

# Lift3D: Zero-Shot Lifting of Any 2D Vision Model to 3D

Mukund Varma T<sup>1</sup>, Peihao Wang<sup>2</sup>, Zhiwen Fan<sup>2</sup>, Zhangyang Wang<sup>2</sup>, Hao Su<sup>1</sup>, Ravi Ramamoorthi<sup>1</sup>  
<sup>1</sup>University of California San Diego, <sup>2</sup> University of Texas at Austin

{tmukund, haosu, ravir}@ucsd.edu, {peihaowang, zhiwenfan, atlaswang}@utexas.edu

[mukundvarmat.github.io/Lift3D/](https://mukundvarmat.github.io/Lift3D/)

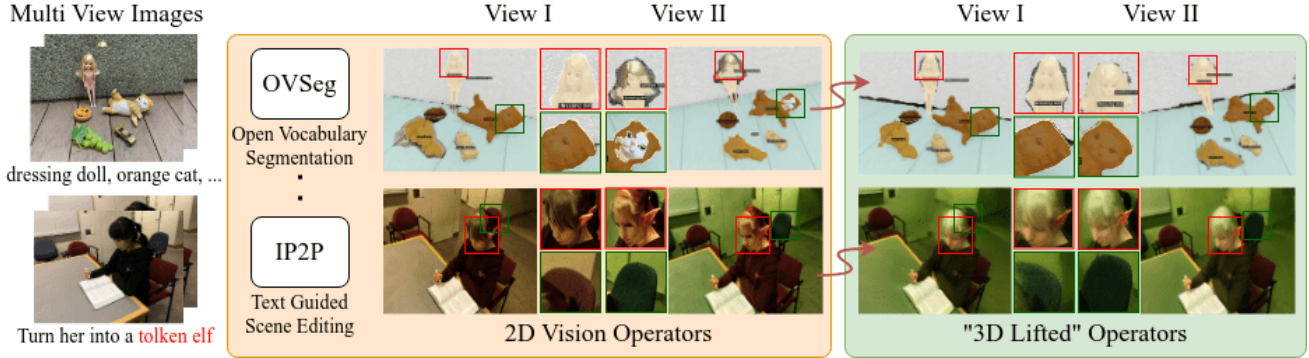


Figure 1. Imagine we are using a 2D vision operator, such as semantic segmentation or scene editing, on multiple-view input images. This often leads to inconsistent predictions across different views (as shown in the middle column). To address this, we introduce *Lift3D*, a framework designed to transform these inconsistent 2D outputs into view-consistent 3D predictions (illustrated in the right column). Our approach is both scene and operator-agnostic, meaning it can adapt to any downstream task or scene without additional adjustments. We demonstrate how Lift3D effectively resolves inconsistencies in multi-view predictions across open vocabulary segmentation and text-driven scene editing. Notice the color discrepancies in the same rightmost chair across two views (varying from reddish to greenish) in the 2D results at the bottom row, and the inconsistencies in facial and hair color. For a clearer comparison between the 2D and 3D outcomes, we recommend zooming into the electronic version of this image.

## Abstract

In recent years, there has been an explosion of 2D vision models for numerous tasks such as semantic segmentation, style transfer or scene editing, enabled by large-scale 2D image datasets. At the same time, there has been renewed interest in 3D scene representations such as neural radiance fields from multi-view images. However, the availability of 3D or multiview data is still substantially limited compared to 2D image datasets, making extending 2D vision models to 3D data highly desirable but also very challenging. Indeed, extending a single 2D vision operator like scene editing to 3D typically requires a highly creative method specialized to that task and often requires per-scene optimization. In this paper, we ask the question of whether **any** 2D vision model can be lifted to make 3D consistent predictions. We answer this question in the affirmative; our new Lift3D method trains to predict unseen views on feature spaces generated by a few visual models (i.e. DINO and CLIP), but then generalizes to novel vision operators and tasks, such as style transfer, super-resolution, open vocabulary segmentation and image colorization; for some of these tasks, there is no comparable previous 3D method. In many cases, we even outperform state-of-the-art methods specialized for the task in question.

Moreover, *Lift3D* is a zero-shot method, in the sense that it requires no task-specific training, nor scene-specific optimization.

## 1. Introduction

Recent progress in 2D image understanding has been extraordinary, driven by the assembly of extensive image datasets with intricate labels, and the innovation of varied network architectures. This has led to remarkable advances in diverse tasks such as semantic segmentation [24, 28], style transfer [66], scene editing [3], and super-resolution [62]. Yet, the domain of 3D understanding, essential for sectors like autonomous driving, robotics, and 3D asset creation, lags in the development of versatile and robust neural networks. A common requirement for 3D understanding is the processing of multi-view images, which is hindered by the lack of expansive, well-labeled multi-view image datasets. This limitation raises a critical question: Is it possible to modify existing neural networks, initially intended for single-image analysis, to accommodate multi-view inputs, and in doing so, eliminate the inconsistencies typically encountered when applying 2D operators to each view individually?

In this paper, we set an ambitious goal to universally transfer an *arbitrary* pre-trained 2D feature backbone to a 3D model or vision operator *on the fly*, which produces

*view-consistent* predictions (see Fig. 1) from any arbitrary viewing angle given a set of posed 2D images and their 2D predictions. We observe that the intermediate feature maps are roughly aligned with the input image for modern strong 2D operators. Drawing inspiration from label propagation algorithms [21, 53, 67, 70, 71], for us to produce smooth predictions across views, we only need to rectify inconsistencies and propagate labels from supporting views to novel views.

Motivated by these observations, we propose a novel algorithm *Lift3D* which learns to fix and propagate inconsistent multi-view network outputs to view-consistent predictions. Our architecture design is underpinned by image-based rendering [9, 16, 27], where modern methods essentially learn to aggregate pixels with epipolar constraints to synthesize novel views [45, 46, 49, 55, 61]. By viewing dense features as colors, our method is trained to interpolate novel views on a feature space generated by a pre-trained 2D visual model or operator, incorporating both RGB and feature maps from adjacent views.

Following the widespread success of Neural Radiance Fields or NeRF for view synthesis [35, 49, 55], Lift3D casts a ray for each pixel on the target image plane, samples, and projects points to nearby views to fetch RGB and feature values of the epipolar correspondences. The aggregated RGB information will be used to infer geometry and appearance properties involved in volume rendering [33, 35]. This geometric information enables the interpolation of features from the pre-trained 2D vision model to estimate a target feature map, that is both multi-view consistent and decodable to generate the desired output (Fig 2, Sec. 4).

The whole pipeline is end-to-end differentiable and can be supervised by both ground-truth color and feature maps exported from 2D visual models. An impressive property of Lift3D is that, after we have trained Lift3D on only a few vision operators/models (DINO and CLIP), we discovered that *Lift3D has strong zero-shot ability, enabling any 2D vision operator to be lifted to 3D without any scene-specific or operator-specific training*. On a variety of 3D vision tasks such as semantic segmentation, style transfer, super-resolution, and text-driven editing of 3D scenes, our performance is comparable and sometimes even better than methods that use per-scene optimization and/or have been designed for the specific task (Fig. 1, Sec. 5). In some other tasks like image colorization and open vocabulary segmentation, we present the first 3D extension of 2D vision operators.

In summary, by formulating feature propagation as a rendering process, Lift3D becomes agnostic to 2D models and task domains, demonstrating versatility to generalize across various feature backbones. We are unaware of such abilities reported by other works, though there exist several individual works that have dealt with the extension of each separate 2D vision operator to 3D. The discovery further supports the practical impact of our method, i.e., as long as there exists a 2D feature backbone for a given task, Lift3D helps realize the same in 3D without any extra fine-tuning.

## 2. Related Work

**Progress in 2D Vision Models.** Deep Neural Networks have achieved wide success in a variety of computer vision tasks, ranging from classical object recognition tasks like classification, and segmentation [10, 24] to more complicated tasks including image generation, editing, etc [39, 40, 42, 43]. The 2D image domain boasts of several large, high-quality image datasets [37, 41, 44]. The recent breakthroughs in natural language processing for large-scale model pre-training have opened the way for similar foundation models in 2D computer vision [4, 37, 50], that generate visual features at an image and/or pixel level for several downstream applications. Given the plethora of 2D models, it is highly desirable to lift these 2D vision operators and tasks to 3D, enabling them to work in a 3D consistent way from multi-view images.

**Attempts in 3D Vision Models.** Some recent efforts [18, 59, 69] have attempted to build 3D foundational models by mirroring the successful approach of 2D foundation models. However, this direction is significantly limited by the training data available in 3D. Another line of work tries to lift a pre-trained 2D foundation model into 3D representations (e.g. point clouds or NeRFs [35]) for 3D scene segmentation [13, 31, 32] and editing [12, 17, 51, 64]. A major drawback of these approaches is their computational inefficiency as they often involve a per-scene optimization process. Solutions proposed by [5, 6, 29, 60] distill a 2D vision model into a generalizable rendering pipeline [49, 55], that allows for zero-shot inference across different scenes. Although promising results have been demonstrated on scene segmentation, it remains elusive to generalize these methods to other tasks.

**Novel View Synthesis.** A pioneering work, Neural Radiance Fields (NeRF) [35] synthesizes photo-realistic and consistent novel views by fitting each scene as a continuous 5D radiance field (3D coordinates and 2D viewing directions) parameterized by an MLP. Several follow-up works have further improved NeRF’s rendering quality [1, 2, 54, 57]. However, these methods are rarely applied to general 3D applications beyond novel view synthesis, primarily due to the absence of labeled multi-view data. While some methods have been successful in realizing general computer vision tasks in 3D [17, 30, 31, 51], they still require a substantial amount of creativity to tailor NeRFs for each task and even need scene-specific optimization.

## 3. Preliminary: Generalizable View Synthesis

In this work, we extend the pipeline of Generalizable Novel View Synthesis (GNVS) to render 3D consistent feature maps. Below we equip readers with the necessary background. Given  $N$  calibrated input (or source) views with known pose information  $\{\mathbf{I}_i, \mathbf{P}_i\}_{i=1}^N$ , the goal of GNVS is to synthesize target novel view  $\mathbf{I}_T$ , even for scenes not observed during training, thereby achieving generalizability. These methods first extract deep convolutional features  $\mathbf{F}_i = \mathcal{F}_{\text{conv}}(\mathbf{I}_i)$  for each input view. To render a target view,



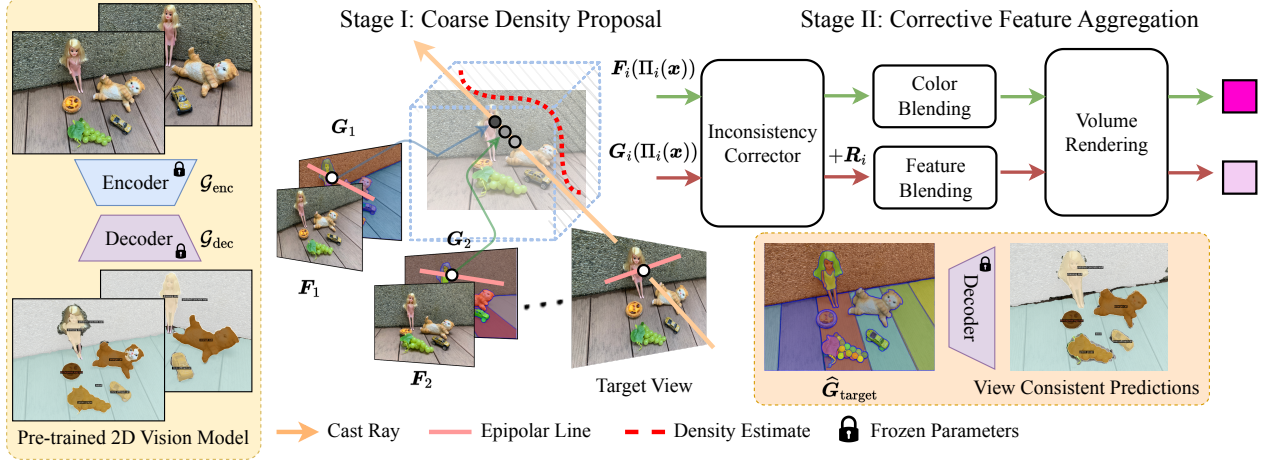


Figure 2. Overview of *Lift3D*: 1) Given multi-view images of a scene, we first extract an intermediate feature map using an encoder for each view independently, 2) Using the source view features, we estimate the target view feature map via an extended generalizable novel view synthesis pipeline that learns to correct feature space inconsistency and fuse information according to 3D geometry priors, and 3) Directly use the decoder from the pre-trained visual model to process the estimated feature map and synthesize the final prediction for downstream tasks.

several rays are cast into the scene, and  $K$  points  $\{x_j\}_{j=1}^K$  as shown in Fig. 2, are sampled along each ray. Each point is then projected onto each source view  $I_i$  using projection function  $\Pi_i$ , and the nearest feature in the image plane is queried. These multi-view features are aggregated into point feature  $f(x)$  (We drop the index  $j$  for simplicity) as follows:

$$f(x) = \mathcal{F}_{\text{view}}(\{F_i(\Pi_i(x))\}_{i=1}^N) \quad (1)$$

where  $\Pi_i(x)$  projects  $x$  onto  $I_i$ ,  $F_i(\Pi_i(x))$  fetches the feature vector at the epipolar projections for input view  $i$ , and  $\mathcal{F}_{\text{view}}$  is a permutation-invariant aggregation mapping that learns to be occlusion-aware to combine epipolar features [49]. In IBNet [55],  $\mathcal{F}_{\text{view}}$  is implemented by estimating blending weights  $w_i$  via a DeepSets-like architecture [63] to fuse features from different views and produce point feature  $f(x) = \sum_i w_i F_i(\Pi_i(x))$ . Then the alpha value  $\alpha(x)$  and color  $c(x)$  can be decoded from feature  $f(x)$  by another network [55, 61]. Finally, using volume rendering [35], the point-wise individual color and density values are composed to estimate ray color.

## 4. Method

**Overview.** We introduce *Lift3D* that applies the effect of any 2D visual models to arbitrary 3D scenes, and illustrate the entire pipeline in Fig. 2. Formally, given multi-view images of a scene and a pretrained 2D visual model  $\mathcal{G}_{2D}$ , our goal is to generate *view-consistent* predictions from an arbitrary angle, as if the 2D visual model is transferred to be a 3D model. We assume our 2D visual model follows an encoder-decoder structure<sup>1</sup>  $\mathcal{G}_{2D} := \mathcal{G}_{\text{dec}} \circ \mathcal{G}_{\text{enc}}$ , where the encoder stage  $\mathcal{G}_{\text{enc}}$  maps the input image to a latent representation and the decoder stage  $\mathcal{G}_{\text{dec}}$  transforms the latent features to

<sup>1</sup>The 2D visual model need not necessarily follow an encoder-decoder design. We simply refer to the first few layers of the model as  $\mathcal{G}_{\text{enc}}$  and the remaining as  $\mathcal{G}_{\text{dec}}$ .

the desired output space. We observe that the intermediate latent representation, denoted as  $G_i = \mathcal{G}_{\text{enc}}(I_i)$ , is robust and spatially aligned with the input image  $I_i$ , despite having inconsistencies across multiple views that need to be fixed. Our method builds upon existing GNVS techniques [55], i.e. by viewing dense features as colors, we learn to interpolate novel views on a feature space generated by the 2D pretrained vision model, by incorporating both RGB and feature maps from input views (Sec. 4.1). However, naively doing this does not work well since per-view features by 2D models are inherently noisy and inconsistent. Rather, our method learns to leverage consistency information on RGB maps to rectify inconsistent artifacts fetched from the feature maps while performing view aggregation (Sec. 4.2). By doing so, our method generates the now, view-consistent feature map for the target view that can be directly decoded using  $\mathcal{G}_{\text{dec}}$  to synthesize the final prediction i.e. the same task the input pretrained 2D vision model was intended for. Such a pipeline can be pre-trained on very few 2D vision models and then directly applied on unseen scenes and 2D vision operators during inference (Sec. 4.3). Below, we provide more details on the individual components.

### 4.1. Generalizable Feature Rendering

Existing GNVS techniques [46, 49, 55] only render target view color, but it is possible to easily extend them to other quantities of interest, e.g. in this case a high-dimensional feature vector  $G_i$ . A straightforward approach is to modify Eq. 1 by replacing  $F_i(\Pi_i(x))$  with  $G_i(\Pi_i(x))$ , i.e. to fetch the epipolar projections on the per-view intermediate features of the 2D vision model and aggregate them across multiple views to obtain  $g(x)$ . But such a pipeline cannot generalize to any unseen 2D vision operator [60], as the input intermediate features differ greatly across different 2D vision models. Moreover, the absence of ground truth features from arbitrary viewing angles that are multi-view

consistent makes training supervision non-trivial. However, we have access to RGB information across multiple source views that are view-consistent and can be used to guide the feature correction process.

Motivated by this observation, we simply share the aggregation weights between the epipolar RGB features  $F_i(\Pi_i(x))$  and similarly constructed epipolar projections on the encoded features from our 2D visual model, denoted by  $G_i(\Pi_i(x))$ . After we obtain the volumetric representations from RGB features denoted by  $f(x)$  and from the pretrained vision model  $g(x)$ . Formally, we rewrite Eq. 1 as:

$$\begin{aligned}\mathcal{F}_{\text{view}} &:= f(x) = \sum_i w_i F_i(\Pi_i(x)) \\ g(x) &= \sum_i w_i G_i(\Pi_i(x)) \\ \text{where } w_i &= \mathcal{F}_w(\{F_i(\Pi_i(x))\}_{i=1}^N)\end{aligned}\quad (2)$$

where  $\mathcal{F}_w$  denotes a fully connected layer that transforms epipolar RGB feature  $F_i(\Pi_i(x))$  to blending weight  $w_i$ . Next,  $f(x)$  will be decoded as pointwise color and density to predict novel view image  $\hat{I}_{\text{target}}$ , while  $g(x)$  is accumulated along the ray using the estimated density to obtain the feature prediction  $\hat{G}_{\text{target}}$ . In Fig. 2, we illustrate the epipolar aggregation and ray marching as blending and volume rendering operations respectively. By viewing feature rendering as a view interpolation task grounded by RGB input, we ensure generalization to several unseen 2D vision models during inference.

## 4.2. Corrective Feature Aggregation

In the previous subsection, we assume the view-consistent RGB feature maps are roughly aligned with the per-view feature derived from the 2D vision model. However, in practice this is not the case and several inconsistencies cause significant noise in the rendered feature map. In Fig. 3 we visualize a three-channel PCA map of the estimated features after interpolation and compare it against a 2D feature encoding of the target view. We can clearly see that the estimated novel view features are noisy, therefore inconsistent and do not accurately represent the scene geometry, making them undesirable for several downstream tasks.

As a solution, we propose a two-stage aggregation strategy that first performs a correction on the epipolar features obtained from the 2D visual model, before the remaining aggregation steps. The correction factor can be obtained from the RGB features, particularly near the points sampled around the actual 3D location of the object intersecting the cast ray, where multi-view consistency of the epipolar projections is expected. We incorporate a coarse-fine stage like NeRF [35] i.e. perform view synthesis [55] using points randomly sampled along the cast ray, and then perform importance sampling using the obtained coarse density function to specifically fetch projections of points near the actual 3D location of the object of interest. The per-view projections on both the RGB and encoder latent feature planes are roughly

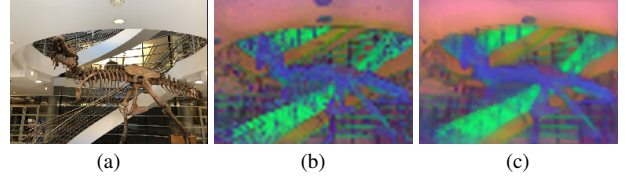


Figure 3. (a) Target View, (b) We visualize the PCA (3-channel) map of the interpolated feature maps from our baseline that naively shares blending and ray marching weights between color rendering and 2D encoder features. (c) A similar PCA map is estimated by directly encoding the target view with the 2D encoder. We can see that naively sharing weights without accounting for the inconsistencies in the source view features results in noisy feature estimation, undesirable for several fine-grained downstream tasks.

expected to be view-consistent since they lie on the same object, and can be used for the feature-lifting process. Therefore, we compute a residual correction term for each view  $R_i$  as:

$$\begin{aligned}\Delta_i &= F_i(\Pi_i(x)) - \mathcal{F}_1(G_i(\Pi_i(x))) \\ R_i &= \mathcal{F}_2(\text{pool}_i(\{G_i(\Pi_i(x))\}_{i=1}^N) \oplus \Delta_i) \\ \tilde{G}_i(\Pi_i(x)) &= G_i(\Pi_i(x)) + R_i\end{aligned}\quad (3)$$

where  $x$  denotes importance sampled points,  $\mathcal{F}_1$  projects  $G_i(\Pi_i(x))$  into the same latent dimension as  $F_i(\Pi_i(x))$ ,  $\oplus$  denotes the concatenation operation,  $\text{pool}_i$  is a max pooling operation applied across input views and  $\mathcal{F}_2$  represents a fully connected layer that unprojects back to the original dimension size. Using these now corrected features, we continue with the aggregation steps detailed in Eq. 2 by replacing  $G_i(\Pi_i(x))$  with  $\tilde{G}_i(\Pi_i(x))$  to synthesize novel view feature (see Fig. 2). We ablate on the individual components of our method and include a discussion in the supplementary (see Sec. A).

## 4.3. Training and Inference

The entire pipeline (both the coarse density proposal, and feature renderer networks) can be end-to-end supervised across multiple scenes with their multi-view images and the affiliated 2D encoded feature maps.

$$\mathcal{L} = \|\hat{I}_{\text{target}} - I_{\text{target}}\|_2^2 + \|\hat{G}_{\text{target}} - G_{\text{target}}\|_2^2 \quad (4)$$

where  $I_{\text{target}}$  and  $G_{\text{target}}$  denote the target ground truth view and its corresponding feature map encoded by a 2D vision model that is not strictly view-consistent. However, since the underlying scene geometry is constrained by view-consistent RGB input views, we expect similar propagation in the rendered feature map from the target view. Moreover, the simple training objective can cancel out random inconsistent patterns, leading to smooth and precise prediction, with a mechanism similar to [26]. More significantly, we find that our model can be pre-trained with only very few (in practice 2-3) 2D visual models and then directly applied to unseen scenes and 2D vision operators during inference. Given an unseen 2D vision model that needs to be lifted to 3D, we

choose an appropriate intermediate layer and set  $\mathcal{G}_{\text{enc}}, \mathcal{G}_{\text{dec}}$  based on the choice. During inference, given an unseen 2D vision model that needs to be lifted to 3D, we choose an appropriate intermediate layer and set  $\mathcal{G}_{\text{enc}}, \mathcal{G}_{\text{dec}}$  based on the choice. Next, the input views are encoded using  $\mathcal{G}_{\text{enc}}$ , and a target feature map from any arbitrary view is rendered  $\hat{\mathbf{G}}_{\text{target}}$ . Finally, the desired view-consistent prediction is decoded as  $\mathcal{G}_{\text{dec}}(\hat{\mathbf{G}}_{\text{target}})$ . Our model can generalize out of the box since it simplifies 3D feature prediction to a simple interpolation task, which is easy to learn. Such interpolation schemes also support truncated or padded features when coming across dimension mismatch with our pre-trained model. We acknowledge that a view interpolation technique like ours can lead to loss in visual quality, arguably very insignificant (see Figures 1, 9). However, in several tasks multi-view consistency is more desirable and repurposing 2D models is a promising direction atleast until the data scarcity problem in 3D is solved.

## 5. Results

We conduct experiments to compare *Lift3D* against state-of-the-art methods for 3D semantic segmentation, style transfer and scene editing. Since our method is agnostic to 2D models and task domains, we further go on to show results on several other tasks (some not looked at previously in the 3D domain) to demonstrate versatility to generalize across various feature backbones. For all the experiments listed below, our method is never fine-tuned to each scene or trained on the target task or the target feature encoder’s features.

### 5.1. Implementation Details

We train our entire pipeline end-to-end on datasets of multi-view posed images. More specifically, we use the training data from IBRNet [55] consisting of synthetic [11] and real data [14, 36, 55, 68]. During training, we randomly sample  $N \in (8, 12)$  source views from a pool of  $k \times N$  (where  $k \in (1, 3)$ ) nearby views from the target view. This sampling strategy simulates varying view densities and therefore helps the network generalize better [55]. At every training iteration, we randomly choose a 2D vision model, one of DINO [4] (ViT-B/8) and CLIP [50] (ViT-B/16) to encode the source and target views. Our method is jointly optimized for view consistent RGB and feature interpolation using the Adam optimizer with an initial learning rate of  $5 \times 10^{-4}$ , decayed over the course of 250,000 training steps. We sample 2048 rays at every iteration, further sampled into 64 coarse and 128 fine points. During inference, we simply apply our lifting process on any given unseen 2D vision model and for any incoming scene without additional optimization.

### 5.2. Semantic Segmentation

**Task Description.** Recent methods such as DFF-DINO [25], N3F [48] and ISRF [15] propose to distill semantic information extracted from 2D models, for example, DINO [4] onto the 3D feature volume that can be effectively queried (for example using user strokes) to segment objects

Models	Chess Table		Color Fountain		Stove		Shoe Rack	
	IoU↑	mAP↑	IoU↑	mAP↑	IoU↑	mAP↑	IoU↑	mAP↑
N3F [48]	0.344	0.334	0.871	0.871	0.416	0.387	0.589	0.582
ISRF [15]	0.912	0.916	0.927	0.927	0.819	0.817	0.861	0.869
Lift3D	0.824	0.83	0.935	0.938	0.817	0.814	0.871	0.875

Table 1. Quantitative results for semantic segmentation using user-guided strokes. We report average scores across the validation views for each scene prescribed by ISRF [15]. The **best** scores and **second best** scores are highlighted with their respective colors.

present in the scene. Specifically, the user provides positive strokes over a region of interest in one view that can be used to determine a high-confidence *seed region* of the object to the segmented. Further, feature matching is performed to match the marked features with the distilled semantics across multiple views. In the context of our method, we lift features from a 2D DINO model (specifically the ViT-S/8 variant) and utilize the same feature-matching strategy used in [15] to segment the desired region. We follow the experiment protocol in ISRF to evaluate the performance of our method.

**Discussion.** We compare our method against recent scene-specific methods N3F [48], and ISRF [15] Table 1 presents the segmentation metrics computed over the validation views of each of the 4 scenes presented in ISRF [15]. It is important to note here that we lift a different variant of the 2D DINO encoder (ViT-S/8) unseen during training to further enforce our “*universality*” claim, while other methods use the much larger variant trained on 768 feature dimensions. In spite of these differences, our method performs as well as other methods and sometimes even outperforms them on specific scenes. This indicates that directly lifting the 2D features to a 3D volume instead of 2D – 3D feature distillation ensures better retention of the original 2D feature backbone’s capabilities. In Fig. 4, we visualize these results qualitatively and we can clearly see that our method shows superior segmentation qualitatively (e.g. near the shoe laces in Fig. 4). Even in worse-performing scenes by IoU and mAP (e.g. chesstable), we argue that our method segments clearer boundaries (near the stem of the chess table in Fig. 4) and perhaps the lower metrics are due to the missing base not included in the user stroke leading to ambiguity.

### 5.3. Style Transfer

**Task Description.** Given multi-view images of a 3D scene and an image capturing a target style, 3D style transfer aims to generate novel views of the 3D scene that have the target style consistently across the generated views. Naively combining 3D novel view synthesis and 2D style transfer often leads to multi-view inconsistency or poor stylization quality. However, our method Lift3D can transfer a pre-trained 2D feature backbone (SOTA 2D style transfer backbone CAST [66]) to a 3D model to produce view-consistent predictions from arbitrary viewing angles. To evaluate our method on this task, we compare against the classical image style transfer method AdaIN [20], SOTA video style transfer methods CCPL [58], ReReVST [56] and more recent 3D style transfer methods Hyper [7], LSNV [19] and StyleRF [30].



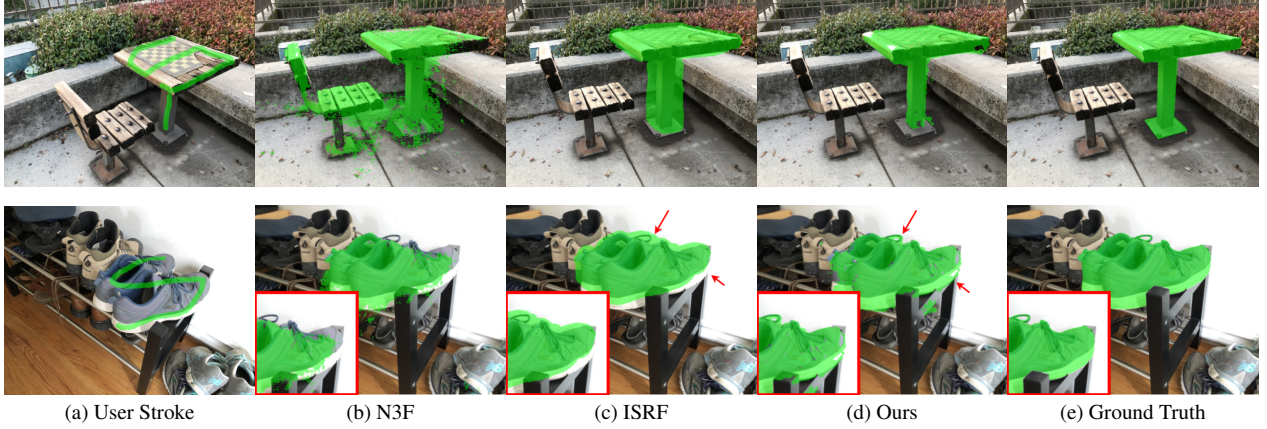


Figure 4. Qualitative results for semantic segmentation using user-provided stroke input. In the chess table scene (row 1), our method can derive fine-grained semantic masks, especially around the stem of the table. In the shoe rack scene (row 2), our method can successfully discover the shoe laces and shoe sole better than other methods.

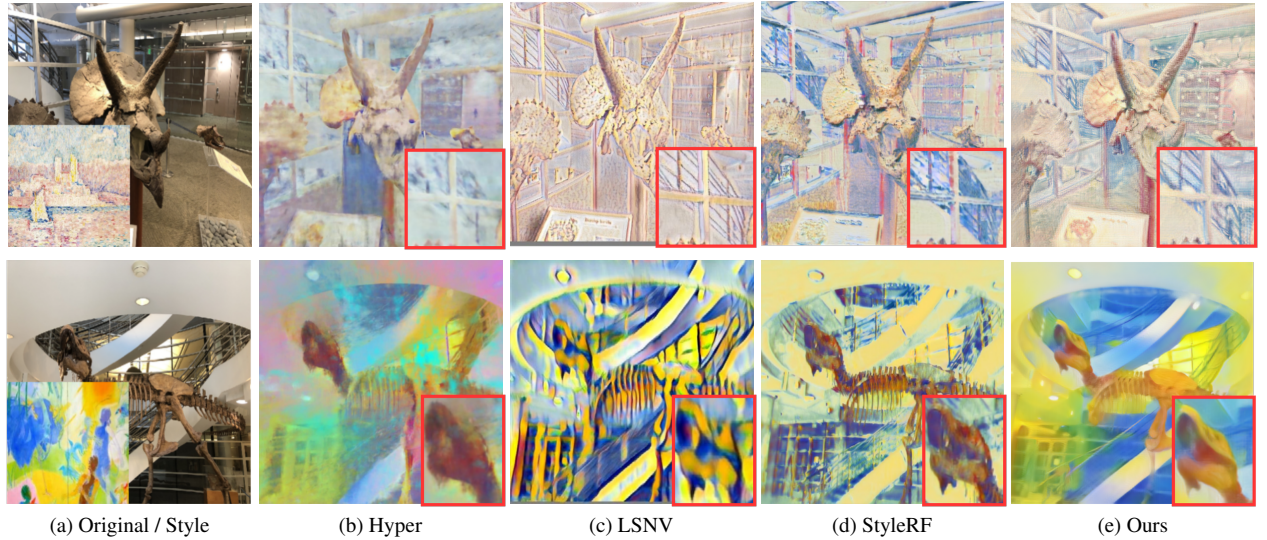


Figure 5. Qualitative results for 3D scene style transfer. Our method achieves better retention of the original geometry (row 1) and is more coherent to the style image (row 2).

Models	Short-Range Consistency		Long-Range Consistency	
	LPIPS↓	RMSE↓	LPIPS↓	RMSE↓
AdaIN [20]	0.279	0.085	0.452	0.167
CCPL [58]	0.250	0.080	0.414	0.167
ReReVST [56]	0.208	0.071	0.379	0.160
LSNV [19]	0.191	0.075	0.344	0.163
Hyper [7]	0.158	0.066	0.310	0.145
StyleRF [30]	0.138	0.071	0.333	0.165
Lift3D	0.181	0.053	0.327	0.138

Table 2. Quantitative results for 3D scene style transfer. We compute the short and long-term consistency metrics across all 8 scenes from the LLFF dataset.

**Discussion.** We evaluate our method over the LLFF dataset [34], containing real forward-facing scenes with complex geometry structures. Following the protocol in StyleRF, we measure the multi-view stylization consistency by warping one view to the other based on optical flow [47] and then compute the masked RMSE and LPIPS score [65] across nearby views (short-range consistency) and far-away views

(long-range consistency). It can be seen from Table 2 that our method achieves the best RMSE score and competitive LPIPS scores when compared to other methods. Please note that these metrics do not adequately measure the stylization quality and over-smoothed results (e.g. in the case of Hyper) tend to produce better scores [30]. Fig. 5 visualizes qualitative results and we can clearly see that our method achieves better stylization (row 2) while retaining the original scene geometry better (row 1). Furthermore, our method Lift3D generalizes to the style transfer task in a zero-shot manner unlike other stylization-specific methods that still require per-scene tuning.

## 5.4. Scene Editing.

**Task Description.** In this task, we look at editing a 3D implicit representation (in this case a NeRF) using text instructions. We attempt to lift a 2D diffusion model Instruct-Pix2Pix [3] to construct a view-consistent feature volume





Figure 6. Qualitative results on the text-driven scene editing task. Our method consistently achieves superior editing quality while still retaining the original scene geometry compared to previous work. In row 3, we can clearly see that our method retains the structure of the house, and clothes worn by the people while still adhering to the text prompt.

that can be decoded for text-based 3D editing. To evaluate our method on this task, we compare against CLIP-NeRF [51], NeRF-Art [52] and Masked SDS [38] that fine-tune the parameters of a NeRF model using 2D supervision either from a pretrained CLIP or Diffusion model, and more recent methods like Instruct NeRF-to-NeRF [17] that directly manipulate the training images of NeRF.

**Discussion.** As seen in Fig. 6, our method shows the capability to edit scenes with superior quality compared to other methods. Prior work that leverages CLIP-based optimization [51, 52] struggles to produce reasonable edited images due to the inherent limitations in the CLIP model. More recent diffusion-based approaches [17, 38] produce higher quality results but either alter the image too severely or too little (see Figs. 6d, 6e), while our method consistently does well across several edits (we present more results in the supplementary). Further, we provide more examples that compare our method against InstructNeRF2NeRF, a recent SOTA method for instruction-based editing of 3D scenes, and leverage the same 2D feature backbone as ours. We can clearly see that our method ensures retention of the original scene geometry (see Figs. 6i, 6l) while still adhering to the editing text prompt. These results clearly show that directly lifting the intermediate latent of the 2D diffusion model outperforms naive optimization strategies that edit NeRF’s training images. We defer quantitative results that measure view-consistency and text-similarity of the edited scenes to the supplementary (see Table 5).

## 5.5. Other Tasks

In the previous subsections, we discussed the generality of our proposed method to lift several 2D feature backbones to generate 3D consistent predictions without any extra training. This enables us to realize tasks that have not been previously looked at in the 3D domain. This underscores the generality of our method i.e. as long as there exists a 2D feature backbone for a given task, Lift3D helps realize the same in 3D without any extra training or efforts on data collection.

**Open Vocabulary Segmentation.** Given the multi-view images of a 3D scene and the open-vocabulary text descriptions, we wish to derive accurate object boundaries for the scene. Some recent work [23] distills large 2D vision-language models onto a 3D radiance field, which can be effectively used to obtain relevancy maps for language queries. However, the obtained semantic masks are very coarse and cannot be used for localization. Therefore, the task of open vocabulary segmentation in the 3D domain is still an open problem. Towards this end, we transfer the intermediate features of a SOTA 2D open vocabulary segmentation method OVSeg [28] to a 3D feature volume to derive multi-view consistent and accurate semantic masks. Figures 7a and 7b visualize predictions obtained at multiple viewpoints by 1) naively applying 2D OVSeg on multi-view data and 2) using the 3D “Lifted” OVSeg features estimated by Lift3D. We can clearly see that the lifted features ensure multi-view consistent predictions (also see Fig. 1) and surprisingly even improve segmentation quality (precise keyboard boundaries).



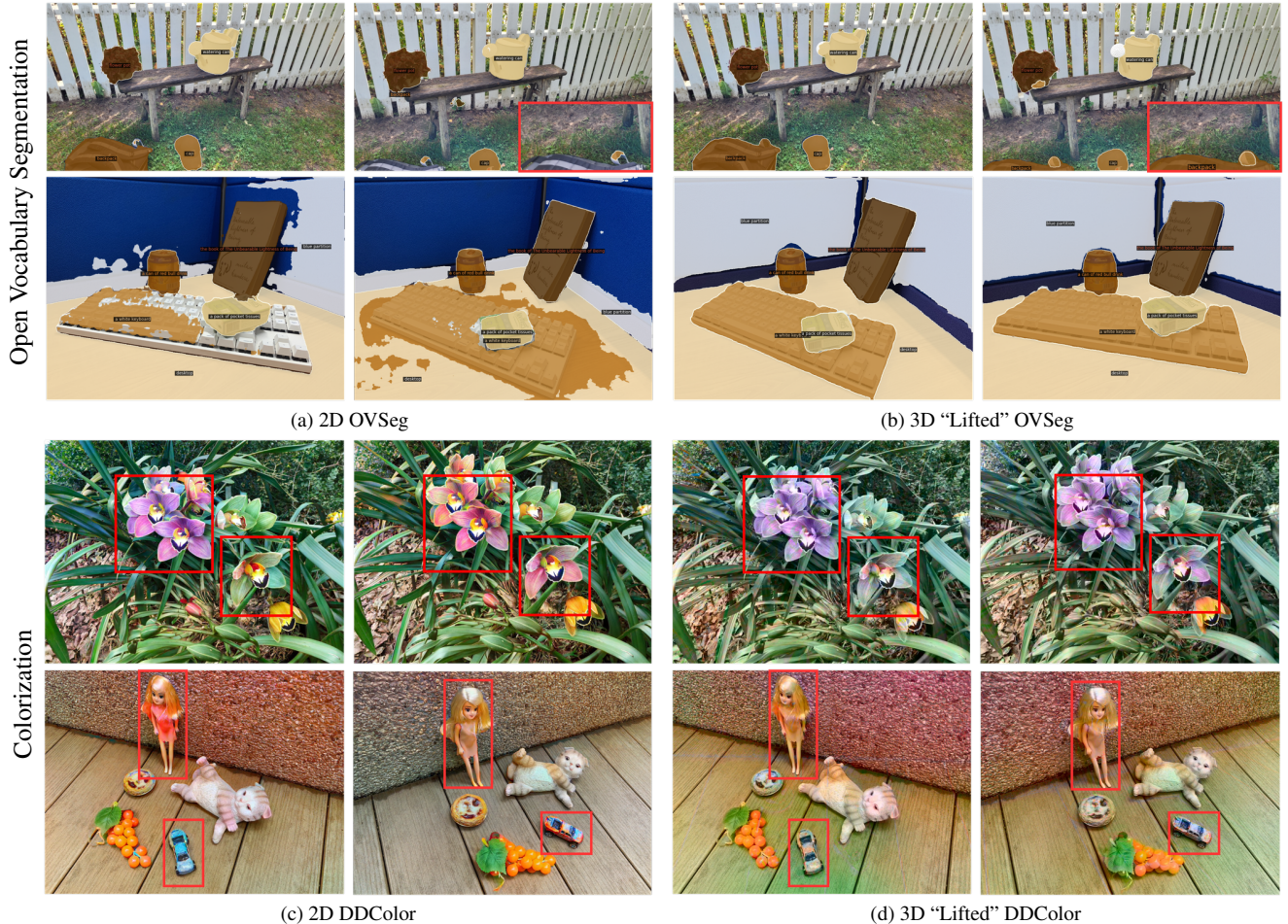


Figure 7. Qualitative results on open vocabulary segmentation and colorization, across two views. Our 3D “lifted” features can successfully segment occluded objects (row 1) and with superior quality (row 2). Moreover, they also show improved view consistency compared to the 2D operation in both tasks.

By using multi-view information our model can reason about occlusions and other variabilities across arbitrary viewpoints, resulting in enhanced segmentation masks.

**Image Colorization.** The image colorization task utilizes grayscale multi-view images to generate plausible colored NeRF scenes. The ambiguity in 2D image colorization leads to severe multi-view inconsistencies and we verify if our proposed strategy can rectify them by adhering to multi-view geometry. We leverage SOTA single view colorization technique DDColor [22] and construct a 3D feature volume that can be decoded to estimate the colored image. Figs. 7c and 7d compare the lifted 3D features against 2D image colorization of each view individually. We can clearly see that our method ensures consistent colors across different views while preserving a similar output quality as the original network. These results clearly verify the adaptability of our proposed pipeline to several tasks.

## 6. Conclusion

We present *Lift3D*, a generalizable system that can lift any 2D visual model to synthesize view-consistent feature predic-

tions without training with data from downstream tasks. Our method essentially learns to rectify and propagate predicted feature maps of source views to synthesize the feature map for the novel view. Our algorithm mitigates inconsistencies among source view predictions, and generates view-smooth predictions at the target view. We demonstrate that Lift3D is merely pre-trained on DINO and CLIP features but can directly generalize to a wide range of 2D vision models, empowering various applications, including semantic segmentation, stylization, instructed scene editing, and many others. All empirical observations endorse that Lift3D can be a crucial component in bringing the recent advancement of 2D vision models to the 3D domain.

## Acknowledgements

This work was funded in part by a Jacobs School of Engineering Fellowship, ONR grant (N00014-23-1-2526), NSF grants (2100237, 2120019), NSF CAREER Award (2240160), NSF TILOS AI Institute (2112665), the Ronald L. Graham Chair and the UC San Diego Center for Visual Computing. We also acknowledge support from IARPA via contract 140D0423C0076 and gifts from Adobe, Google, Qualcomm and Rembrand.



## References

- [1] Jonathan T Barron, Ben Mildenhall, Matthew Tancik, Peter Hedman, Ricardo Martin-Brualla, and Pratul P Srinivasan. Mip-nerf: A multiscale representation for anti-aliasing neural radiance fields. In *IEEE International Conference on Computer Vision (ICCV)*, 2021. 2
- [2] Jonathan T Barron, Ben Mildenhall, Dor Verbin, Pratul P Srinivasan, and Peter Hedman. Mip-nerf 360: Unbounded anti-aliased neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5470–5479, 2022. 2
- [3] Tim Brooks, Aleksander Holynski, and Alexei A Efros. Instructpix2pix: Learning to follow image editing instructions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18392–18402, 2023. 1, 6
- [4] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 9650–9660, 2021. 2, 5
- [5] Hanlin Chen, Chen Li, Mengqi Guo, Zhiwen Yan, and Gim Hee Lee. Gnesf: Generalizable neural semantic fields. *arXiv preprint arXiv:2310.15712*, 2023. 2
- [6] Runnan Chen, Youquan Liu, Lingdong Kong, Nenglu Chen, ZHU Xinge, Yuexin Ma, Tongliang Liu, and Wenping Wang. Towards label-free scene understanding by vision foundation models. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. 2
- [7] Pei-Ze Chiang, Meng-Shiun Tsai, Hung-Yu Tseng, Wei-Sheng Lai, and Wei-Chen Chiu. Stylizing 3d scene via implicit representation and hypernetwork. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 1475–1484, 2022. 5, 6
- [8] Wenyan Cong, Hanxue Liang, Peihao Wang, Zhiwen Fan, Tianlong Chen, Mukund Varma, Yi Wang, and Zhangyang Wang. Enhancing nerf akin to enhancing llms: Generalizable nerf transformer with mixture-of-view-experts. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3193–3204, 2023. 2
- [9] Paul E Debevec, Camillo J Taylor, and Jitendra Malik. Modeling and rendering architecture from photographs: A hybrid geometry and image-based approach. In *Proceedings of the 23rd annual conference on Computer graphics and interactive techniques*, pages 11–20, 1996. 2
- [10] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations (ICLR)*, 2020. 2
- [11] Laura Downs, Anthony Francis, Nate Koenig, Brandon Kinman, Ryan Hickman, Krista Reymann, Thomas B McHugh, and Vincent Vanhoucke. Google scanned objects: A high-quality dataset of 3d scanned household items. *arXiv e-prints*, pages arXiv-2204, 2022. 5
- [12] Zhiwen Fan, Yifan Jiang, Peihao Wang, Xinyu Gong, Dejia Xu, and Zhangyang Wang. Unified implicit neural stylization. *arXiv preprint arXiv:2204.01943*, 2022. 2
- [13] Zhiwen Fan, Peihao Wang, Yifan Jiang, Xinyu Gong, Dejia Xu, and Zhangyang Wang. Nerf-sos: Any-view self-supervised object segmentation on complex scenes. *arXiv preprint arXiv:2209.08776*, 2022. 2
- [14] John Flynn, Michael Broxton, Paul Debevec, Matthew Duvall, Graham Fyffe, Ryan Overbeck, Noah Snavely, and Richard Tucker. Deepview: View synthesis with learned gradient descent. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2367–2376, 2019. 5
- [15] Rahul Goel, Dhawal Sirikonda, Saurabh Saini, and PJ Narayanan. Interactive segmentation of radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4201–4211, 2023. 5
- [16] Steven Gortler, Radek Grzeszczuk, Richard Szeliski, and Michael F Cohen. The lumigraph. In *Proceedings of the 23rd annual conference on computer graphics and interactive techniques (SIGGRAPH 1996)*, pages 43–54, 1996. 2
- [17] Ayaan Haque, Matthew Tancik, Alexei A Efros, Aleksander Holynski, and Angjoo Kanazawa. Instruct-nerf2nerf: Editing 3d scenes with instructions. *arXiv preprint arXiv:2303.12789*, 2023. 2, 7
- [18] Di Huang, Sida Peng, Tong He, Honghui Yang, Xiaowei Zhou, and Wanli Ouyang. Ponder: Point cloud pre-training via neural rendering. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 16089–16098, 2023. 2
- [19] Hsin-Ping Huang, Hung-Yu Tseng, Saurabh Saini, Maneesh Singh, and Ming-Hsuan Yang. Learning to stylize novel views. In *ICCV*, 2021. 5, 6
- [20] Xun Huang and Serge Belongie. Arbitrary style transfer in real-time with adaptive instance normalization. In *Proceedings of the IEEE international conference on computer vision*, pages 1501–1510, 2017. 5, 6
- [21] Ahmet Iscen, Giorgos Tolias, Yannis Avrithis, and Ondrej Chum. Label propagation for deep semi-supervised learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5070–5079, 2019. 2
- [22] Xiaoyang Kang, Tao Yang, Wenqi Ouyang, Peiran Ren, Lingzhi Li, and Xuansong Xie. Ddcolor: Towards photo-realistic image colorization via dual decoders. *arXiv preprint arXiv:2212.11613*, 2022. 8, 1
- [23] Justin Kerr, Chung Min Kim, Ken Goldberg, Angjoo Kanazawa, and Matthew Tancik. Lrf: Language embedded radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 19729–19739, 2023. 7
- [24] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. *arXiv preprint arXiv:2304.02643*, 2023. 1, 2
- [25] Sosuke Kobayashi, Eiichi Matsumoto, and Vincent Sitzmann. Decomposing nerf for editing via feature field distillation. *Advances in Neural Information Processing Systems*, 35:23311–23330, 2022. 5
- [26] Jaakko Lehtinen, Jacob Munkberg, Jon Hasselgren, Samuli Laine, Tero Karras, Miika Aittala, and Timo Aila. Noise2noise: Learning image restoration without clean data. *arXiv preprint arXiv:1803.04189*, 2018. 4
- [27] Marc Levoy and Pat Hanrahan. Light field rendering. In *Proceedings of the 23rd annual conference on Computer graphics and interactive techniques*, pages 31–42, 1996. 2

- [28] Feng Liang, Bichen Wu, Xiaoliang Dai, Kunpeng Li, Yanan Zhao, Hang Zhang, Peizhao Zhang, Peter Vajda, and Diana Marculescu. Open-vocabulary semantic segmentation with mask-adapted clip. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7061–7070, 2023. 1, 7
- [29] Fangfu Liu, Chubin Zhang, Yu Zheng, and Yueqi Duan. Semantic ray: Learning a generalizable semantic field with cross-reprojection attention. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 17386–17396, 2023. 2
- [30] Kunhao Liu, Fangneng Zhan, Yiwen Chen, Jiahui Zhang, Yingchen Yu, Abdulmotaleb El Saddik, Shijian Lu, and Eric P Xing. Stylef: Zero-shot 3d style transfer of neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8338–8348, 2023. 2, 5, 6
- [31] Kunhao Liu, Fangneng Zhan, Jiahui Zhang, Muyu Xu, Yingchen Yu, Abdulmotaleb El Saddik, Christian Theobalt, Eric Xing, and Shijian Lu. Weakly supervised 3d open-vocabulary segmentation. *arXiv preprint arXiv:2305.14093*, 2023. 2
- [32] Youquan Liu, Lingdong Kong, Jun Cen, Runnan Chen, Wenwei Zhang, Liang Pan, Kai Chen, and Ziwei Liu. Segment any point cloud sequences by distilling vision foundation models. *arXiv preprint arXiv:2306.09347*, 2023. 2
- [33] Nelson Max. Optical models for direct volume rendering. *IEEE Transactions on Visualization and Computer Graphics (TVCG)*, 1995. 2
- [34] Ben Mildenhall, Pratul P Srinivasan, Rodrigo Ortiz-Cayon, Nima Khademi Kalantari, Ravi Ramamoorthi, Ren Ng, and Abhishek Kar. Local light field fusion: Practical view synthesis with prescriptive sampling guidelines. *ACM Transactions on Graphics (TOG)*, 38(4):1–14, 2019. 6
- [35] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In *European conference on computer vision*, pages 405–421. Springer, 2020. 2, 3, 4
- [36] Ben Mildenhall, Peter Hedman, Ricardo Martin-Brualla, Pratul Srinivasan, and Jonathan T Barron. Nerf in the dark: High dynamic range view synthesis from noisy raw images. *arXiv preprint arXiv:2111.13679*, 2021. 5
- [37] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning robust visual features without supervision. *arXiv preprint arXiv:2304.07193*, 2023. 2
- [38] Ben Poole, Ajay Jain, Jonathan T Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d diffusion. In *The Eleventh International Conference on Learning Representations*, 2022. 7
- [39] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *International Conference on Machine Learning*, pages 8821–8831. PMLR, 2021. 2
- [40] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 1(2): 3, 2022. 2
- [41] Tal Ridnik, Emanuel Ben-Baruch, Asaf Noy, and Lihi Zelnik-Manor. Imagenet-21k pretraining for the masses. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 1)*, 2021. 2
- [42] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10684–10695, 2022. 2
- [43] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic text-to-image diffusion models with deep language understanding. *Advances in Neural Information Processing Systems*, 35:36479–36494, 2022. 2
- [44] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. Laion-5b: An open large-scale dataset for training next generation image-text models. *Advances in Neural Information Processing Systems*, 35:25278–25294, 2022. 2
- [45] Mohammed Suhail, Carlos Esteves, Leonid Sigal, and Amesh Makadia. Light field neural rendering. *arXiv preprint arXiv:2112.09687*, 2021. 2
- [46] Mohammed Suhail, Carlos Esteves, Leonid Sigal, and Amesh Makadia. Generalizable patch-based neural rendering. *arXiv preprint arXiv:2207.10662*, 2022. 2, 3
- [47] Zachary Teed and Jia Deng. Raft: Recurrent all-pairs field transforms for optical flow. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part II 16*, pages 402–419. Springer, 2020. 6
- [48] Vadim Tschernezki, Iro Laina, Diane Larlus, and Andrea Vedaldi. Neural feature fusion fields: 3d distillation of self-supervised 2d image representations. In *2022 International Conference on 3D Vision (3DV)*, pages 443–453. IEEE, 2022. 5
- [49] Mukund Varma, Peihao Wang, Xuxi Chen, Tianlong Chen, Subhashini Venugopalan, and Zhangyang Wang. Is attention all that nerf needs? In *The Eleventh International Conference on Learning Representations*, 2022. 2, 3, 1
- [50] Can Wang, Menglei Chai, Mingming He, Dongdong Chen, and Jing Liao. Clip-nerf: Text-and-image driven manipulation of neural radiance fields. *arXiv preprint arXiv:2112.05139*, 2021. 2, 5
- [51] Can Wang, Menglei Chai, Mingming He, Dongdong Chen, and Jing Liao. Clip-nerf: Text-and-image driven manipulation of neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3835–3844, 2022. 2, 7
- [52] Can Wang, Ruixiang Jiang, Menglei Chai, Mingming He, Dongdong Chen, and Jing Liao. Nerf-art: Text-driven neural radiance fields stylization. *IEEE Transactions on Visualization and Computer Graphics*, 2023. 7
- [53] Fei Wang and Changshui Zhang. Label propagation through linear neighborhoods. In *Proceedings of the 23rd international conference on Machine learning*, pages 985–992, 2006. 2
- [54] Peng Wang, Lingjie Liu, Yuan Liu, Christian Theobalt, Taku Komura, and Wenping Wang. Neus: Learning neural implicit surfaces by volume rendering for multi-view reconstruction. *arXiv preprint arXiv:2106.10689*, 2021. 2

- [55] Qianqian Wang, Zhicheng Wang, Kyle Genova, Pratul P Srinivasan, Howard Zhou, Jonathan T Barron, Ricardo Martin-Brualla, Noah Snavely, and Thomas Funkhouser. Ibrnet: Learning multi-view image-based rendering. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 2, 3, 4, 5
- [56] Wenjing Wang, Shuai Yang, Jizheng Xu, and Jiaying Liu. Consistent video style transfer via relaxation and regularization. *IEEE Trans. Image Process.*, 2020. 5, 6
- [57] Suttisak Wizadwongsa, Pakkapon Phongthawee, Jiraphon Yenphraphai, and Supasorn Suwajanakorn. Nex: Real-time view synthesis with neural basis expansion. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 2
- [58] Zijie Wu, Zhen Zhu, Junping Du, and Xiang Bai. Ccpl: Contrastive coherence preserving loss for versatile style transfer. In *European Conference on Computer Vision*, pages 189–206. Springer, 2022. 5, 6
- [59] Saining Xie, Jiatao Gu, Demi Guo, Charles R Qi, Leonidas Guibas, and Or Litany. Pointcontrast: Unsupervised pre-training for 3d point cloud understanding. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part III 16*, pages 574–591. Springer, 2020. 2
- [60] Jianglong Ye, Naiyan Wang, and Xiaolong Wang. Featurenerf: Learning generalizable nerfs by distilling foundation models. *arXiv preprint arXiv:2303.12786*, 2023. 2, 3, 1
- [61] Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. pixelnerf: Neural radiance fields from one or few images. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 2, 3
- [62] Zongsheng Yue, Jianyi Wang, and Chen Change Loy. Resshift: Efficient diffusion model for image super-resolution by residual shifting. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. 1
- [63] Manzil Zaheer, Satwik Kottur, Siamak Ravanbakhsh, Barnabas Poczos, Russ R Salakhutdinov, and Alexander J Smola. Deep sets. *Advances in neural information processing systems*, 30, 2017. 3
- [64] Kai Zhang, Nick Kolkin, Sai Bi, Fujun Luan, Zexiang Xu, Eli Shechtman, and Noah Snavely. Arf: Artistic radiance fields. In *European Conference on Computer Vision*, pages 717–733. Springer, 2022. 2
- [65] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *CVPR*, 2018. 6
- [66] Yuxin Zhang, Fan Tang, Weiming Dong, Haibin Huang, Chongyang Ma, Tong-Yee Lee, and Changsheng Xu. Domain enhanced arbitrary image style transfer via contrastive learning. In *ACM SIGGRAPH 2022 Conference Proceedings*, pages 1–8, 2022. 1, 5
- [67] Dengyong Zhou, Olivier Bousquet, Thomas Lal, Jason Weston, and Bernhard Schölkopf. Learning with local and global consistency. *Advances in neural information processing systems*, 16, 2003. 2
- [68] Tinghui Zhou, Richard Tucker, John Flynn, Graham Fyffe, and Noah Snavely. Stereo magnification: Learning view synthesis using multiplane images. *arXiv preprint arXiv:1805.09817*, 2018. 5
- [69] Haoyi Zhu, Honghui Yang, Xiaoyang Wu, Di Huang, Sha Zhang, Xianglong He, Tong He, Hengshuang Zhao, Chunhua Shen, Yu Qiao, et al. Ponderv2: Pave the way for 3d foundation model with a universal pre-training paradigm. *arXiv preprint arXiv:2310.08586*, 2023. 2
- [70] Xiaojin Zhu and Zoubin Ghahramani. Learning from labeled and unlabeled data with label propagation. *ProQuest Number: INFORMATION TO ALL USERS*, 2002. 2
- [71] Xiaojin Zhu, Zoubin Ghahramani, and John D Lafferty. Semi-supervised learning using gaussian fields and harmonic functions. In *Proceedings of the 20th International conference on Machine learning (ICML-03)*, pages 912–919, 2003. 2



# Lift3D: Zero-Shot Lifting of Any 2D Vision Model to 3D

## Supplementary Material

Model	ResShift [62]	DDColor [22]	OVSeg [28]
Feature Prediction	3.121	12.282	10.117
Feature Interpolation			
w/o Inconsistency Correction	0.042	7.491	0.375
w/ Inconsistency Correction			
Single Stage	0.0395	3.39	0.371
Ours	0.021	3.240	0.370

Table 3. Ablation study of several components of Lift3D on lifting three unseen 2D feature encoders to 3D. The indent indicates the studied setting is added upon the upper-level ones.

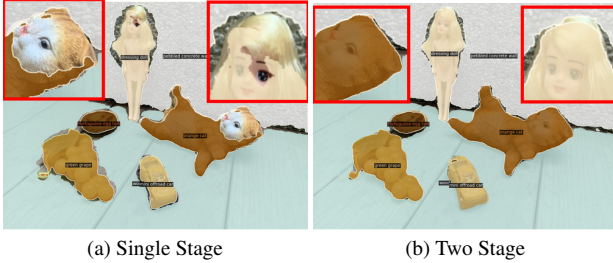


Figure 8. Qualitative comparison on the open vocabulary segmentation task between a Single Stage and Two Stage feature correction pipeline.

### A. Ablation Studies

We primarily ablate on the following key design considerations of our method.

**Feature Prediction.** Inspired by [60], we propose a baseline that directly predicts the target view feature and color. To do so, we leverage a powerful generalizable NVS method GNT [49], and add an additional feature head that is supervised by the 2D vision model.

**w/o Inconsistency Correction.** Next, We remove the inconsistency correction module on the source view features and simply share the aggregation weights between the epipolar RGB features and similarly constructed epipolar projections on the encoder features from the 2D vision model.

**Single Stage.** Lastly, we convert our two-stage pipeline into a “one-stage” feature correction pipeline using the randomly sampled points along the ray (in contrast to our two-stage, density proposal followed by feature propagation only on the importance-sampled points).

We report performance on the above investigations in Table 3, specifically the mean squared error distance between the estimated feature and ground truth feature obtained by naively encoding the target view using the 2D vision model (or  $G_{\text{target}}$ ). We follow the training strategy discussed in Sec. 5.1 *i.e.* trained on DINO and CLIP features and evaluated on unseen 2D vision models - ResShift [62], DDColor [22], and OVSeg [28].

A straightforward extension of [60], that directly predicts the target view feature cannot generalize to unseen feature encoders resulting in a very high error rate. Instead,

2D DINO	3D “Lifted” Dino		
N.A.	3 views	6 views	10 views
0.39 / 0.90 / 0.36	0.76 / 0.97 / 0.76	0.80 / 0.98 / 0.80	0.82 / 0.98 / 0.83

Table 4. Effect of the number of input views for semantic segmentation of scenes. Metrics are ordered as IoU / Acc / mAP (higher is better)

our method uses a feature interpolation strategy that derives consistency information from RGB to rectify and propagate inconsistent feature maps. We verify that our inconsistency correction module on the source view features  $\{G_i\}_{i=1}^N$  is essential and ensures better blending of feature maps. This becomes even more apparent in the case of severely view-inconsistent input features, *e.g.* in the case of the colorization task (DDColor [22]). Finally, our method is two-stage *i.e.* derives a coarse density proposal and performs corrective feature aggregation only on the importance-sampled points. This is necessary and leads to slight deviations in the estimated novel view feature otherwise (and even worse decoded outputs, see Fig. 8).

### B. Computational Efficiency

Naively applying a 2D vision model on each rendered view yields multi-view inconsistent predictions and can be quite inefficient. On the other hand, our method *Lift3D* uses the 2D vision model to encode only the training views (can be pre-computed) and relies on nearest source views when estimating the features from arbitrary viewing angles, that are further decoded to obtain the desired output. This is significantly efficient and faster especially when performing the downstream task from a large number of arbitrary viewpoints.

Formally, let’s assume the time to encode, decode, and render each view as  $t_{\text{enc}}$ ,  $t_{\text{dec}}$ ,  $t_{\text{rend}}$  respectively. To perform the desired task on 100 rendered views, the 2D baseline would roughly take time  $t_{2D} = 100 \times (t_{\text{rend}} + t_{\text{enc}} + t_{\text{dec}})$ . Assuming we have around 15 training views for each scene, the time taken by Lift3D  $t_{3D} = 15 \times (t_{\text{enc}}) + 100 \times (t_{\text{rend}} + t_{\text{dec}})$ . We can clearly see that  $t_{3D} < t_{2D}$ , and the difference is quite significant in the case of lifting diffusion features like InstructPix2Pix [3] that requires a time-consuming multistep denoising process during encoding *i.e.*  $t_{\text{enc}} \gg t_{\text{dec}}$ . Therefore, when performing downstream tasks on several arbitrary viewing angles, our method also boasts of superior efficiency along with multi-view consistency when compared to its corresponding 2D counterpart.

### C. Limitations

We acknowledge that an interpolation technique like ours from input views does result in some loss in quality, arguably not very significant. However, in several practical

Method	Text-Image Direction Similarity $\uparrow$	Direction Consistency $\uparrow$
InstructNeRF2NeRF [17]	0.180	0.966
Ours	0.193	0.982

Table 5. Quantitative results for text-guided scene editing.

applications multi-view consistent outputs are usually even more desirable, and even mild deviations from the current viewpoint yield significantly different outputs when naively applying a 2D vision model (see Figures 1, 6, 7, 9). Although our method successfully lifts many 2D features to be multi-view consistent, its potential remains capped by the epipolar-based rendering. For example, our method may not handle sparse 360-degree scenes or objects with complex light transport where epipolar geometry no longer holds and drops in performance with limited number of input views (see Table 4). Interesting future directions include scaling up training of Lift3D to unbounded scenes [8] or combining extant pre-trained 3D models with 2D models.

## D. Gallery of Tasks

In Fig. 9, we qualitatively compare the decoded outputs using our 3D lifted features against the original 2D operation on two views and across different tasks. We can clearly see that our method yields multi-view consistent predictions, unlike the 2D operator. In some cases, we even see that our method yields improved predictions perhaps due to multi-view information. For example in Figs. 9k and 9k, we can see that our method is able to segment the hair dryer along with its cable as per the input prompt “*hair dryer with cable*”. Similarly in Figs. 9n and 9o, in the case of super-resolution, we see that our lifted features preserve the original scene geometry with higher detail, unlike its 2D counterpart. For the sake of completeness, we also provide quantitative results for 3D scene editing in Table 5. Following the evaluation protocol in [17], we compute the text-image direction similarity and consistency scores across 6 scenes and report average metrics. Our method outperforms the SOTA scene-specific 3D editing technique InstructNeRF2NeRF [17], both in terms of editing quality and multi-view consistency. Fig. 9 only represents a few tasks and in practice, our method can be extended to any 2D vision operator without any extra tuning.

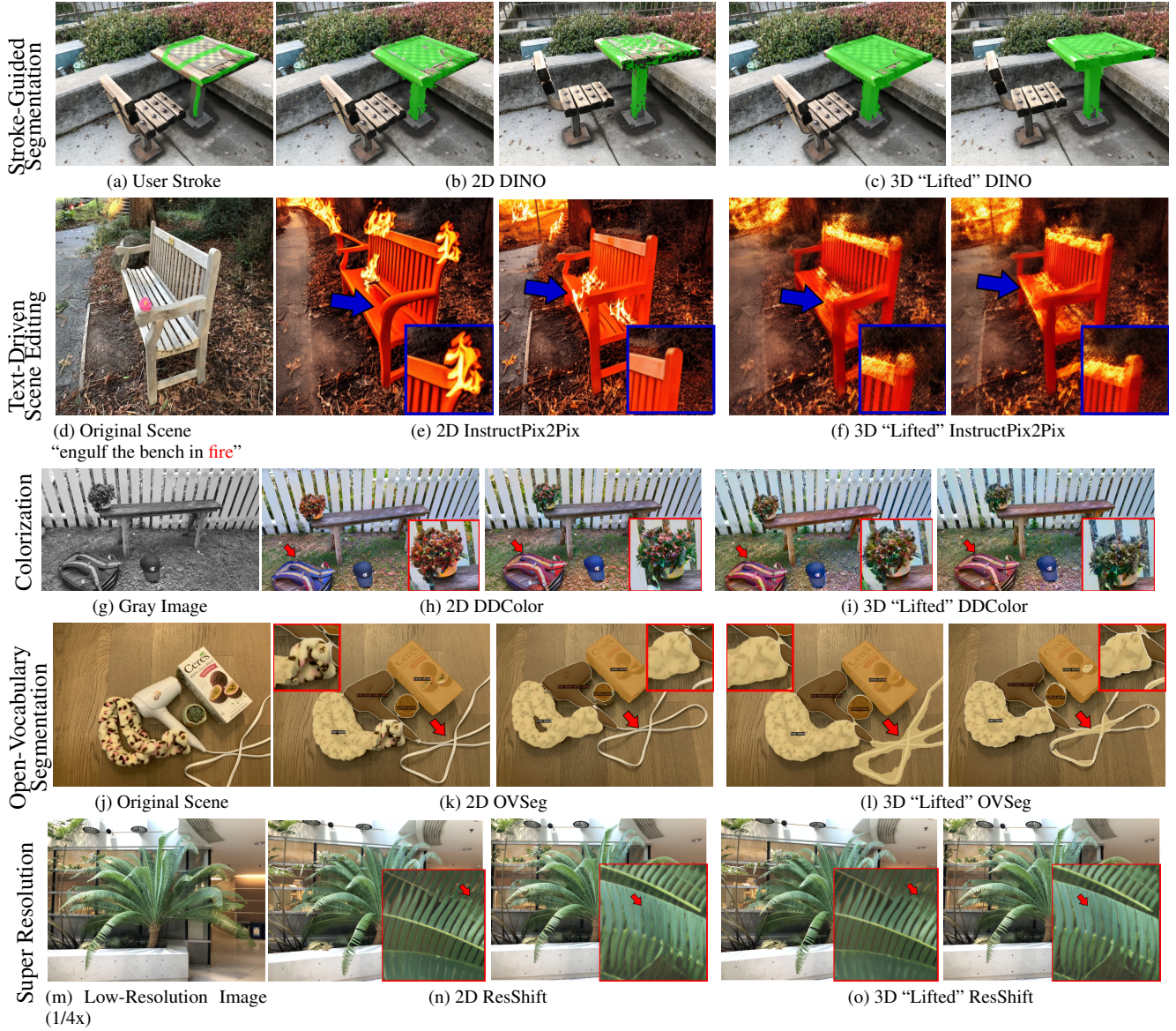


Figure 9. Qualitative comparisons of the 3D "Lifted" features against its corresponding 2D counterpart on two different views and across several tasks. We observe that our 3D-corrected features are more multi-view consistent and sometimes even improve prediction quality. For clearer comparison between the 2D and 3D outcomes, we recommend zooming into the electronic version of this image.