

Common Pets in 3D: Dynamic New-View Synthesis of Real-Life Deformable Categories

Supplementary material

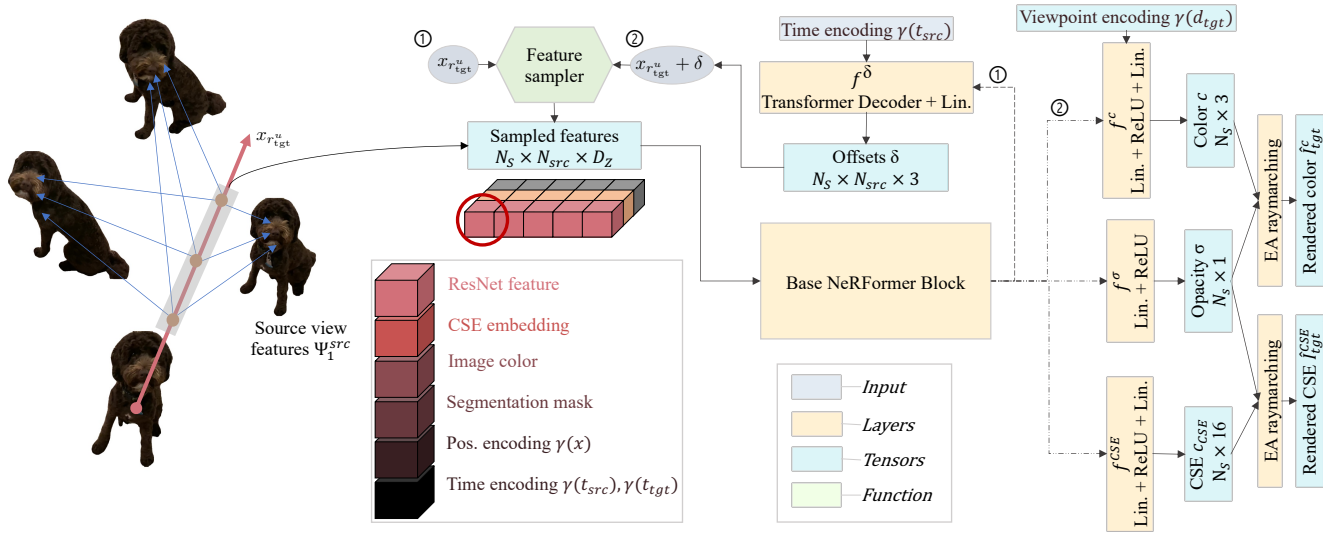


Figure I. **The architecture of Tracker-NeRF** First, source image features are bilinearly sampled given the points on the target ray \mathbf{r}_{tgt}^u . The sampled features are then processed with the Base NeRFormer Block (with the same architecture as in [38]) to generate intermediate features. The latter enters the offset prediction head \mathcal{D}_{TR} that generates per-point offsets δ (① in the figure). The source images are then sampled again at the locations of projected offset-adjusted ray-points. The resampled features enter Base NeRFormer Block again to predict a new intermediate feature grid which enters 3 final heads that predict colour, opacity, and CSE embedding for each ray point (② in the figure). The final color and CSE render is formed by Emission-Absorption ray marching over the predicted opacities, colors, and CSE embeddings of the ray-points.

A. Network architecture

The network architecture is summarized in Figure I. The sampled grid of TWCE encodings $Z_{TWCE}^{\mathbf{r}_u}$, or the adjusted tokens $\tilde{Z}_{TWCE}^{\mathbf{r}_u}$, are processed with Base NeRFormer Block which has the same architecture as the vanilla NeRFormer from [39]. The output of the NeRFormer block is a set of intermediate features that are converted to either: 1) scene flow δ with the offset predictor \mathcal{D}_{TR} during the first pass or, 2) are converted to colors \mathbf{c} , opacities σ , or CSE embeddings \mathbf{C} with the final head f_{TR}^l during the second pass. The offset predictor \mathcal{D}_{TR} is implemented with a single-layer transformer decoder block that takes as input the intermediate features from the base NeRFormer block together with the encoding of the source time of each sampled point $\gamma(t^{src})$, and outputs the per-point offsets.

B. Mask annotations for CoP3D

To generate the object masks M_i^j , we input each frame I_i^j to the PointRend semantic segmentation network [16],

extracting N_i^j candidate masks $\hat{M}_i^j \in [0, 1]^{N_i^j \times H \times W}$. To track a single foreground object across the video, we use a Hidden Markov Model to generate the set $\{M_i^t | M_i^t \in [0, 1]^{H \times W}\}_{j=1}^{N_{V_i}}$ containing a single mask per frame, using intersection-over-union (IoU) of masks' bounding boxes as a pairwise potential.