\mathbf{x}_i^t) \right) - \mathbf{x}_i^t \right\|_2^2, \quad (19)$$

Method



Method



5: Total Loss $\mathcal{L} = \underbrace{(\mathcal{L}_{\text{sil}} + \mathcal{L}_{\text{rgb}} + \mathcal{L}_{\text{OF}})}_{\text{reconstruction losses}} + \underbrace{(\mathcal{L}_{\text{match}} + \mathcal{L}_{\text{2D-cyc}})}_{\text{feature registration losses}} + \mathcal{L}_{\text{3D-cyc}}$.

Experiments

- Dataset

- Casual –videos dataset
- AMA dataset
- Animated Objects dataset



- Evaluation metrics

- CD
- F-Score

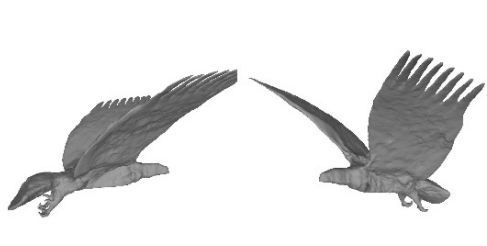
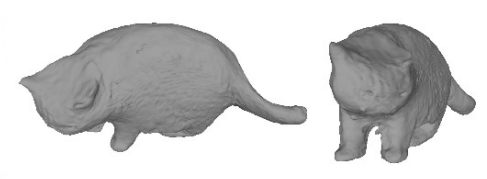
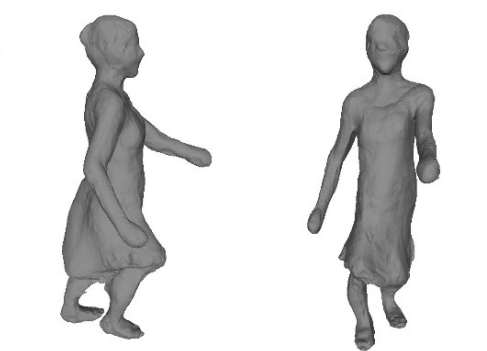


Experiments

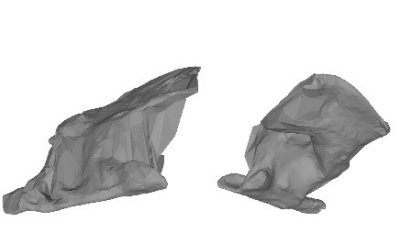
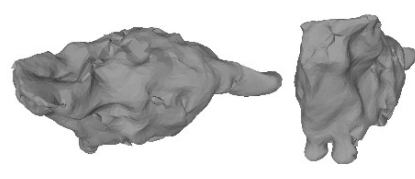
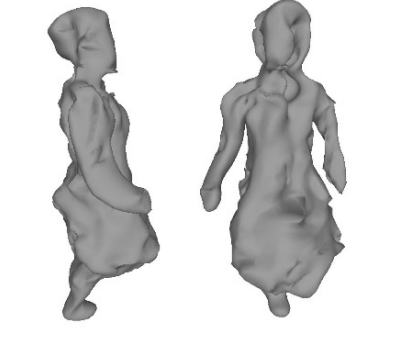
Qualitative comparison of our method with prior work



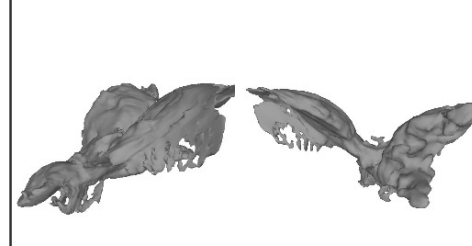
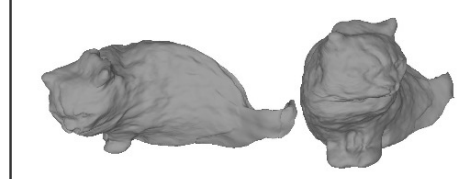
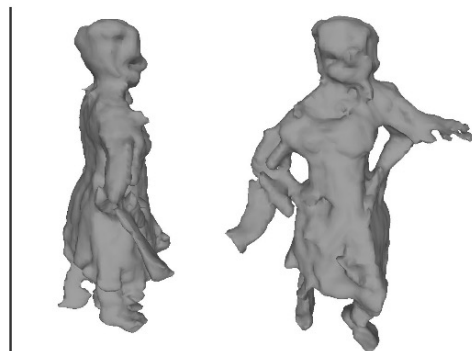
Reference image



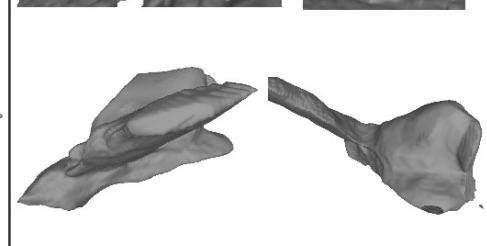
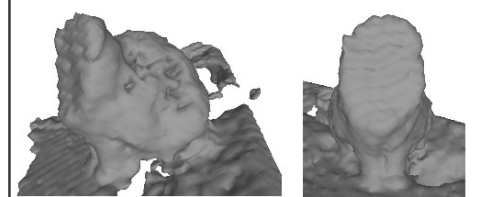
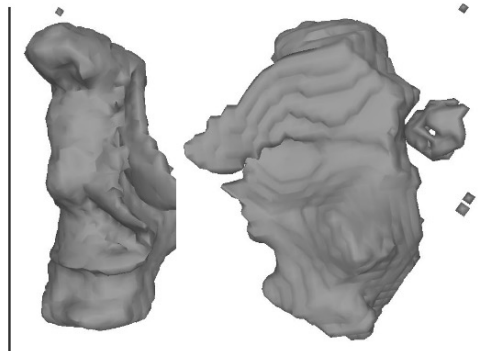
Ours (multi-videos)



ViSER (multi-videos)



Ours (single video)

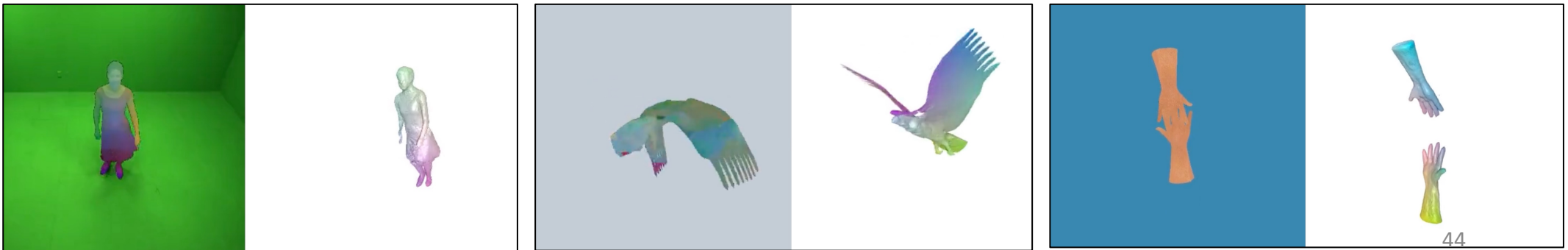


Nerfies (single video)

Experiments

Quantitative results on AMA and Animated Objects

Method	AMA-swing		Eagle*		Hands*	
	CD	F@2%	CD	F@2%	CD	F@2%
Ours	9.1	57.0	8.1	56.7	7.5	49.6
ViSER	15.7	52.2	23.0	20.6	16.8	21.3
Ours ^S	9.4	56.8	10.8	48.6	10.5	35.2
Nerfies ^S	22.6	13.2	18.4	18.0	24.4	14.9



Experiments

Reconstruction Results



25% total iter.



50% total iter.



75% total iter.



100% total iter.

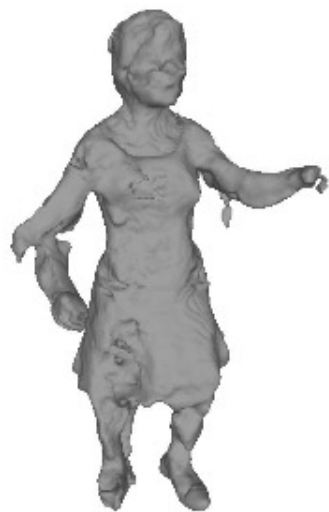
Compliance to topology changes in optimization.

- BANMo incorrectly reconstructs a single rear leg of the dog, but automatically corrects the topology with gradient updates.

Experiments

Diagnostics (Ablation study)

single coherent recon.



Reference

samba recon.

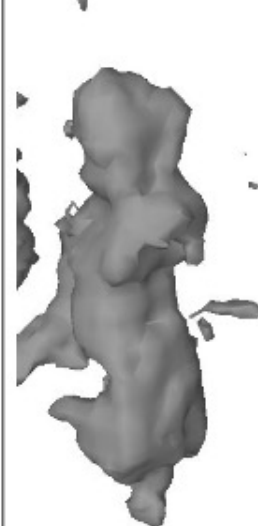


w/o feature registration (Sec. 3.3)

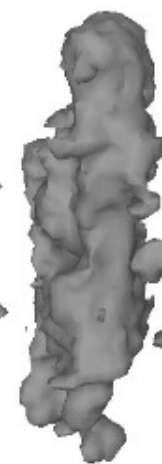
swing recon.



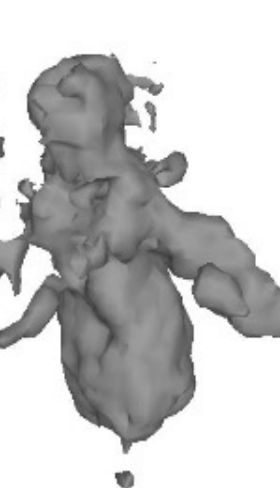
recon. 1



recon. 2



recon. 3



further remove flow loss (Eq. 14)

Experiments

Diagnostics (Deformation modeling)



Reference
image



Neural blend skinning
(Ours)



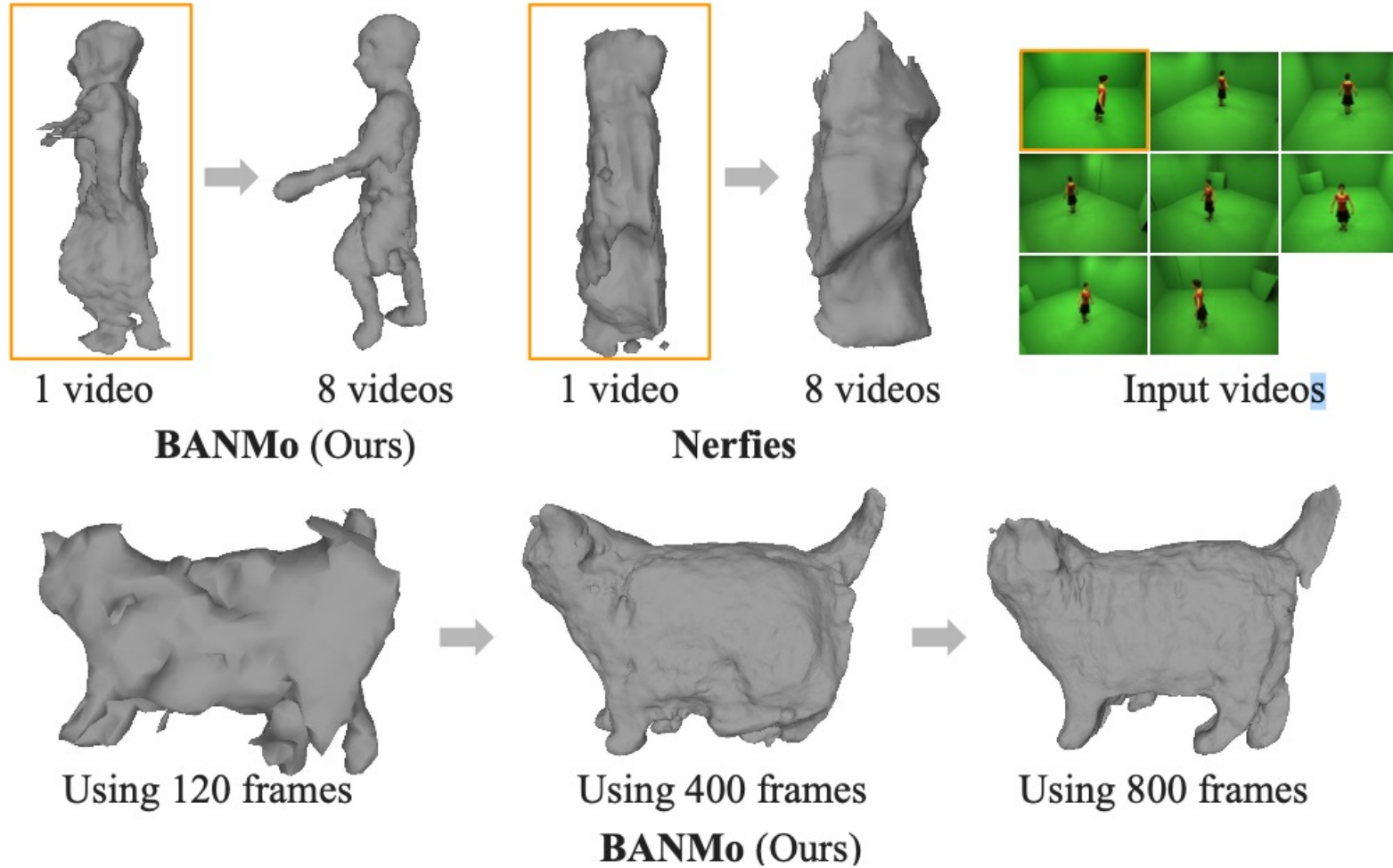
MLP-SE(3)
(Nerfies)



MLP-translation
(NSFF, D-NeRF)

Experiments

Reconstruction completeness vs number of input videos and video frames



Experiments

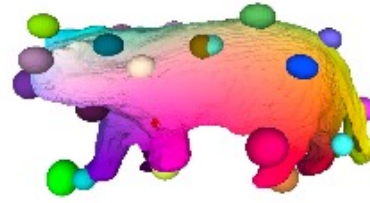
Motion retargeting



Driving frame



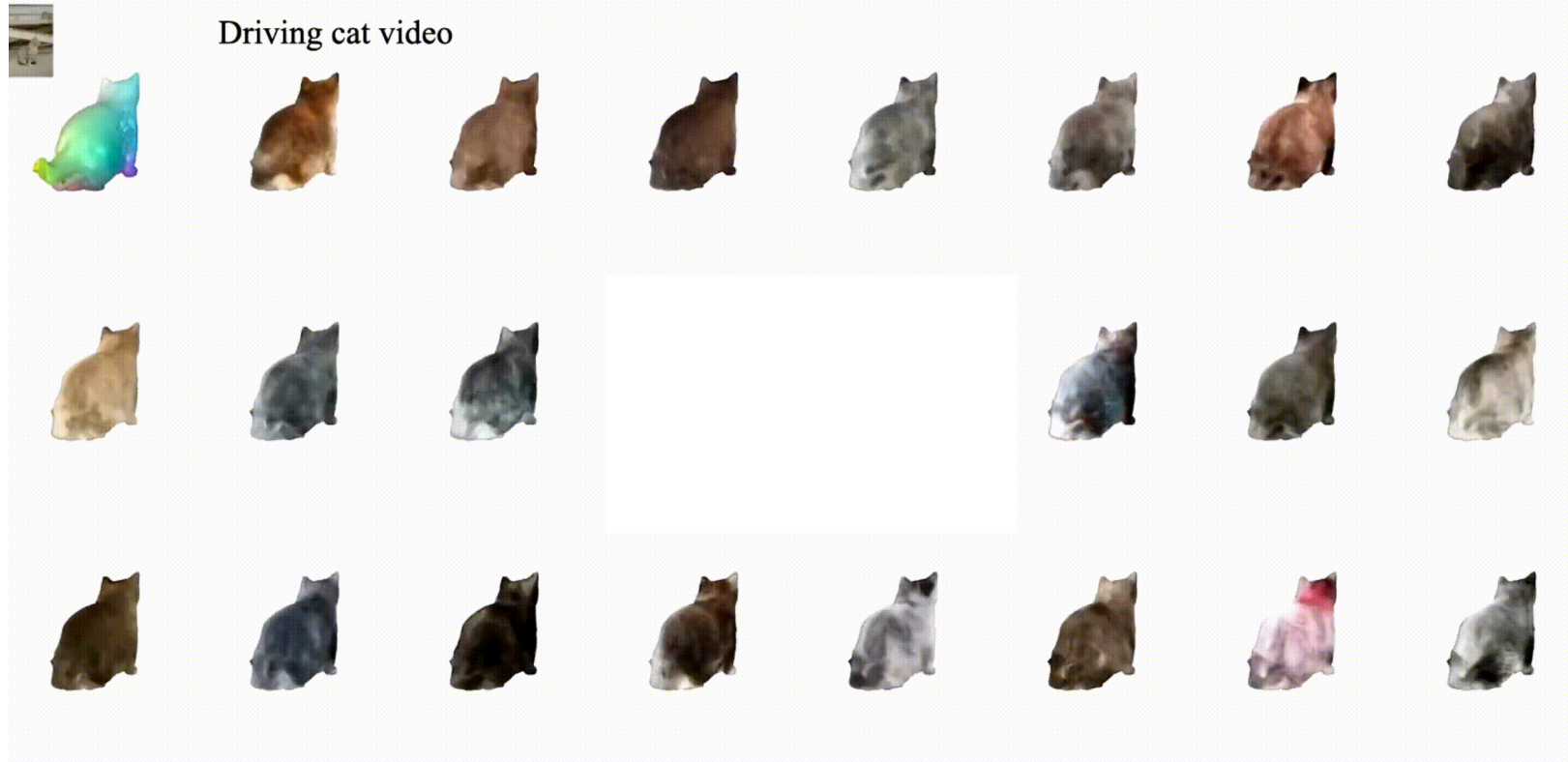
Re-targeted pose rendering



Re-targeted pose



Source model (cat)



Conclusion

