# Unsupervised Learning of Category-Level 3D Pose from Object-Centric Videos

**Leonhard Sommer[1], Artur Jesslen[1], Eddy Ilg[2], Adam Kortylewski[1,3]**

[1]University of Freiburg [2]Saarland University [3]Max Planck Institute for Informatics

universität freiburg

MAX PLANCK INSTITUTE FOR INFORMATICS
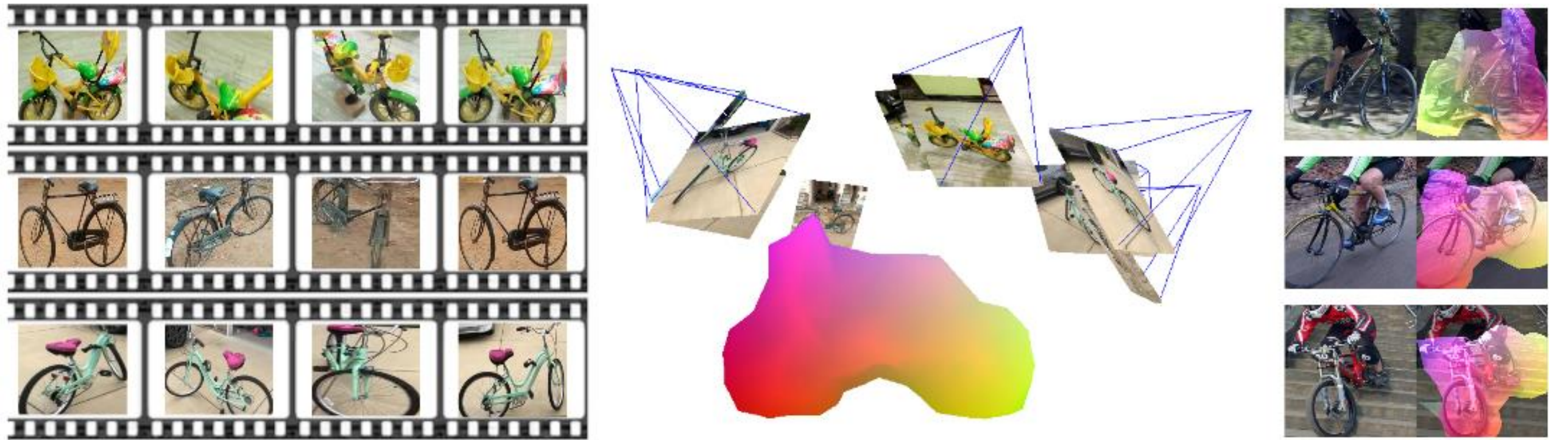
CVPR SEATTLE, WA JUNE 17-21, 2024

## Motivation

➢ **Goal: Unsupervised Learning of Category-Level 3D Pose from Object-Centric Videos**



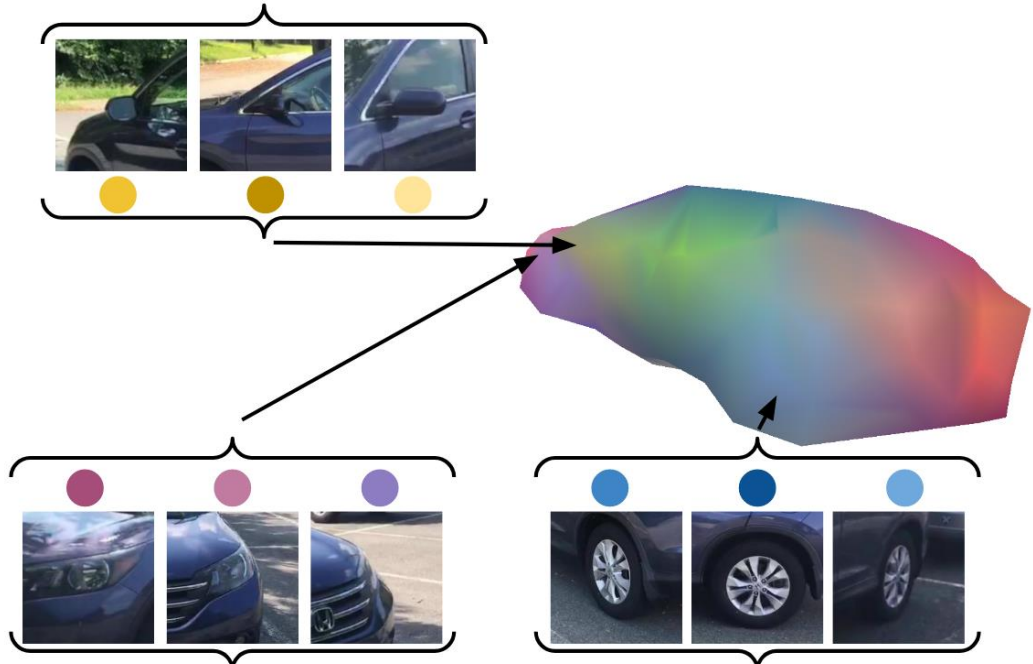Unaligned Object-Centric Videos — Self-Supervised Alignment & Feature Learning — In-the-Wild 3D Pose Estimation
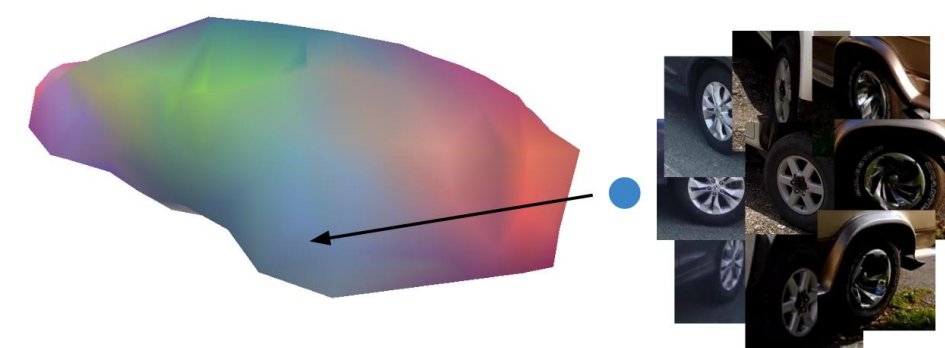
➢ **Challenges**

(1) Structure-from-Motion aligns only camera poses in one video.

→ Novel alignment across videos via self-supervised surface features.

(2) Domain gap between object-centric videos and in-the-wild images.

→ Learning category-level surface features and impose compositionality.

## Our Representations: Surface Features

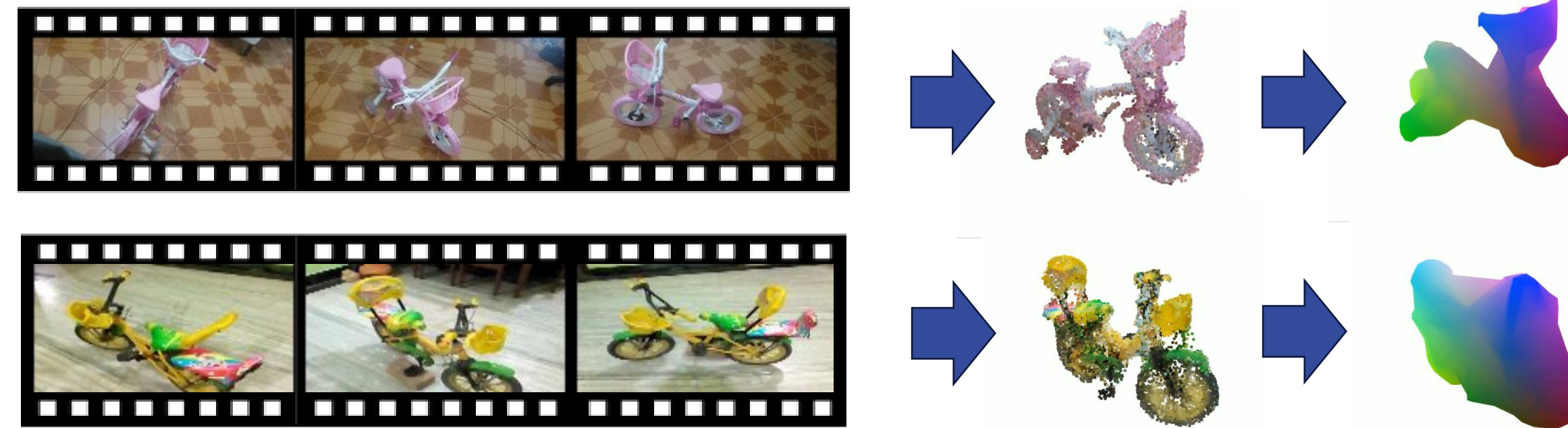### (1) From Single Video

### (2) Category-Level from Many Videos



➢ Viewpoint-dependent features (DINO).

➢ Viewpoint-invariant features facilitate 3D pose estimation.

## (1) Method: Video Alignment

**(1) Obtain surface features representation per video** $S = \{V, F\}$.



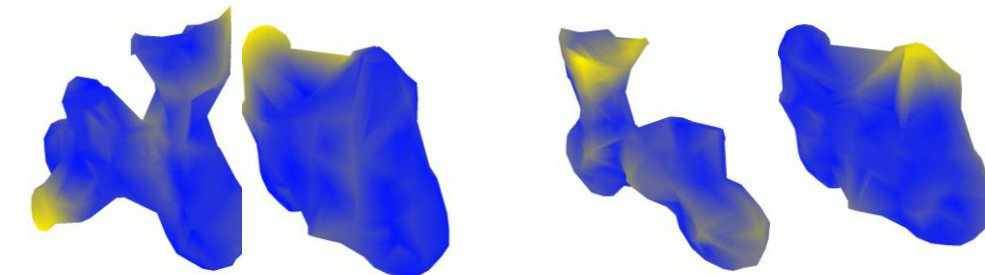**(2) Divide many-to-one alignment into many one-to-one alignments.**

**(3) Minimize one-to-one weighted distance.**

$$\min_{T} \mathcal{D}(S, \bar{S}, T) = (1-\alpha)\mathcal{D}_{geo}(S, \bar{S}, \tau) + \alpha\mathcal{D}_{app}(S, \bar{S}, \tau)$$
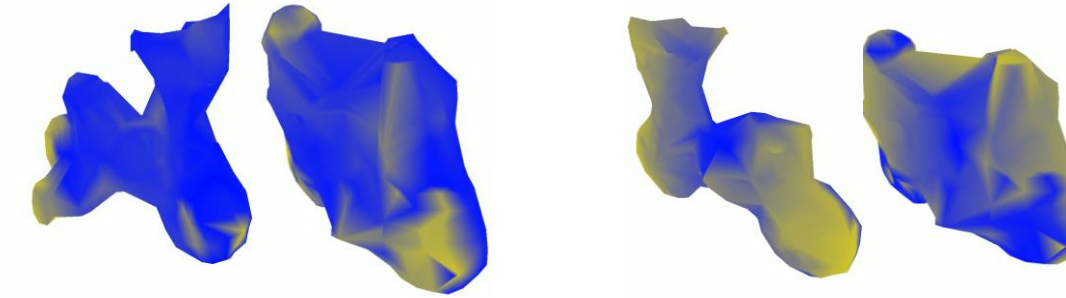
**Geometric Distance**

➢ Weighted Chamfer Distance.

➢ $\mathcal{D}_{geo}(S, \bar{S}, \tau) = \sum_{v_i \in V \cup \bar{V}} \sigma(i, \tau) \, ||v_i - v_{\chi(v_i)}||_2$
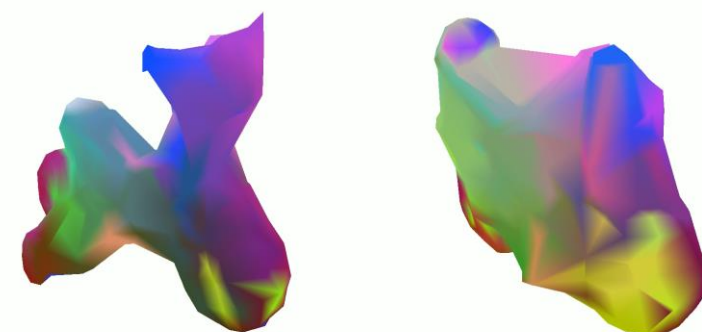
**Appearance Distance**

➢ Weighted appearance correspondences distance.

➢ $\mathcal{D}_{app}(S, \bar{S}, \tau) = \sum_{v_i \in V \cup \bar{V}} \sigma(i, \tau) \, ||v_i - v_{\psi(v_i)}||_2$

**Appearance Correspondences**
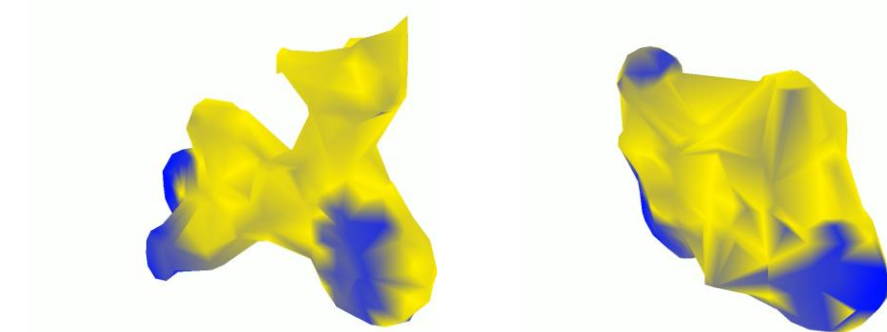
➢ $\psi(v_i, f_i) = \begin{cases} \arg\min_{j \in 1...|\bar{V}|} \min_{k,l} ||f_j^k - f_i^l||, & \text{for } v_i \in V \\ \arg\min_{j \in 1...|V|} \min_{k,l} ||f_j^k - f_i^l||, & \text{for } v_i \in \bar{V} \end{cases}$
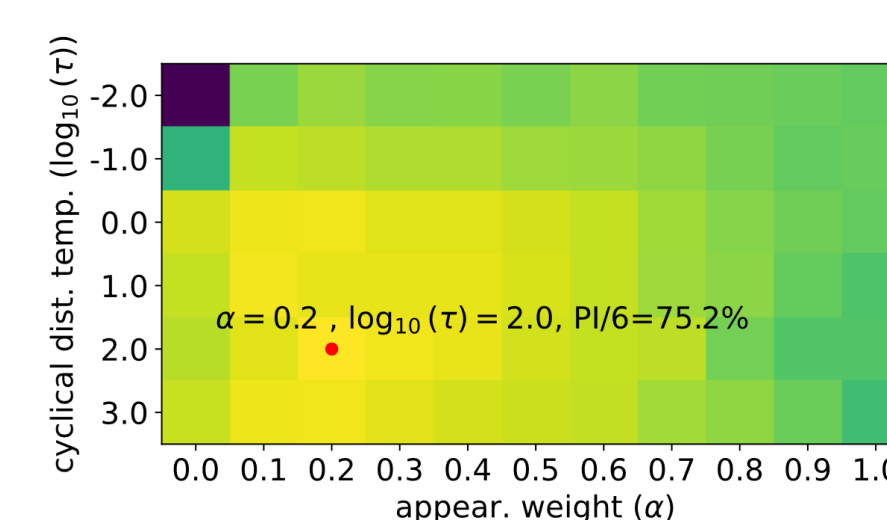
**Cycle-Distance Weighting**

➢ $d_{cycle}(v_i, f_i) = ||v_i - v_{\psi(v_j, f_j)}||_2$, with $j = \psi(v_i, f_i)$

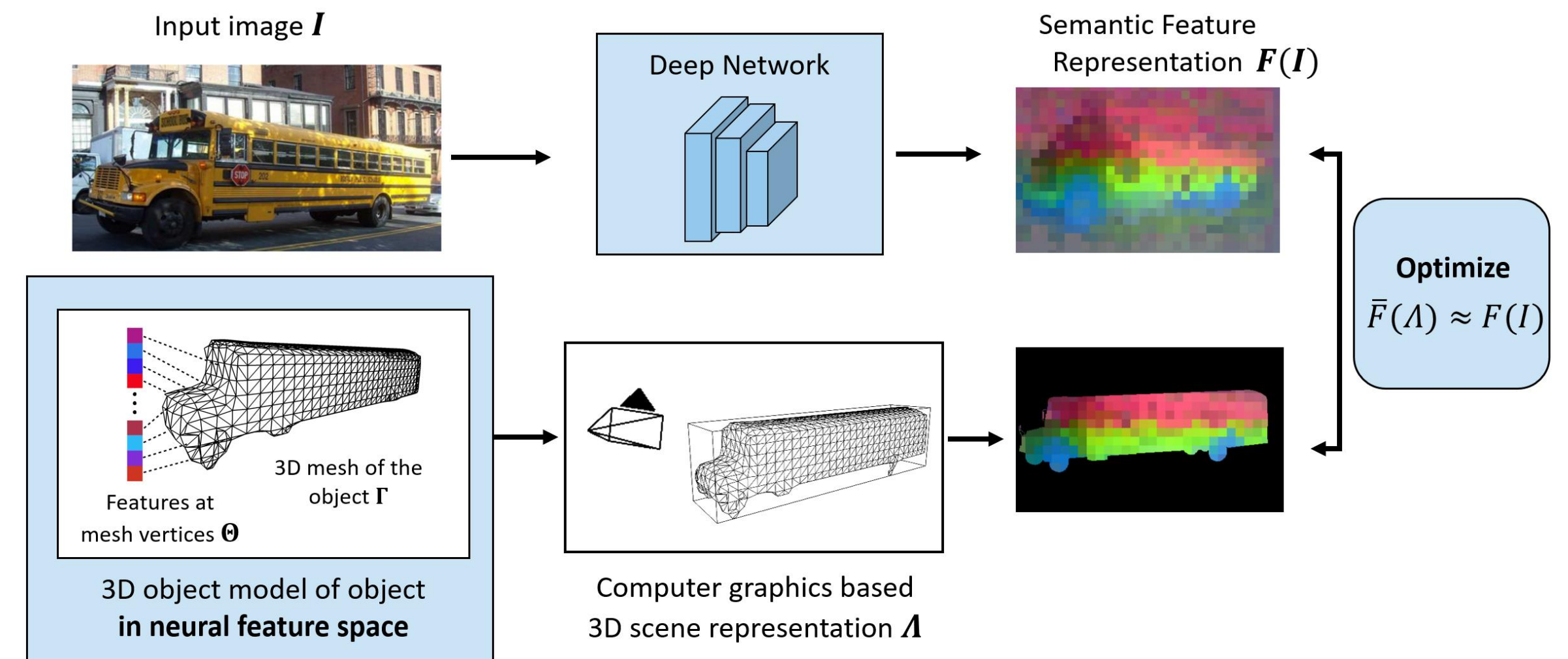➢ $\sigma(i, \tau) = \text{Softmax}\left(-\frac{d_{cycle}(v_i, f_i)}{\tau \max_{v_i, v_j \in V} ||v_i - v_j||_2}\right)$

**Optimum**

➢ Geometric correspondences weighted four times as much as the appearance correspondences.
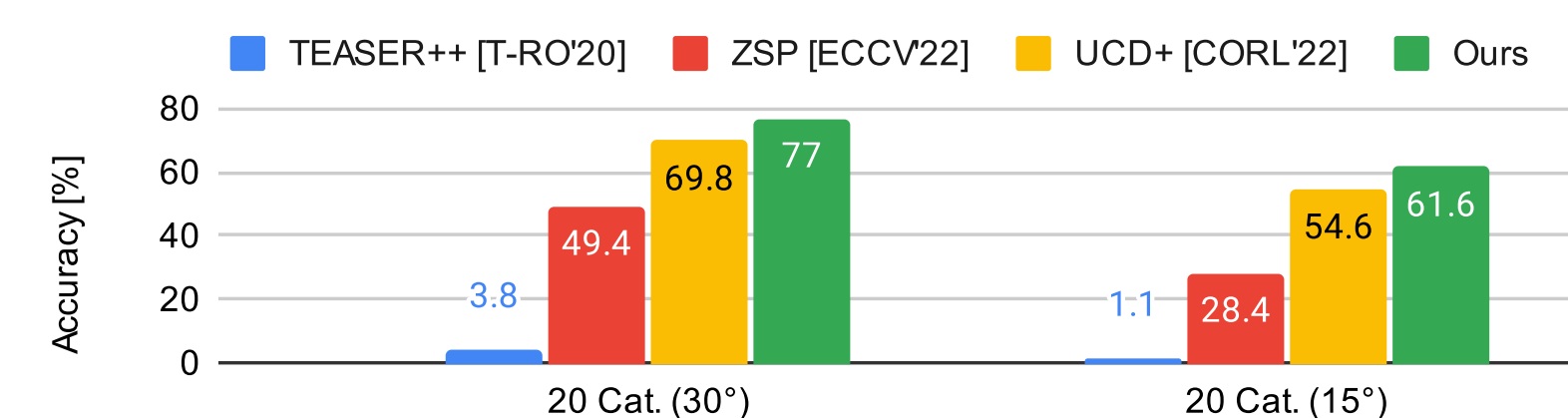
➢ Only weak cycle-distance weighting.



$\alpha = 0.2$, $\log_{10}(\tau) = 2.0$, PI/6=75.2%

appear. weight ($\alpha$)
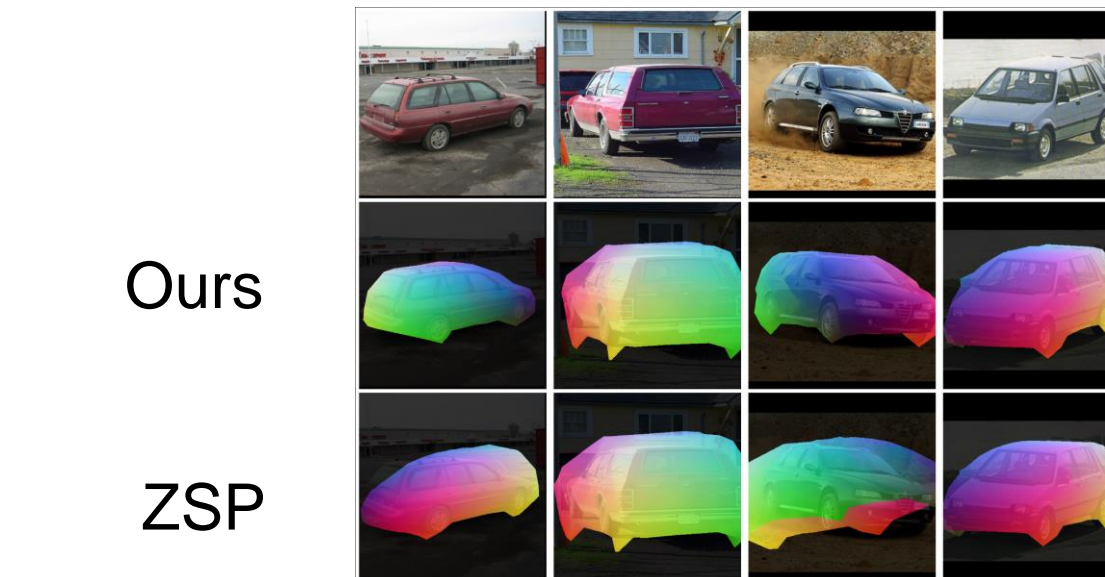
## (2) Method: Learning Category-Level Surface Features



Input image $I$ — Deep Network — Semantic Feature Representation $F(I)$

Optimize $\bar{F}(\Lambda) \approx F(I)$

Features at mesh vertices $\Theta$ — 3D object model of object **in neural feature space** — 3D mesh of the object $\Gamma$ — Computer graphics based 3D scene representation $\Lambda$

## Results

### Video Alignment – CO3D



Legend: TEASER++ [T-RO'20], ZSP [ECCV22], UCD+ [CORL'22], Ours

20 Cat. (30°): 3.8, 49.4, 69.8, 77
20 Cat. (15°): 1:1, 28.4, 54.6, 61.6

### In-the-Wild 3D Pose Estimation – PASCAL3D+



Ours / ZSP

Legend: ZSP [ECCV22], Ours, VoGE [ICLR'23] (supervised)

Mean (7 Cat.): 46, 69.2, 90.93

Legend: ZSP [ECCV22], Ours (Regression), Ours, VoGE [ICLR'23] (supervised)

| | 5 Videos | 10 Videos | 20 Videos | 50 Videos | >10K Images with 3D Pose Labels |
|---|---|---|---|---|---|
| ZSP | 46.00% | 46.00% | 46.00% | 46.00% | |
| Ours (Regression) | 45.90% | 52.50% | 60.10% | 61.80% | |
| Ours | 56.30% | 58.20% | 65.20% | 69.20% | |
| VoGE | | | | | 90.93% |