# *E-RayZer:* Self-supervised 3D Reconstruction as Spatial Visual Pre-training

Qitao Zhao[1]    Hao Tan[2]    Qianqian Wang[3]    Sai Bi[2]

Kai Zhang[2]    Kalyan Sunkavalli[2]    Shubham Tulsiani[1*]    Hanwen Jiang[2*]

[1]Carnegie Mellon University  [2]Adobe Research  [3]Harvard University  *Equal advising

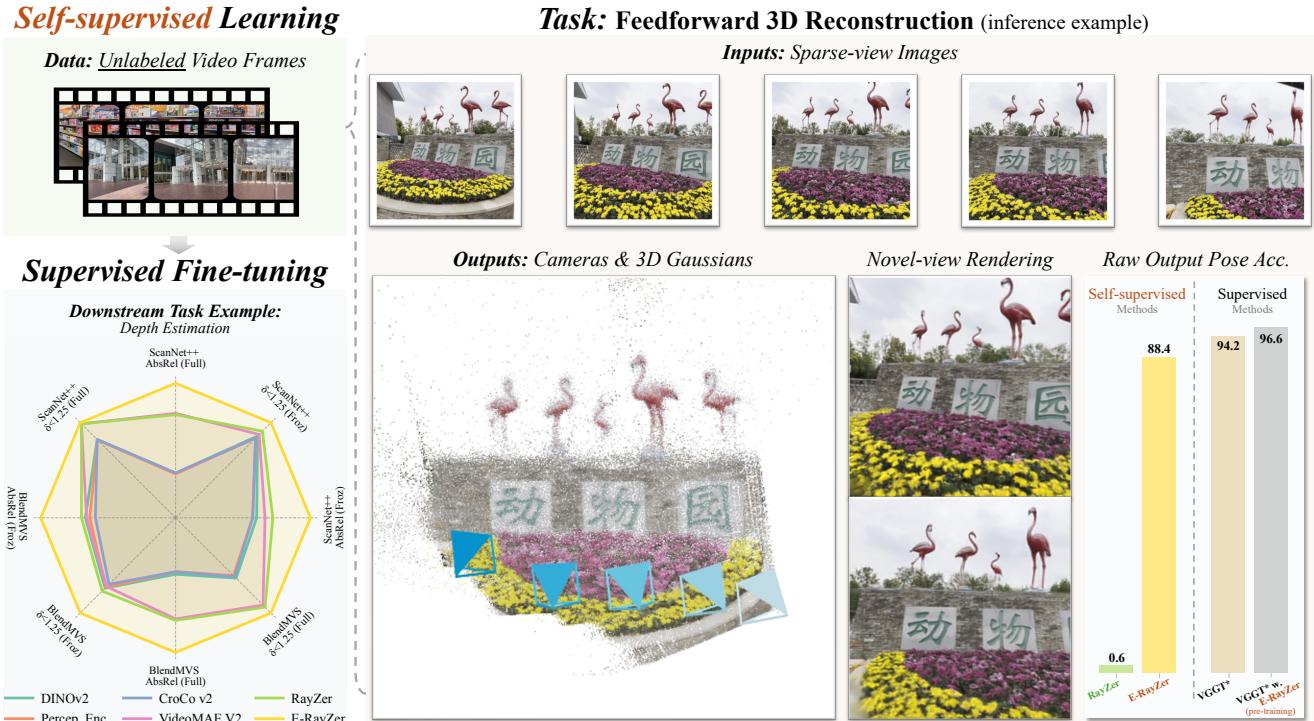*Project & Code:* qitaozhao.github.io/E-RayZer

Figure 1. **E-RayZer**, a **self-supervised** 3D Vision model predicting camera poses and scene geometry as 3D Gaussians. The use of **explicit** 3D geometry yields more geometrically grounded poses compared to its implicit counterpart, RayZer [23]: They are comparable to and sometimes even better than those from supervised state-of-the-art model VGGT. Furthermore, E-RayZer also serves as a self-supervised **visual pre-training** framework, with learned representations that transfer effectively to downstream tasks requiring 3D understanding. It outperforms previous visual representation learners, *e.g.*, CroCo v2 [62], VideoMAE V2 [57], DINOv3 [47], Perception Encoder [7], *etc*.

## Abstract

*Self-supervised pre-training has revolutionized foundation models for languages, individual 2D images and videos, but remains largely unexplored for learning 3D-aware representations from multi-view images. In this paper, we present E-RayZer, a self-supervised large 3D Vision model that learns truly 3D-aware representations directly from unlabeled images. Unlike prior self-supervised methods such as RayZer that infer 3D indirectly through latent-space view synthesis, **E-RayZer** operates directly in 3D space, performing self-supervised 3D reconstruction with **Explicit** geometry. This formulation eliminates short-cut solutions and yields representations that are geometrically grounded. To ensure convergence and scalability, we introduce a novel fine-grained learning curriculum that organizes training from easy to hard samples and harmonizes heterogeneous data sources in an entirely unsupervised manner. Experiments demonstrate that E-RayZer significantly outperforms RayZer on pose estimation, matches or sometimes surpasses fully supervised reconstruction models such as VGGT. Furthermore, its learned representations outperform leading visual pre-training models (e.g., DINOv3, CroCo v2, VideoMAE V2, and RayZer) when transferring to 3D downstream tasks, establishing E-RayZer as a new paradigm for 3D-aware visual pre-training.*

1

## 1. Introduction

**Pre-training with self-supervision** forms the foundation of frontier models, allowing them to learn *meaningful representations* on vast amounts of unlabeled data. This paradigm has proven to be effective for text [8, 11], 2D image [17, 36] and video [2, 50] domains, where large models manage to capture language semantics, visual concepts, and temporal dynamics. However, we argue that one essential component is still missing – **learning 3D-aware representations from unlabeled multi-view images**, as 3D spatial understanding is fundamental for perceiving and interacting with the 3D physical world we live in. Yet, current 3D Vision models mostly rely on a different route: *fully-supervised learning* using 3D pseudo-labels estimated by COLMAP [44], which is inherently inefficient, imperfect, and ultimately unscalable. To move forward, we need a self-supervised pre-training framework that can learn 3D-aware representations from abundant raw visual observations.

In this paper, we present **E-RayZer**, the first **truly self-supervised 3D Gaussian splatting reconstruction** model that learns 3D-aware representations from unlabeled data, thereby establishing a new paradigm for *3D spatial visual pre-training* (Fig. 1). Unlike its predecessor RayZer [23], which exhibits only *superficial* 3D awareness by learning the *proxy task* of self-supervised view synthesis in *latent space*, E-RayZer operates directly in the *3D space*, learning self-supervised 3D reconstruction. Concretely, E-RayZer predicts camera parameters and 3D Gaussians [29] from inputs, and renders them back for photometric self-supervision under the constraints of physical rendering rules. By grounding representations in **explicit** scene geometry, E-RayZer learns features that are genuinely 3D-aware and free from RayZer's shortcut solutions such as frame interpolation (see Sec. 3.1). This design not only yields a *camera space* that is more geometrically grounded and interpretable than RayZer's, but also produces *latent representations* that are truly **3D-aware**, effectively benefiting downstream 3D Vision tasks.

Although using explicit 3D Gaussians offers clear advantages, it also introduces substantial training challenges. As reported in RayZer (Tab. 7), training with explicit 3D leads to non-convergence. To address this key challenge, we propose a **fine-grained learning curriculum**, built on the concept of *visual overlap* between input views. To stabilize training, we begin with samples of *high visual overlap*, allowing the pose estimator to be initialized from predicting near-identity poses, and gradually reduce overlap to promote general 3D understanding. When *scaling* to heterogeneous training resources, visual overlap provides a natural and unified metric to adaptively align varying camera motion distributions, improving data consistency. Notably, we approximate visual overlap in an *unsupervised* way, keeping the framework entirely free from any 3D annotations.

We systematically study the performance of E-RayZer with different training data scales. We highlight key conclusions and summarize our contributions as follows:

- E-RayZer is the first **self-supervised feedforward 3DGS reconstruction** model, trained **from scratch** with zero 3D annotation.

- E-RayZer **outperforms prior visual representation learners**, *e.g.*, DINOv3 [47], CroCo v2 [62], Video-MAE V2 [57], and Perception Encoder [7] on downstream 3D tasks (Tab. 3-4), establishing E-RayZer as a strong paradigm for *spatial visual pre-training*.

- Compared with previous self-supervised 3D Vision models, E-RayZer showcases **stronger 3D understanding capability**, as evidenced by its significantly improved unsupervised camera pose estimation accuracy (Tab. 1) and 3D downstream task fine-tuning results (Tab. 3).

- Compared with state-of-the-art supervised models, *e.g.*, VGGT [55], E-RayZer achieves **on-par or sometimes superior performance** (Tab. 2) and exhibits *similar scaling patterns* (Tab. 5), despite being purely self-supervised.

## 2. Related Work

**Supervised Pose Estimation and 3D Reconstruction.** Early learning-based methods estimated relative camera poses from image pairs [3, 4, 9, 41], while later approaches explored multi-view reasoning across multiple inputs [21, 22, 32, 48, 54, 74, 75]. Given posed images, 3D representations can be reconstructed either by direct regression [21, 72, 76] or by optimization-based mode-seeking with diffusion models [78, 82]. Recent work has unified pose estimation and 3D reconstruction by predicting pixel-aligned pointmaps [12, 55, 58, 60, 79], exhibiting strong robustness under sparse inputs and generalizing well across diverse domains [52]. Nevertheless, training such supervised models still relies on camera pose and dense depth annotations, which are typically obtained from traditional SfM systems (*e.g.*, COLMAP [44]) and can be inaccurate, limiting performance of supervised models.

Recent work has also investigated predicting 3D Gaussians [29] with photometric losses as (part of the) supervision. However, these methods are still *de facto* supervised by 3D annotations, as they rely on ground-truth intrinsics [18, 26, 68] and/or target-view camera poses during training [26, 49, 68], or require initialization and/or regularizzation from 3D-supervised models [19, 24, 49]. In contrast, E-RayZer can be trained from scratch without any 3D supervision, and is therefore *truly self-supervised*, and can achieve even better performance.

**Self-supervised Novel-view Synthesis.** To alleviate the dependence on 3D supervision, another line of research investigates learning scene representations directly from 2D images using novel-view synthesis. Early works predict scene features from a single viewpoint and renders target

views as supervision [15, 30, 63, 80]. Recently, RUST [42], RayZer [23] and others [35, 53, 59] adopt learning-based latent rendering from multi-view inputs. However, these methods demonstrate limited 3D awareness, *e.g.*, RayZer learns view interpolation within an uninterpretable pose space. We build on RayZer but differ by adopting an explicit 3D representation (*i.e.*, 3D Gaussians [29]), a more fine-grained learning curriculum, and larger-scale training. We show that explicit 3D modeling leads to more geometrically grounded representations, establishing it as a promising pre-training framework for downstream tasks that require 3D understanding.

**Visual Pre-training for Representation Learning.** Prior works have made substantial progress in learning global image semantics by image-language association [1, 38, 51], learning 2D spatial priors via contrastive and completion losses [10, 16, 17], and via capturing temporal correlations with video-level self-supervision [5, 13, 50]. However, learning *3D-aware and geometrically grounded representations* remains underexplored, despite its strong potential to benefit 3D-related tasks where supervision is scarce. Recent efforts explore 3D awareness through *proxy tasks* of latent-space novel-view synthesis [23, 61, 62], but the degree to which these methods enforce true 3D understanding remains ambiguous. In this work, E-RayZer tackles the problem with explicit 3D modeling and introduces a learning curriculum that enables effective scaling, making the learned representations 3D-grounded and generalizable.

## 3. Approach

From unlabeled multi-view image sets, E-RayZer learns to predict camera (poses & intrinsics) and **explicit** 3D scene geometry under self-supervision. E-Rayzer's internal self-supervised representations can be further leveraged for downstream tasks, showing E-RayZer's potential as a 3D-aware visual pre-training framework.

In the following, we first revisit RayZer [23], the *implicit* predecessor, and discuss its limitations (Sec. 3.1). Building on RayZer's core design while addressing these issues by leveraging *Explicit* 3D modeling, we introduce *E*-RayZer (Sec. 3.2). Finally, we present a sequence-level curriculum learning strategy based on visual overlap between frames to improve performance and scalability (Sec. 3.3).

### 3.1. Preliminaries: RayZer with Implicit 3D

RayZer splits all input images into two *non-overlapping* subsets: an "observed" reference set ($\mathcal{I}_{\text{ref}}$) for latent scene inference, and a "hidden" target set ($\mathcal{I}_{\text{tgt}}$) for providing self-supervision. RayZer uses predicted cameras of target views ($\mathcal{I}_{\text{tgt}}$) to render the scene predicted from reference views ($\mathcal{I}_{\text{ref}}$), and applies photometric loss as self-supervision:

$$\mathcal{L} = \Sigma_{(I,\hat{I}) \in (\mathcal{I}_{\text{tgt}}, \hat{\mathcal{I}}_{\text{tgt}})} \big(\texttt{MSE}(I, \hat{I}) + \lambda \cdot \texttt{Percep}(I, \hat{I})\big), \quad (1)$$

where $\texttt{Percep}$ denotes perceptual loss [25].

RayZer leverages transformers for pose estimation, latent (implicit) scene reconstruction, and rendering. It first predicts camera intrinsics and extrinsics for all input images $\mathcal{I} \in \mathbb{R}^{V \times H \times W \times 3}$ using a multi-view transformer $f_{\boldsymbol{\theta}}^{\text{cam}}$, as:

$$(\mathbf{K}, \mathbf{T}) = f_{\boldsymbol{\theta}}^{\text{cam}}(\mathcal{I}), \quad \mathbf{T}_i = [\mathbf{R}_i \,|\, \mathbf{t}_i] \in SE(3), \quad (2)$$

where $\mathbf{K} \in \mathbb{R}^{3 \times 3}$ is the intrinsics shared by all views, $\mathbf{T} \in \mathbb{R}^{V \times 4 \times 4}$ denotes the extrinsics, and $i = 1, \ldots, V$ indexes the input images. Each camera $(\mathbf{K}, \mathbf{T}_i)$ is then converted into a pixel-aligned Plücker ray map $\mathbf{R}_i^{\text{plk}}$ [37, 75].

To infer latent scene representations, RayZer tokenizes the concatenation (along the feature dimension) of image and rays for $\mathcal{I}_{\text{ref}}$ and updates a set of learnable scene tokens $\mathbf{z}_0^{\text{scene}}$ through a transformer $f_{\boldsymbol{\psi}}^{\text{scene}}$, as:

$$\mathbf{z}_{\text{ref}}^{\text{scene}} = f_{\boldsymbol{\psi}}^{\text{scene}}\big(\mathbf{z}_0^{\text{scene}}, \text{Linear}(\mathcal{I}_{\text{ref}}, \mathbf{R}_{\text{ref}}^{\text{plk}})\big), \quad (3)$$

where $\text{Linear}(\cdot)$ denotes a patch-wise linear projection for fusing and tokenizing RGB and ray information. The resulting $\mathbf{z}_{\text{ref}}^{\text{scene}}$ represents the latent scene features.

For rendering, the self-predicted target-view Plücker ray maps are likewise tokenized and concatenated with the scene representation $\mathbf{z}_{\text{ref}}^{\text{scene}}$ (along the token dimension). These target-view ray tokens are refined via transformer $f_{\boldsymbol{\phi}}^{\text{rend}}$ and finally decoded to RGB images, as:

$$\hat{\mathcal{I}}_{\text{tgt}} = f_{\boldsymbol{\phi}}^{\text{rend}}\big(\mathbf{z}_{\text{ref}}^{\text{scene}}, \text{Linear}(\mathbf{R}_{\text{tgt}}^{\text{plk}})\big). \quad (4)$$

Then RayZer applies photometric self-supervision (Eq. 1).

**Limitations of RayZer's Implicit 3D.** RayZer achieves high-fidelity novel-view synthesis. However, **RayZer is not fully 3D-grounded**. Since its camera estimation ($f_{\boldsymbol{\theta}}^{\text{cam}}$), latent scene reconstruction ($f_{\boldsymbol{\psi}}^{\text{scene}}$), and rendering ($f_{\boldsymbol{\phi}}^{\text{rend}}$) modules are jointly learned from scratch, they only need to remain *mutually compatible*, but are not guaranteed to be physically or spatially meaningful. This issue is further amplified by RayZer's pure transformer-based architecture, which contains almost *no 3D inductive bias* and thus possesses excessive flexibility to learn undesirable shortcut solutions. As evidenced by its imperfect camera pose distribution, RayZer relies on a mixture of true 3D understanding and video-interpolation priors to achieve high-quality synthesis. While this design suffices for novel-view synthesis, it limits RayZer's potential as a *spatial pre-training* framework for learning genuinely 3D-aware representations.

### 3.2. E-RayZer: Explicit 3D with Self-supervision

**Our Insights.** We argue that **3D inductive biases remain essential** for 3D representation learning but they must be introduced correctly in ways that preserve **learning scalability**.
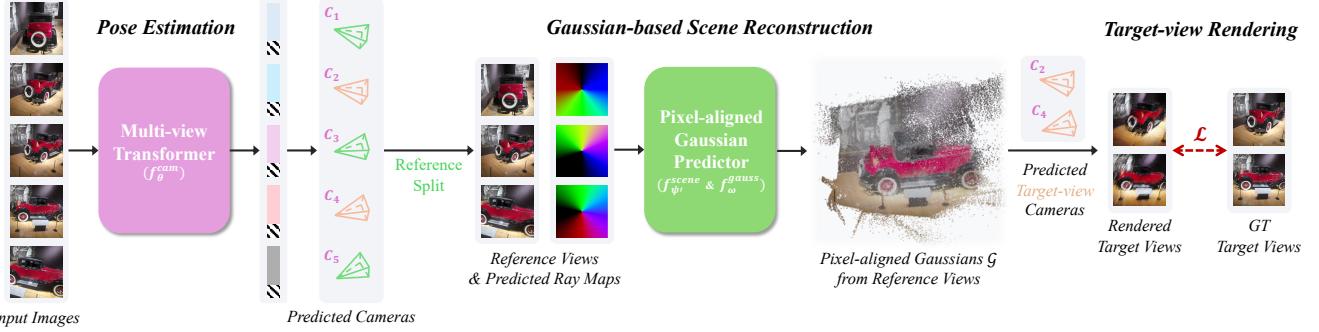
Figure 2. **E-RayZer Model & Training.** E-RayZer first predicts camera poses and intrinsics for all images. Then it follows RayZer [23] to split images into two sets. E-RayZer predicts explicit 3D Gaussians as scene representation from the reference views ($\mathcal{I}_{\text{ref}}$), and renders the scene using self-predicted target-view ($\mathcal{I}_{\text{tgt}}$) cameras. Finally, E-RayZer is trained with self-supervised photometric losses on target views.

Thus, we propose to inject *lightweight* 3D inductive bias through model design, while keeping the training fully self-supervised, striking a better balance between *3D awareness* and *scalability*. Specifically, E-RayZer replaces RayZer's latent scene representation with *explicit* 3D geometry (*i.e.*, 3D Gaussians [29]), providing *geometric regularization* to learn geometrically grounded pose estimation, scene reconstruction, and latent representations.

**Overview.** As shown in Fig. 2, E-RayZer first predicts the camera parameters for all images, and then infers pixel-aligned 3D Gaussians $\mathcal{G}_{\text{ref}}$ from the reference views subset ($\mathcal{I}_{\text{ref}}$). Then E-RayZer predicts the target views subset ($\mathcal{I}_{\text{tgt}}$), by rendering the 3D Gaussians predicted from $\mathcal{I}_{\text{ref}}$ under self-predicted cameras of $\mathcal{I}_{\text{tgt}}$. Since 3D Gaussians support closed-form differentiable rendering, the latent rendering decoder used in RayZer (*i.e.*, $f_\phi^{\text{rend}}$ in Eq. 4) is no longer required. We now describe our *key differences* from RayZer while elaborating on details.

**Gaussian-based Scene Reconstruction.** E-RayZer first predicts cameras of all views in a similar way with RayZer (besides differences in model architecture that will be detailed later). Then, E-RayZer directly transforms the "posed" reference views to pixel-aligned 3D Gaussians. We first encode posed reference views into latent tokens:

$$\mathbf{s}_{\text{ref}} = f_{\boldsymbol{\psi}'}^{\text{scene}}\big(\text{Linear}(\mathcal{I}_{\text{ref}}, \mathbf{R}_{\text{ref}}^{\text{plk}})\big) \tag{5}$$

where $\mathbf{s}_{\text{ref}} \in \mathbb{R}^{K_{\text{ref}}hw \times C}$ denotes the updated image tokens of reference views after multi-view aggregation. In detail, $K_{\text{ref}}$ is the number of views in $\mathcal{I}_{\text{ref}}$, $h = H/p$ and $w = W/p$ are token number along height and width dimensions using a patch size of $p$, and $C$ is channel dimension of the latent space. Note that the complexity of global attention in Eq. 5 is $\mathcal{O}((K_{\text{ref}}hw)^2)$, while it is $\mathcal{O}((K_{\text{ref}}hw + n_{\mathbf{z}})^2)$ for RayZer (Eq. 3), where $n_{\mathbf{z}}$ is the size for RayZer's scene token set.

Then, we use a lightweight decoder to transform the updated image tokens $\mathbf{s}_{\text{ref}}$ into per-pixel 3D Gaussian parame-

ters along each camera ray across all reference views, as:

$$\mathcal{G}_{\text{ref}} = f_{\boldsymbol{\omega}}^{\text{gauss}}(\mathbf{s}_{\text{ref}}), \quad \text{where}$$
$$\mathcal{G}_{\text{ref}} = \big\{\, g_i = (d_i, \mathbf{q}_i, \mathbf{C}_i, \mathbf{s}_i, \alpha_i) \,\big\}_{i=1}^{K_{\text{ref}} \times H \times W}. \tag{6}$$

These parameters include the distance along the ray $d_i \in \mathbb{R}$, orientation represented as a quaternion $\mathbf{q}_i \in \mathbb{R}^4$, spherical harmonic coefficients $\mathbf{C}_i \in \mathbb{R}^{(d_{\text{SH}}+1)^2 \times 3}$, scale $\mathbf{s}_i \in \mathbb{R}^3$, and opacity $\alpha_i \in \mathbb{R}$. The predicted 3D Gaussians collectively represent the scene geometry.

We then use E-RayZer's self-predicted target views cameras, denoted as $\mathcal{C}_{\text{tgt}} = \big\{\, (\mathbf{K}, \mathbf{T}_i) \mid i \in \mathcal{I}_{\text{tgt}} \,\big\}$, to render the 3D Gaussians $\mathcal{G}_{\text{ref}}$ and get prediction of target views, as:

$$\hat{\mathcal{I}}_{\text{tgt}} = \pi(\mathcal{G}_{\text{ref}}, \mathcal{C}_{\text{tgt}}), \tag{7}$$

where $\pi$ denotes the differentiable rendering equation of 3D Gaussians. Note that we modify gsplat [70] to support gradient back-propagation to camera intrinsics $\mathbf{K}$. Compared with RayZer, this design improves both *rendering efficiency* and *3D-awareness* by removing the need to learn a transformer-based renderer. Finally, we apply photometric loss on render target views as Eq. 1.

**Avoiding Undesirable View Interpolation.** As discussed in Sec. 3.1, RayZer tends to learn undesirable frame interpolation cues as shortcut solutions. We identify a main cause as its use of *image index embeddings* to associate image tokens with corresponding camera tokens for camera estimation, which provides a strong cue for learning interpolation.

In E-RayZer, we remove the image index embeddings entirely. We adopt a VGGT-style [55] multi-view transformer with alternating local-global attention, where the local attention boundary naturally defines the association relationship. Different from the original VGGT, E-RayZer performs *pairwise pose prediction*: camera tokens from a canonical view and a target view are concatenated to regress their relative camera pose. Consequently, E-RayZer does not require different camera register tokens for canonical and non-canonical views. This architectural design is ap-
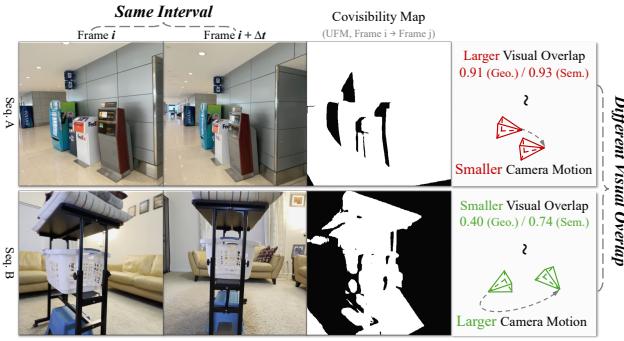
4

Figure 3. **Different Visual Overlaps under the Same Frame Interval.** Two sequences from DL3DV [33] share the same frame interval yet exhibit drastically different levels of visual overlap. Our proposed semantic and geometric overlap metrics more accurately capture the true difficulty (or camera motion) of each sequence.

plied to both the transformers used for camera estimation ($f_{\theta}^{\text{cam}}$) and that for scene reconstruction ($f_{\psi'}^{\text{scene}}$).

### 3.3. Sequence Curriculum Based on Visual Overlap

As E-RayZer leverages explicit scene representation, it suffers from harder convergence when **trained from scratch**. To stabilize training, we propose a learning curriculum based on the concept of *visual overlap* between input views, providing *fine-grained control* over training data difficulty. This curriculum also adaptively aligns the data distributions across diverse data sources, making E-RayZer more scalable to heterogeneous training resources.

We highlight that E-RayZer's learning curriculum fundamentally differs from that of RayZer, which is based on fixed frame-index intervals. As illustrated in Fig. 3, RayZer's interval-based sampling provides only an inaccurate and inflexible approximation of visual overlap, is hardcoded and thus not scalable to heterogeneous resources.

We then describe the two key steps for constructing our learning curriculum: *data labeling* and *sampling*. We then introduce two variants of visual-overlap labeling tools: a *geometric* version that computes actual covisibility, and a *semantic* version as an unsupervised approximation of it.

***Labeling.*** For each training sequence $u$ (from any data resource), we compute a spacing profile by *uniformly* sampling a small set of frame triplets for each spacing $\Delta t$, as $\mathcal{T}_u(\Delta t) = \{(i, i + \Delta t, i + 2\Delta t)\}$, and averaging the two pairwise overlaps $o(\cdot, \cdot)$ per triplet:

$$o_{\text{tri}}(i, \Delta t) = \frac{1}{2}\Big(o(i, i+\Delta t) + o(i+\Delta t, i+2\Delta t)\Big). \quad (8)$$

Averaging $o_{\text{tri}}(i, \Delta t)$ over all sampled triplets yields the per-sequence profile $O_u(\Delta t)$, characterizing how overlap (and consequently difficulty) varies with frame index spacing.

***Training-time Sampling.*** Given curriculum progress $s \in [0, 1]$, we use a visual overlap lower limit of $o(s) = s\, o_{\min} + (1-s)\, o_{\max}$, so that it decreases over training. We

then obtain the sequence-specific spacing $\Delta t_u(s)$ by *looking up* the precomputed table $\{(\Delta t_k, O_u(\Delta t_k))\}$ and *linearly interpolating* between the nearest entries. Finally, the sequence length follows $t = (V - 1)\,\Delta t_u(s)$.

***Instantiations.*** We instantiate $o$ with two alternatives – *geometric overlap* (UFM [77] covisibility, which is trained with 3D annotations) and *semantic overlap* (DINOv2 [36] cosine similarity, which is trained w. self-supervision):

$$\begin{aligned} o_{\text{sem}}(i, j) &= \cos\big(\phi_{\text{DINO}}(I_i), \phi_{\text{DINO}}(I_j)\big), \\ o_{\text{geo}}(i, j) &= \text{Cov}_{\text{UFM}}(I_i, I_j). \end{aligned} \quad (9)$$

In Sec. 4.4, we show that both the semantic and geometric curricula outperform RayZer's interval-based curriculum, and that the two variants perform comparably.

## 4. Experiments

We first describe the experimental setups in Sec. 4.1. We then evaluate E-RayZer in two aspects: as a self-supervised model for pose estimation and 3D reconstruction (Sec. 4.2), and as a spatial visual pre-training framework for downstream tasks (Sec. 4.3). Finally, we ablate the key design choices of E-RayZer (Sec. 4.4).

### 4.1. Experimental Setup

**Implementation Details.** E-RayZer is trained with 10 input images, where 5 are used as reference views and 5 as target views. During training, we follow a linear decay in visual-overlap scores: $1.0 \rightarrow 0.5$ for geometric-overlap scheduling and $1.0 \rightarrow 0.75$ for semantic-overlap scheduling. For a fair comparison, we align RayZer with E-RayZer using the better model architecture and the novel training curriculum. For other baselines, we use official checkpoints and provide specific implementation details in the corresponding subsections. See more details in the supplementary material.

**Metrics.** For pose estimation, we report relative pose accuracy (RPA) at thresholds of $5°$, $15°$, and $30°$, which jointly reflects rotation and translation accuracy. For novel-view synthesis, we use standard PSNR. For depth estimation, we evaluate absolute relative error (AbsRel) and $\delta < 1.25$, following Depth Anything [66]. For pairwise flow prediction, we report the average end-point error (EPE) and the proportion of outlier flow predictions under thresholds of 1px, 2px, and 5px, following UFM [77].

**Datasets.** *Training.* We present results of E-RayZer trained on both single-dataset and multi-dataset settings. The single-dataset variants are trained exclusively on RealEstate10K [43] or DL3DV [33], while the multi-dataset variant is trained on a mixture of seven datasets: DL3DV [33], CO3Dv2 [40], RealEstate10K [81], MVImgNet [73], ARKitScenes [6], WildRGB-D [64], and ACID [34], covering diverse indoor and outdoor sequences.

*Evaluation.* We primarily evaluate pose estimation and novel-view synthesis on WildRGB-D, DL3DV test set, and
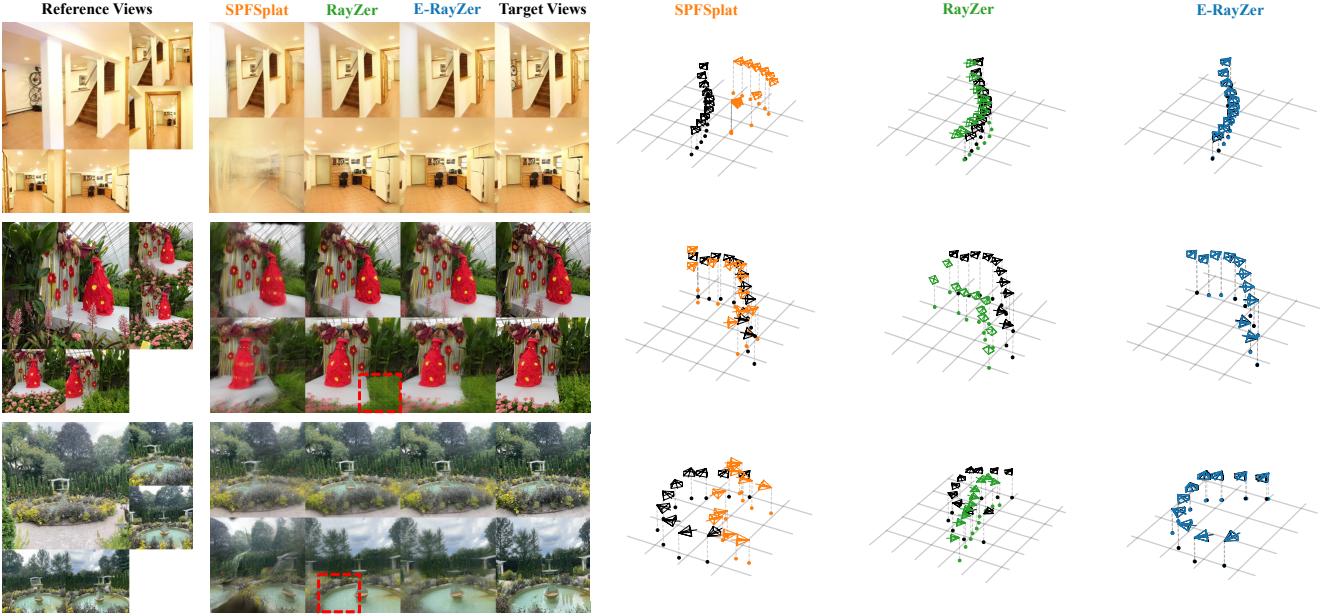
Figure 4. **Visual Comparison with (Partially) Self-supervised Methods.** We include results on both novel view synthesis (left) and pose estimation (right), where E-RayZer outperforms baselines on pose accuracy, showing its grounded 3D understanding. E-RayZer also outperforms RayZer on low-texture regions (highlighted w/ red box) on NVS, a case where RayZer's view interpolation cannot handle.

Table 1. **Comparison with (Partially) Self-supervised Methods on Novel-view Synthesis (NVS) and Pose Estimation.** We report PSNR for NVS and RPA↑@5°/15°/30° for pose estimation. RayZer [23] and E-RayZer are fully self-supervised methods trained from scratch, while SPFSplat [18] is initialized from MASt3R [12], which itself is trained under dense 3D supervision on 14 datasets.

| Method | Self-supervised? | Training Data | WildRGB-D [64] | | | | ScanNet++ [71] | | | | DL3DV [33] | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | PSNR↑ | @5°↑ | @15°↑ | @30°↑ | PSNR↑ | @5°↑ | @15°↑ | @30°↑ | PSNR↑ | @5°↑ | @15°↑ | @30°↑ |
| SPFSplat [19] | ✗ (MASt3R ini.) | RE10K [81] (+ extra) | 16.7 | 31.5 | 58.0 | 69.8 | 14.0 | **2.5** | 11.8 | 30.3 | 15.1 | 19.5 | 40.6 | 50.5 |
| E-RayZer (ours) | ✓ | RE10K [81] | **21.0** | **40.3** | **89.4** | **96.5** | **17.5** | 1.1 | **13.3** | **37.3** | **17.3** | **21.2** | **55.0** | **72.7** |
| RayZer [23] | ✓ | DL3DV [33] | **25.9** | 0.0 | 0.2 | 6.5 | **20.5** | 0.0 | 0.7 | 6.2 | **21.4** | 0.0 | 0.6 | 6.2 |
| E-RayZer (ours) | ✓ | | 24.3 | **84.5** | **98.4** | **99.3** | 20.1 | **7.7** | **33.6** | **63.0** | 20.3 | **72.0** | **88.4** | **93.5** |
| RayZer [23] | ✓ | 7 datasets | **26.7** | 0.2 | 9.3 | 43.6 | **21.5** | 0.0 | 0.9 | 9.0 | **20.8** | 0.0 | 1.9 | 17.0 |
| E-RayZer (ours) | ✓ | | 24.9 | **90.8** | **98.6** | **99.3** | 20.7 | **5.7** | **34.8** | **63.7** | 19.7 | **59.9** | **82.9** | **90.2** |

the out-of-distribution (OOD) ScanNet++ [71]. To assess the generalization of the learned representations (Sec. 4.3), we evaluate on OOD ScanNet++ and BlendedMVS [67] for pose and depth estimation, and StaticThings3D [45] for pairwise flow prediction.

## 4.2. Pose Estimation and Novel-view Synthesis

**Baselines and Setups.** We compare against SPFSplat [19] and RayZer [23]. Notably, SPFSplat is initialized from the supervised MASt3R [31] model, and thus is not truly self-supervised; while E-RayZer and RayZer are trained from scratch under self-supervision. We evaluate pose accuracy on all images and assess novel-view synthesis quality on the target views rendered with predicted camera poses.

**Results.** As shown in Tab. 1, E-RayZer consistently outperforms SPFSplat [19] on most metrics, despite being truly self-supervised. Moreover, E-RayZer significantly surpasses RayZer [23] in pose estimation under all setups and achieve comparable novel-view synthesis quality. The results suggest that the *explicit* 3D modeling strategy of E-

RayZer leads to more geometrically meaningful pose representations, whereas RayZer's *implicit* method is overly optimized for high-quality view synthesis and is not truly 3D-aware, making the pose space less interpretable. The numbers are also verified by the visuals in Fig. 4.

## 4.3. E-RayZer as Self-supervised Pre-training

We validate E-RayZer as a self-supervised spatial visual pre-training framework. First, we show that its performance is comparable to the supervised VGGT and that E-RayZer pre-training further enhances VGGT (Sec. 4.3.1). We then probe the learned features on downstream tasks to verify E-RayZer's representation quality (Sec. 4.3.2).

### 4.3.1. E-RayZer Benefits Supervised Model

**Baselines and Setups.** We compare with the state-of-the-art supervised model VGGT [55]. Note that we train it using the same data and architecture with E-RayZer for an apple-to-apple comparison, denoted as VGGT*.

**E-RayZer is Comparable with Supervised VGGT*.** First

Table 2. **Comparison with Supervised VGGT [55] on Pose Estimation. E-RayZer's pre-training improves VGGT performance** (last row), forming an effective self-supervised pre-training and supervised post-training paradigm. We report pose accuracy RPA↑@5°/15°. Both models are trained on DL3DV [33] and evaluated on DL3DV & eight out-of-domain datasets for zero-shot testing. Models are labeled as self-supervised or supervised. VGGT* denotes our re-implement. with E-RayZer's pairwise camera head. Results are color-ranked from red to yellow, and we <u>underline</u> the results that our self-supervised E-RayZer surpasses supervised VGGT*. See Tab. 8 for more results.

| Method | In-domain DL3DV [40] | | RE10K [81] | | CO3Dv2 [40] | | WildRGB-D [64] | | 7-Scenes [46] | | CamLand [27] | | BlendedMVS [67] | | NAVI [20] | | ScanNet++ [71] | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | @5°↑ | @15°↑ | @5°↑ | @15°↑ | @5°↑ | @15°↑ | @5°↑ | @15°↑ | @5°↑ | @15°↑ | @5°↑ | @15°↑ | @5°↑ | @15°↑ | @5°↑ | @15°↑ | @5°↑ | @15°↑ |
| E-RayZer (ours) | 72.0 | 88.4 | <u>83.0</u> | 96.8 | <u>19.1</u> | 61.8 | <u>51.1</u> | <u>82.3</u> | <u>38.8</u> | 78.0 | <u>18.1</u> | <u>62.9</u> | 22.9 | <u>46.8</u> | <u>20.7</u> | <u>57.8</u> | <u>7.7</u> | 33.6 |
| VGGT* | 79.6 | 94.2 | 80.4 | 97.9 | 16.0 | 64.3 | 32.5 | 76.2 | 34.7 | **83.6** | 11.1 | 49.8 | 17.0 | 42.8 | 14.3 | 54.5 | 6.7 | 39.8 |
| E-RayZer→VGGT* | 87.3 | 96.6 | 85.3 | 98.4 | 25.3 | 72.2 | 56.2 | 91.4 | 43.8 | 82.8 | 30.2 | 75.6 | 29.2 | 52.2 | 26.9 | 64.3 | 14.3 | 53.8 |

Table 3. **Probing 3D Spatial Awareness of Learned Representations on Multi-view Depth and Pose Estimation.** We evaluate the learned representations via both frozen-backbone and fully supervised finetuning on ScanNet++ [71] and BlendedMVS [67], which are not included in pre-training for any model. The best results are shown in **bold**, and the second-best are <u>underlined</u>. The experiments only use the *encoders* of RayZer [23] and E-RayZer.

| | | Method | Depth | | Camera Pose | |
|---|---|---|---|---|---|---|
| | | | AbsRel↓ | δ<1.25↑ | RPA@5°↑ | RPA@15°↑ |
| ScanNet++ [71] | Frozen | DINOv2 [36] | 0.193 | 74.9 | 0.8 | 9.5 |
| | | DINOv3 [47] | 0.201 | 73.2 | 0.4 | 10.0 |
| | | Percep. Encoder [7] | 0.203 | 73.2 | 0.5 | 8.5 |
| | | CroCo v2 [62] | 0.203 | 73.0 | 1.4 | 15.1 |
| | | VideoMAE V2 [57] | 0.175 | 76.3 | 0.1 | 6.6 |
| | | RayZer [23] | <u>0.161</u> | <u>79.3</u> | <u>4.7</u> | <u>27.4</u> |
| | | E-RayZer (ours) | **0.116** | **87.1** | **13.8** | **49.5** |
| | Full-finetune | DINOv2 [36] | 0.178 | 78.2 | 3.3 | 19.6 |
| | | DINOv3 [47] | 0.176 | 78.7 | 4.0 | 22.3 |
| | | Percep. Encoder [7] | 0.181 | 77.8 | 2.9 | 20.0 |
| | | CroCo v2 [62] | 0.177 | 78.2 | 3.8 | 20.8 |
| | | VideoMAE V2 [57] | <u>0.076</u> | <u>93.9</u> | 12.8 | 51.4 |
| | | RayZer [23] | 0.077 | <u>93.9</u> | <u>21.5</u> | <u>60.6</u> |
| | | E-RayZer (ours) | **0.059** | **95.1** | **22.7** | **64.3** |
| BlendedMVS [67] | Frozen | DINOv2 [36] | 0.366 | 50.5 | 1.1 | 8.0 |
| | | DINOv3 [47] | 0.397 | 49.1 | 1.2 | 6.8 |
| | | Percep. Encoder [7] | 0.385 | 49.9 | 1.2 | 6.2 |
| | | CroCo v2 [62] | 0.412 | 47.7 | 1.6 | 12.6 |
| | | VideoMAE V2 [57] | 0.371 | 49.4 | 1.0 | 6.2 |
| | | RayZer [23] | <u>0.351</u> | <u>52.6</u> | <u>16.7</u> | <u>34.5</u> |
| | | E-RayZer (ours) | **0.245** | **68.3** | **26.5** | **45.8** |
| | Full-finetune | DINOv2 [36] | 0.353 | 52.5 | 1.8 | 12.8 |
| | | DINOv3 [47] | 0.349 | 52.1 | 1.7 | 15.3 |
| | | Percp. Encoder [7] | 0.370 | 50.3 | 2.1 | 11.6 |
| | | CroCo v2 [62] | 0.369 | 51.2 | 2.8 | 15.9 |
| | | VideoMAE V2 [57] | 0.197 | 75.9 | 17.3 | 45.5 |
| | | RayZer [23] | <u>0.194</u> | <u>77.7</u> | <u>26.1</u> | <u>50.2</u> |
| | | E-RayZer (ours) | **0.148** | **82.8** | **36.2** | **58.8** |

Table 4. **Probing 2.5D Spatial Awareness of Learned Representations on Pairwise Flow Estimation.** We evaluate on Static-Things3D [45], an out-of-distribution synthetic dataset. All models are fully finetuned under flow supervision. The best results are shown in **bold**, and the second-best are <u>underlined</u>.

| Method | Error | Outlier Ratio | | |
|---|---|---|---|---|
| | EPE↓ | @1px↓ | @2px↓ | @5px↓ |
| CroCo v2 [62] | 1.273 | 17.7 | 8.7 | 3.8 |
| VideoMAE V2 [57] | 2.028 | 42.7 | 22.1 | 6.9 |
| RayZer [23] | **1.105** | **13.4** | **6.6** | **2.8** |
| E-RayZer (ours) | <u>1.254</u> | <u>16.9</u> | <u>7.8</u> | <u>3.1</u> |

E-RayZer almost consistently achieves higher accuracy on RPA@5°, a stricter metric, suggesting better precision in pose prediction. The results demonstrate the strong performance of E-RayZer as a self-supervised method without using any 3D annotations for training.

**Effectiveness of Pre-training.** As shown in last two rows of Tab. 2, initializing VGGT* with E-RayZer weights yields significant improvements over training from scratch, confirming that E-RayZer serves as an effective pre-training framework for visual geometry learning. The results also suggest that the learned knowledge of our self-supervised and supervised methods are highly complementary (they are trained on same data but pre-training still helps), showing the great potential of spatial visual pre-training.

### 4.3.2. Probing Representations on Downstream Tasks

**Baselines and Setups.** To further assess the spatial awareness, we probe and compare the feature representations of E-RayZer against widely-used vision encoders: DINO-series [36, 47], CroCo v2 [62], VideoMAE V2 [57], Perception Encoder [7], and RayZer [23]. We only use the backbones and train the prediction heads **from scratch**. We compare performance under both frozen-backbone and full-finetuning settings on downstream tasks, including:

• *Multi-view Depth and Pose Estimation (3D Tasks).* For depth estimation, we apply a DPT head [39] on top of the backbones. For pose estimation, we attach VGGT's [55] camera head to each backbone, using either the class token or averaged patch features as camera tokens. These tokens are aggregated across views via transformer layers, enabling even single-view models to reason over multi-view geometry. We note the camera estimation heads of RayZer and E-RayZer in their pre-training stage are not used.

• *Pairwise Flow Estimation (2.5D Task).* We consider backbones that encode binocular geometry, including CroCo v2 [62], VideoMAE V2 [57], RayZer [23], and E-RayZer. We follow the settings of UFM [77].

**Results on 3D Downstream Tasks.** Tab. 3 shows that E-RayZer achieves the best performance across all datasets and settings, demonstrating strong 3D-awareness in its feature representations. Under the frozen-backbone setting, E-RayZer notably outperforms all baselines. With full finetuning, E-RayZer further improves across all metrics, surpass-

two rows of Tab. 2 show that E-RayZer outperforms VGGT* on several out-of-domain datasets (*e.g.*, WildRGB-D [64], CamLand [27], and BlendedMVS [67]). Moreover,

Table 5. **Ablation on Data Mixing and Scaling.** We compare our E-RayZer with supervised VGGT* [55] on varying training data settings. We color-rank the results from red to yellow **for each model itself** across training data, thus **the color distribution reflect their scaling behavior**. We also <u>underline</u> the results where self-supervised E-RayZer outperforms supervised VGGT* (for each training data).

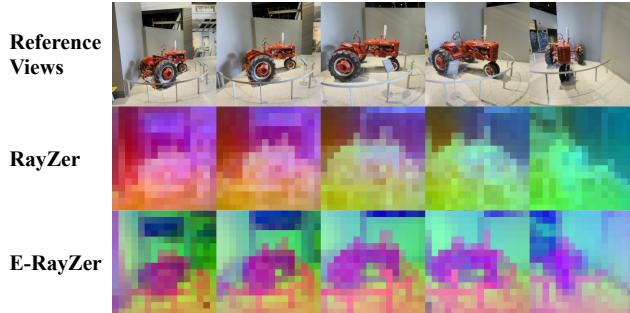| Training Data | Method | NAVI [20] | | | | CO3Dv2 [40] | | | | ScanNet++ [71] | | | | DL3DV [33] | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | PSNR↑ | @5°↑ | @15°↑ | @30°↑ | PSNR↑ | @5°↑ | @15°↑ | @30°↑ | PSNR↑ | @5°↑ | @15°↑ | @30°↑ | PSNR↑ | @5°↑ | @15°↑ | @30°↑ |
| RE10K [81] | VGGT* | / | 0.4 | 8.4 | 22.5 | / | 0.1 | 3.7 | 15.5 | / | 0.6 | 10.0 | 30.7 | / | 17.8 | 50.9 | 69.4 |
| | E-RayZer | 17.2 | <u>1.8</u> | <u>16.9</u> | <u>34.0</u> | 19.1 | <u>0.6</u> | <u>8.3</u> | <u>26.0</u> | 17.5 | <u>1.1</u> | <u>13.3</u> | <u>37.3</u> | 17.3 | <u>21.2</u> | <u>55.0</u> | <u>72.7</u> |
| DL3DV [33] | VGGT* | / | 14.3 | 54.5 | 75.7 | / | 16.0 | 64.3 | 82.1 | / | 6.7 | 39.8 | 71.5 | / | 79.6 | 94.2 | 97.1 |
| | E-RayZer | 20.5 | <u>20.7</u> | <u>57.8</u> | 69.6 | 22.9 | <u>19.1</u> | 61.8 | 78.8 | 20.1 | <u>7.7</u> | 33.6 | 63.0 | 20.3 | 72.0 | 88.4 | 93.5 |
| 7-dataset Mix | VGGT* | / | 28.8 | 67.3 | 84.4 | / | 43.4 | 83.5 | 91.8 | / | 13.1 | 54.8 | 78.5 | / | 66.1 | 88.9 | 95.6 |
| | E-RayZer | 20.6 | 24.6 | 56.1 | 69.2 | 24.3 | 30.3 | 74.2 | 83.7 | 20.7 | 5.7 | 34.8 | 63.7 | 19.7 | 59.9 | 82.9 | 90.2 |



**Reference Views**

**RayZer**

**E-RayZer**

Figure 5. **Comparison with RayZer [23] on Learned Features**, visualized with their top-3 PCA components. The feature maps produced by E-RayZer exhibit more pronounced and spatially consistent patterns aligned with the main scene structures (*e.g.*, the tractor, the surrounding curved metal railing, and the wall).

ing RayZer [23] and VideoMAE V2 [57] by a large margin. The consistently strong results highlight the generalization ability of its geometrically grounded representations, showing its potential as a pre-training framework.

**Results on Pairwise Flow Estimation.** Tab. 4 shows that E-RayZer achieves competitive performance on pairwise flow prediction, closely following RayZer [23], despite not being trained directly for tasks that optimize image correspondences (*e.g.*, masked image modeling in CroCo v2 [62] and VideoMAE V2 [57], or view interpolation in RayZer). Compared to E-RayZer, RayZer holds a slight advantage due to its implicit 3D formulation, naturally suited for low-level motion estimation. Nevertheless, E-RayZer outperforms other baselines, demonstrating that its explicit 3D representation learning captures meaningful spatial correspondences even for 2.5D tasks.

**Visualization.** Fig. 5 shows the multi-view features for RayZer and E-RayZer. We observe that the features from E-RayZer more clearly capture the major 3D scene structures and remain consistent across different views.

## 4.4. Ablation Study

**Data Mixing / Scaling.** We investigate the behavior of self-supervised E-RayZer and supervised VGGT* (Sec. 4.3.1) under varying data scales and quality. In Tab. 5, E-RayZer and VGGT* demonstrate a similar scaling behavior: training on data with broader distributions improves generalization (*e.g.*, models trained on 7 datasets outperform those

Table 6. **Ablation on Curriculum Learning.** We compare four curriculum strategies when training E-RayZer on DL3DV (top) and a seven-dataset mixture (bottom). The proposed visual-overlap-based curriculum consistently outperforms baselines.

| | Curriculum Variant | PSNR↑ | RPA@5°↑ | RPA@15°↑ | RPA@30°↑ |
|---|---|---|---|---|---|
| DL3DV | No Curriculum | 16.1 | 4.0 | 27.8 | 47.2 |
| | Frame Interval | 19.8 | 56.1 | 79.3 | 86.0 |
| | Semantic Overlap | **20.4** | **73.2** | **88.7** | **93.7** |
| | Geometric Overlap | <u>20.3</u> | <u>72.0</u> | 88.4 | <u>93.5</u> |
| 7-dataset | No Curriculum | 15.9 | 2.1 | 21.6 | 40.7 |
| | Frame Interval | 19.1 | 43.8 | 72.1 | 82.9 |
| | Semantic Overlap | **19.7** | <u>58.7</u> | <u>81.0</u> | <u>89.8</u> |
| | Geometric Overlap | 19.7 | **59.9** | **82.9** | **90.2** |

trained on DL3DV alone). However, reducing the sampling frequency of a particular domain slightly degrades performance on its corresponding test set (*e.g.*, 7-dataset models perform worse on DL3DV than DL3DV-only models), a trend consistently observed in prior work [14, 65, 69]. Besides, data quality also plays a key role, as training on DL3DV yields better results than that on RE10K.

Moreover, again, the self-supervised model (E-RayZer) achieves performance on par with the supervised VGGT* (while VGGT* holds advantage when trained on large data), demonstrating that large-scale self-supervision alone can yield geometrically grounded 3D understanding. This result underscores that data diversity and quality, rather than explicit 3D supervision, are the true drivers of scalability in large 3D Vision models. Together, these results highlight the great potential of self-supervised 3D learning when scaled to internet-scale data, and provide valuable guidance for future data selection and curation strategies.

**Curriculum Learning.** In Tab. 6, we compare against two baselines with (1) no curriculum, and (2) a frame-interval-based curriculum, where frame intervals are specified for each dataset. Across two training regimes (*i.e.*, DL3DV-only and the seven-dataset mixture), the proposed visual-overlap curricula consistently outperform both baselines, with the two variants performing comparably. These results demonstrate that our fine-grained curriculum strategy significantly improves self-supervised pose estimation and reconstruction, while eliminating the need for manual tuning for each training dataset and benefiting scaling.

## 5. Conclusion

We propose E-RayZer, a multi-view 3D model for learning geometrically grounded representations via self-supervised 3D reconstruction. E-RayZer demonstrates better performance against prior unsupervised methods and is even comparable with supervised methods. Extensive experimental results demonstrate E-RayZer pre-training benefits supervised models and other 3D downstream tasks, establishing it as a scalable 3D-aware visual pre-training framework.

# E-RayZer: Self-supervised 3D Reconstruction as Spatial Visual Pre-training

## Supplementary Material

## Overview

This supplementary material is organized as follows:

## A. Additional Implementation Details

This section includes more implementation details.

**Training.** E-RayZer is trained on 8 A100 GPUs with a global batch size of 192 (24 per GPU) for 152K iterations. During the first 86K iterations, the learning curriculum progresses linearly according to different metrics, *i.e.*, geometric (default) and semantic visual overlap, as well as the frame intervals described in Sec. 4.4. Our learning rate (LR) schedule includes a 3K-iteration linear warm-up (peak LR of 4e-4), followed by a cosine decay. We use the AdamW optimizer ($\beta_1$=0.9, $\beta_2$=0.95) and apply gradient clipping at 1.0. We further skip optimization steps if the gradient norm exceeds 5.0 before clipping.

For our 7-dataset model (Sec. 4.1), we train on a mixture of datasets with the following sampling ratios: DL3DV [33]: 1.0, CO3Dv2 [40]: 0.25, RealEstate10K [81]: 0.5, MVImgNet [73]: 0.25, ARKitScenes [6]: 0.5, WildRGB-D [64]: 0.25, and ACID [34]: 0.5. These ratios follow a simple heuristic: we downweight object-centric datasets and assign a slightly larger weight to DL3DV, which offers the most diverse and high-quality samples.

Experiments on supervised finetuning are conducted on 8 A100 GPUs as well, but with a smaller global batch size of 96. The finetuning stage runs for 50K iterations.

**Architecture.** E-RayZer uses a patch size of 16 and an image resolution of 256. As described in Sec. 3.2, we replace RayZer's [23] vanilla global attention with VGGT's [55] local-global alternating transformer layers for both pose estimation ($f_\theta^{cam}$) and scene reconstruction ($f_{\psi'}^{scene}$). Both modules use 8 layers, each composed of one global attention layer and one frame-attention layer. Our feature dimension is 768, and we use 12 attention heads. For image and Plücker ray map tokenization, as well as for the Gaussian decoder ($f_\omega^{gauss}$), we simply use a single linear layer.

For a fair comparison with RayZer, all RayZer models used in this paper are trained with our proposed curriculum and the improved architecture.

**Evaluation.** For pose estimation and novel-view synthesis, we use fixed sequence lengths for the test sequences of each dataset and sample views with equal temporal spacing. Following RayZer, we ensure that the first and last images of each sequence are always included in the reference set. The sequence lengths are as follows: WildRGB-D [64]: 96 (Tab. 1) and 192 (Tab. 2), ScanNet++ [71]: 48, DL3DV [33]: 96, RealEstate10K [81]: 256, CO3Dv2 [40]: 96, 7-Scenes [46]: 256, Cambridge Landmarks [28]: 96, BlendedMVS [67]: 24, and NAVI [20]: 24. For (training and) evaluating pairwise flow prediction on Static-Things3D [45], we adopt the pre-computed image pairs provided by the DUSt3R [60] GitHub repository.

## B. More Details on Supervised Finetuning

Here we provide additional details on the supervised finetuning experiments in Sec. 4.3.

**Supervised Finetuning with E-RayZer.** E-RayZer's backbone does not distinguish between the first view and the other views in the input, as it adopts a pairwise pose estimation strategy (see Sec. 3.2). In contrast, supervised pose estimation typically assumes a first-view coordinate frame (*e.g.*, DUSt3R [60] and VGGT [55]). To incorporate this inductive bias into our backbone, we introduce an additional camera token dedicated to the first image (in addition to the existing learned camera token) and train it from scratch. The camera tokens are processed by E-RayZer's pose estimation module ($f_\theta^{cam}$) and subsequently passed to VGGT's camera head for supervised pose estimation. For depth estimation and pairwise flow prediction, the DPT head takes as input the intermediate feature maps generated by the Gaussian-based scene reconstruction module ($f_{\psi'}^{scene}$). For E-RayZer and all other baselines, the DPT head uses four feature maps extracted from equally spaced transformer layers. Note that our Gaussian-based scene reconstruction module takes the predicted reference-view Plücker ray maps as input, but only in the pose and depth estimation experiments are the predicted camera poses supervised. For pairwise flow prediction, the predicted poses produced by the pose head remain unsupervised to ensure a fair comparison with other baselines.

**Details on Other Baselines.** For baselines that use different spatial or temporal patch sizes (*e.g.*, E-RayZer uses a temporal batch size of 1, whereas VideoMAE V2 [57] uses 2),

Table 7. **Comparison with a Pose-supervised Baseline on Novel-view Synthesis (NVS) and Pose Estimation.** We report PSNR for NVS and RPA$_\uparrow$@5°/15°/30° for pose estimation. While the pose-supervised baseline generally outperforms the self-supervised model on coarse pose accuracy (RPA$_\uparrow$@15°/30°), its novel-view synthesis quality is consistently lower.

| Method | Training Data | NAVI [20] | | | | ScanNet++ [71] | | | | DL3DV [33] | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | PSNR$_\uparrow$ | @5°$_\uparrow$ | @15°$_\uparrow$ | @30°$_\uparrow$ | PSNR$_\uparrow$ | @5°$_\uparrow$ | @15°$_\uparrow$ | @30°$_\uparrow$ | PSNR$_\uparrow$ | @5°$_\uparrow$ | @15°$_\uparrow$ | @30°$_\uparrow$ |
| Pose-sup. Baseline | DL3DV [33] | 13.4 | 12.8 | 51.1 | **72.5** | 16.7 | 4.4 | **33.7** | **64.5** | 15.0 | **78.1** | **94.7** | **97.8** |
| E-RayZer (ours) | | **20.5** | **20.7** | **57.8** | 69.6 | **20.1** | **7.7** | 33.6 | 63.0 | **20.3** | 72.0 | 88.4 | 93.5 |
| Pose-sup. Baseline | 7 datasets | 13.5 | 18.9 | **61.6** | **80.6** | 17.3 | **6.4** | **35.7** | **67.4** | 14.9 | 53.0 | **85.0** | **93.2** |
| E-RayZer (ours) | | **20.6** | **24.6** | 56.1 | 69.2 | **20.7** | 5.7 | 34.8 | 63.7 | **19.7** | **59.9** | 82.9 | 90.2 |

Table 8. **Comparison with RayZer [23] as a Pre-trained Backbone.** The top block reports results for models trained on DL3DV [33], and the bottom block reports results for models trained on a mixture of seven datasets. Note that pre-training and supervised finetuning are performed on the same data (*i.e.*, DL3DV or the 7-dataset mixture). We report pose accuracy RPA$_\uparrow$@5°/15°. Models are labeled as self-supervised or supervised. VGGT* denotes our re-implementation with E-RayZer's pairwise camera head. The top-three results are color-ranked from red to yellow. E-RayZer provides stronger pre-training than RayZer.

| | Method | DL3DV [33] | | RE10K [81] | | CO3Dv2 [40] | | WildRGB-D [64] | | 7-Scenes [46] | | CamLand [27] | | BlendedMVS [67] | | NAVI [20] | | ScanNet++ [71] | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | @5°$_\uparrow$ | @15°$_\uparrow$ | @5°$_\uparrow$ | @15°$_\uparrow$ | @5°$_\uparrow$ | @15°$_\uparrow$ | @5°$_\uparrow$ | @15°$_\uparrow$ | @5°$_\uparrow$ | @15°$_\uparrow$ | @5°$_\uparrow$ | @15°$_\uparrow$ | @5°$_\uparrow$ | @15°$_\uparrow$ | @5°$_\uparrow$ | @15°$_\uparrow$ | @5°$_\uparrow$ | @15°$_\uparrow$ |
| DL3DV | RayZer [23] | 0.0 | 0.6 | 0.0 | 0.2 | 0.0 | 0.6 | 0.0 | 0.0 | 0.0 | 0.2 | 0.0 | 0.3 | 0.0 | 0.5 | 0.0 | 0.6 | 0.0 | 0.7 |
| | E-RayZer (ours) | 72.0 | 88.4 | 83.0 | 96.8 | 19.1 | 61.8 | 51.1 | 82.3 | 38.8 | 78.0 | 18.1 | 62.9 | 22.9 | 46.8 | 20.7 | 57.8 | 7.7 | 33.6 |
| | VGGT* | 79.6 | 94.2 | 80.4 | 97.9 | 16.0 | 64.3 | 32.5 | 76.2 | 34.7 | 83.6 | 11.1 | 49.8 | 17.0 | 42.8 | 14.3 | 54.5 | 6.7 | 39.8 |
| | RayZer→VGGT* | 84.4 | 95.3 | 85.7 | 98.4 | 24.9 | 71.2 | 43.9 | 86.4 | 38.0 | 83.6 | 27.3 | 73.0 | 24.0 | 45.8 | 25.5 | 58.3 | 12.2 | 49.6 |
| | E-RayZer→VGGT* | 87.3 | 96.6 | 85.3 | 98.4 | 25.3 | 72.2 | 56.2 | 91.4 | 43.8 | 82.8 | 30.2 | 75.6 | 29.2 | 52.2 | 26.9 | 64.3 | 14.3 | 53.8 |
| 7 datasets | RayZer [23] | 0.0 | 1.9 | 0.0 | 0.9 | 0.0 | 1.6 | 0.0 | 1.1 | 0.0 | 2.0 | 0.0 | 0.6 | 0.0 | 1.6 | 0.0 | 1.6 | 0.0 | 0.9 |
| | E-RayZer (ours) | 59.9 | 82.9 | 84.1 | 97.5 | 30.3 | 74.2 | 63.1 | 85.3 | 26.0 | 76.5 | 9.8 | 47.3 | 22.3 | 45.5 | 24.6 | 56.1 | 5.7 | 34.8 |
| | VGGT* | 66.1 | 88.9 | 85.2 | 98.5 | 43.4 | 83.5 | 76.8 | 96.0 | 31.1 | 78.0 | 22.9 | 66.3 | 19.0 | 49.9 | 28.8 | 67.3 | 13.1 | 54.8 |
| | RayZer→VGGT* | 72.8 | 91.7 | 88.1 | 98.6 | 53.8 | 85.1 | 81.5 | 96.3 | 37.7 | 84.9 | 28.3 | 65.7 | 24.3 | 52.7 | 34.6 | 70.4 | 15.0 | 58.7 |
| | E-RayZer→VGGT* | 78.8 | 92.8 | 91.0 | 99.1 | 58.9 | 86.3 | 86.4 | 96.7 | 42.7 | 88.3 | 35.2 | 64.4 | 31.5 | 57.7 | 41.5 | 73.7 | 22.0 | 65.2 |

we first resize or repeat the input so that the number of output tokens matches that of our model. For these methods, we generally adopt the "base" model checkpoints provided in their official GitHub repositories, as they roughly match the computational budget of our model.

# C. Additional Details on Curriculum Ablation

In this section, we provide additional details on the baseline setups used in Tab. 6. We compare our visual-overlap-based curricula to two baseline strategies: (1) Non-curriculum baseline, where we do not progressively increase the difficulty of training samples. Concretely, the geometric visual-overlap score remains fixed within the range [0.5, 1.0] throughout training, without any linear decay. As a result, the model encounters challenging samples (*e.g.*, wide-baseline views) from the very beginning. (2) Frame-interval-based curriculum, where geometric-overlap scores are converted into frame intervals that linearly increase over training. To construct the interval schedule for each dataset, we pre-sample 10K sequences with geometric-overlap scores in [0.5, 1.0] and set the maximum frame interval to the 95th percentile of these sequences. This heuristic implicitly defines dataset-specific hyperparameters that would otherwise need to be *manually tuned*.

# D. A Pose-supervised Baseline

We introduce a pose-supervised baseline whose pose estimation module is trained using ground-truth camera poses (typically obtained from running Structure-from-Motion systems [44]), following prior supervised methods (*e.g.*, DUSt3R [60] and VGGT [55]). In this baseline, the Gaussian-based scene reconstruction module is still optimized with a photometric loss; however, gradients from this loss are not propagated back to the pose estimation module. The results are shown in Tab. 7.

We observe that while the pose-supervised baseline usually outperforms E-RayZer on coarse pose accuracy (RPA@15°/30°), it consistently achieves lower PSNR for novel-view synthesis. We attribute this weaker NVS performance to a misalignment between the predicted poses and the Gaussian prediction. To supervise pose estimation, the ground-truth camera poses are normalized to a predefined scale (*e.g.*, 1.0), and the pose estimation module learns to predict camera poses at this scale. However, the Gaussian prediction module does not necessarily follow the same scale. In practice, we observe many training instances where the rendered Gaussians fall outside the image plane, providing little or no useful photometric supervision.

In contrast, with our curriculum design, E-RayZer learns pose estimation and Gaussian prediction jointly, allowing both components to automatically align to the same scale. This avoids the scale-misalignment issue and leads to more

Table 9. **Additional Results on Data Mixing and Scaling.** We train E-RayZer with different combinations of datasets. Compared to Tab. 5, we additionally include SpatialVID [56], a large in-the-wild video dataset. Results are color-ranked from red to yellow. Mixing datasets improves distribution coverage, whereas simply using larger datasets does not necessarily yield better performance – both diversity and data quality play critical roles.

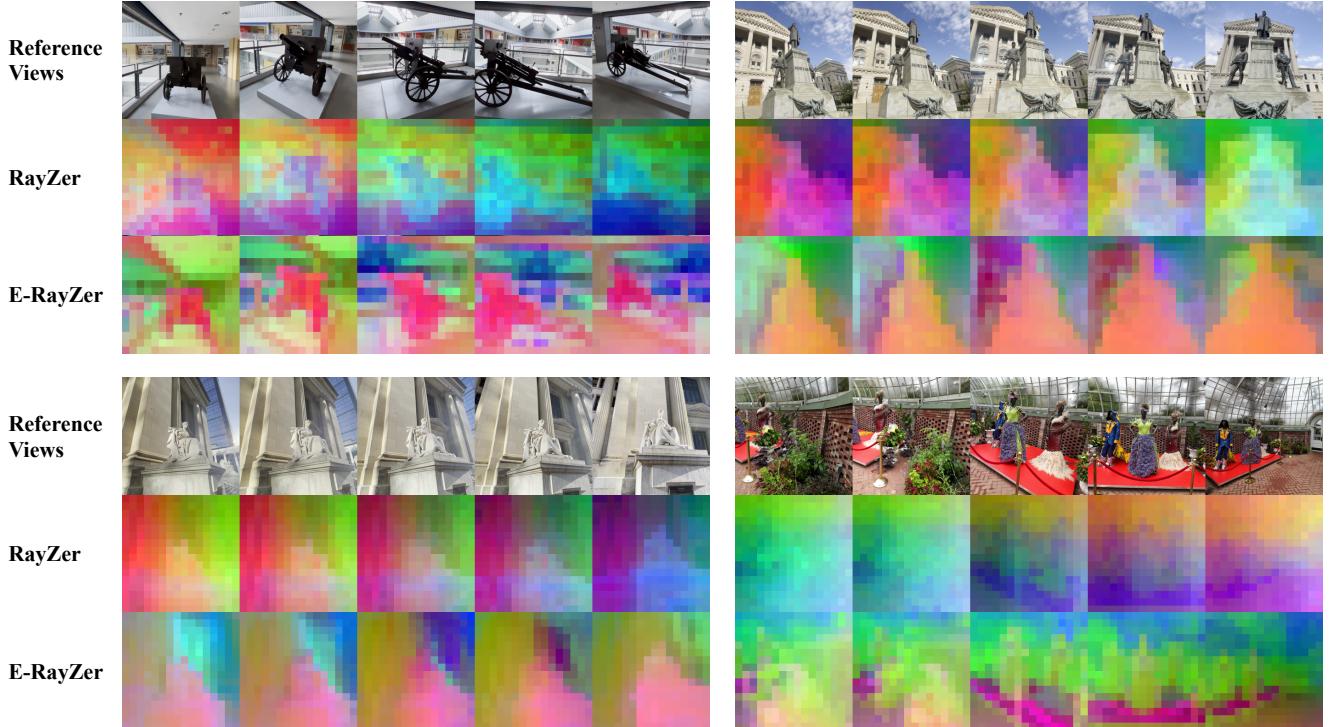| Training Data | # Seq. | NAVI [20] | | | | CO3Dv2 [40] | | | | ScanNet++ [71] | | | | DL3DV [33] | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | PSNR↑ | @5°↑ | @15°↑ | @30°↑ | PSNR↑ | @5°↑ | @15°↑ | @30°↑ | PSNR↑ | @5°↑ | @15°↑ | @30°↑ | PSNR↑ | @5°↑ | @15°↑ | @30°↑ |
| RE10K [81] | 66K | 17.2 | 1.8 | 16.9 | 34.0 | 19.1 | 0.6 | 8.3 | 26.0 | 17.5 | 1.1 | 13.3 | 37.3 | 17.3 | 21.2 | 55.0 | 72.7 |
| SpatialVID [56] | 100K | 17.9 | 0.7 | 11.2 | 26.4 | 19.9 | 0.2 | 5.7 | 20.9 | 18.0 | 0.3 | 6.7 | 26.0 | 17.2 | 11.4 | 36.6 | 56.0 |
| DL3DV [33] | 10K | 20.5 | 20.7 | 57.8 | 69.6 | 22.9 | 19.1 | 61.8 | 78.8 | 20.1 | 7.7 | 33.6 | 63.0 | 20.3 | 72.0 | 88.4 | 93.5 |
| 7-dataset Mix | 352K | 20.6 | 24.6 | 56.1 | 69.2 | 24.3 | 30.3 | 74.2 | 83.7 | 20.7 | 5.7 | 34.8 | 63.7 | 19.7 | 59.9 | 82.9 | 90.2 |



Figure 6. **Additional Visual Comparison with RayZer [23] on Learned Features.** We visualize feature maps using their top-three PCA components. The features produced by E-RayZer exhibit stronger and more spatially consistent patterns that align well with the underlying scene structure, whereas RayZer's features show noticeable color shifts across frames.

stable training and stronger novel-view synthesis performance. In short, this experiment further confirms the benefit of our self-supervised 3D reconstruction framework for both camera pose estimation and novel-view synthesis.

# E. Additional Results on Pre-training

We present additional results where E-RayZer is used as a pre-trained backbone for VGGT* (our re-implementation of VGGT [55], matched to our architecture and training data). We compare E-RayZer against RayZer [23] as an alternative pre-training approach and evaluate pose accuracy across multiple datasets.

Tab. 8 summarizes results under two training configurations: using only DL3DV [33] and using a mixture of seven datasets. Note that pre-training and supervised fine-tuning are conducted on the same data (*i.e.*, DL3DV or the

7-dataset mixture). In both settings, VGGT* initialized with E-RayZer outperforms its RayZer-initialized counterpart on most metrics, indicating that the representations learned by E-RayZer provide stronger and more transferable pre-training for downstream supervised pose estimation.

# F. Further Analysis of Training Data

We further analyze how different training datasets affect model performance.

Compared to Tab. 5, Tab. 9 additionally includes E-RayZer results on a static subset of SpatialVID [56], a large in-the-wild video dataset, and reports the number of training sequences used in each setting. We observe that a larger number of training sequences does not necessarily yield higher performance. For example, the model trained on 100K SpatialVID sequences performs comparably to the

RealEstate10K [81] model (which uses 66K sequences), yet significantly underperforms the DL3DV [33] model (which contains only 10K sequences). We conjecture that this gap stems from the noisy nature of in-the-wild data: SpatialVID sequences originate primarily from internet videos, and our training subsets are selected using their coarse dynamic-ratio labels. Also, SpatialVID often features simple or near-static camera motions. In contrast, DL3DV is carefully curated without moving objects and contains high-quality video sequences with diverse camera trajectories. These results support our earlier observations about data quality and highlight the importance of data curation when scaling self-supervised learning to large in-the-wild resources.

We also find that mixing datasets improves distribution coverage and leads to better generalization. For instance, models trained with mixed data perform better on the object-centric CO3Dv2 [40] compared to models trained solely on non-object-centric datasets.

Finally, we note that all experiments are conducted under a fixed computation budget (*i.e.*, 152K iterations with a global batch size of 192). Within this controlled setting, our results consistently suggest that diversity and quality of data matter more than quantity for training self-supervised models. We believe that collecting diverse, high-quality data remains both a key challenge and a promising direction for future work.

## G. More Qualitative Comparisons

**Learned Feature Representations.** In Fig. 6, we provide additional qualitative results comparing the learned feature representations of E-RayZer with those of RayZer [23]. Consistent with our observations in Fig. 5, the feature maps produced by E-RayZer exhibit more stable and coherent patterns across views, while RayZer's feature maps often display noticeable color shifts between frames. These results suggest that E-RayZer learns feature representations that are more geometrically grounded.

**Pose Estimation and Novel-view Synthesis.** We present additional qualitative comparison with baselines in Fig. 7. Compared to SPFSplat [19], E-RayZer consistently achieves better pose accuracy and higher-quality novel-view synthesis, despite being trained entirely from scratch without relying on pretrained priors such as MASt3R [31]. RayZer [23] generally produces high-quality novel views; however, it often exhibits grid-like artifacts in uncertain regions (highlighted with red bounding boxes). Moreover, RayZer's predicted poses are not physically aligned with the scene, whereas the camera poses learned by E-RayZer are geometrically grounded.

## References

[1] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. Flamingo: a visual language model for few-shot learning. In *NeurIPS*, 2022. 3

[2] Mido Assran, Adrien Bardes, David Fan, Quentin Garrido, Russell Howes, Matthew Muckley, Ammar Rizvi, Claire Roberts, Koustuv Sinha, Artem Zholus, et al. V-jepa 2: Self-supervised video models enable understanding, prediction and planning. *arXiv preprint arXiv:2506.09985*, 2025. 2

[3] Vassileios Balntas, Shuda Li, and Victor Prisacariu. Reloc-net: Continuous metric learning relocalisation using neural nets. In *ECCV*, 2018. 2

[4] Mohamed El Banani, Jason J Corso, and David F Fouhey. Novel object viewpoint estimation through reconstruction alignment. In *CVPR*, 2020. 2

[5] Adrien Bardes, Quentin Garrido, Jean Ponce, Xinlei Chen, Michael Rabbat, Yann LeCun, Mahmoud Assran, and Nicolas Ballas. Revisiting feature prediction for learning visual representations from video. *arXiv preprint arXiv:2404.08471*, 2024. 3

[6] Gilad Baruch, Zhuoyuan Chen, Afshin Dehghan, Tal Dimry, Yuri Feigin, Peter Fu, Thomas Gebauer, Brandon Joffe, Daniel Kurz, Arik Schwartz, et al. Arkitscenes: A diverse real-world dataset for 3d indoor scene understanding using mobile rgb-d data. In *NeurIPS D&B*, 2021. 5, 1

[7] Daniel Bolya, Po-Yao Huang, Peize Sun, Jang Hyun Cho, Andrea Madotto, Chen Wei, Tengyu Ma, Jiale Zhi, Jathushan Rajasegaran, Hanoona Rasheed, et al. Perception encoder: The best visual embeddings are not at the output of the network. *arXiv preprint arXiv:2504.13181*, 2025. 1, 2, 7

[8] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. In *NeurIPS*, 2020. 2

[9] Ruojin Cai, Bharath Hariharan, Noah Snavely, and Hadar Averbuch-Elor. Extreme rotation estimation using dense correlation volumes. In *CVPR*, 2021. 2

[10] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *ICCV*, 2021. 3

[11] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *NAACL*, 2019. 2

[12] Bardienus Duisterhof, Lojze Zust, Philippe Weinzaepfel, Vincent Leroy, Yohann Cabon, and Jerome Revaud. Mast3r-sfm: a fully-integrated solution for unconstrained structure-from-motion. In *3DV*, 2025. 2, 6

[13] Christoph Feichtenhofer, Yanghao Li, Kaiming He, et al. Masked autoencoders as spatiotemporal learners. In *NeurIPS*, 2022. 3

[14] Negar Foroutan, Paul Teiletche, Ayush Kumar Tarun, and Antoine Bosselut. Revisiting multilingual data mixtures in language model pretraining. *arXiv preprint arXiv:2510.25947*, 2025. 8
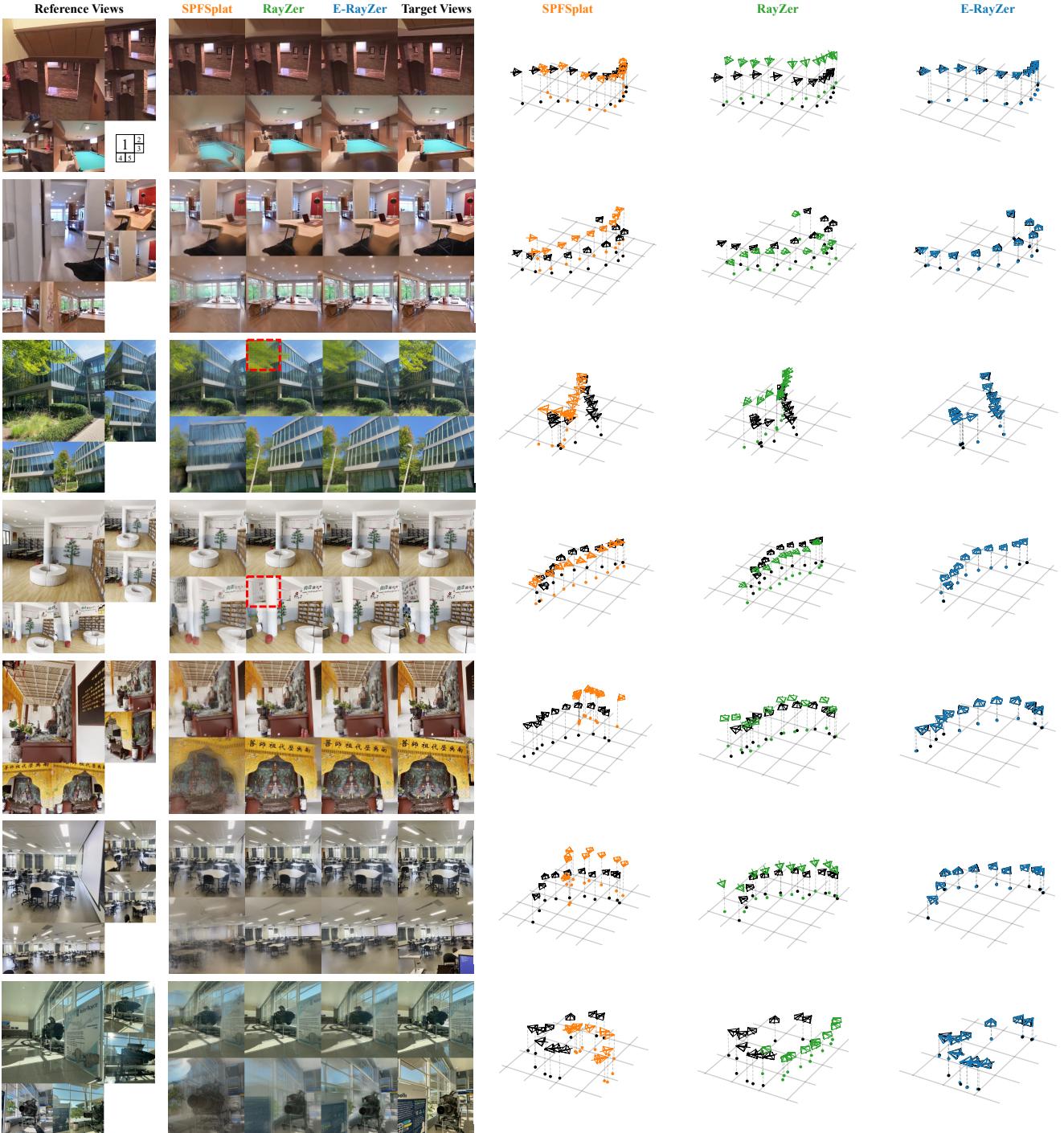
Figure 7. **Additional Visual Comparison with (Partially) Self-supervised Methods.** We show results for both novel-view synthesis (left) and pose estimation (right). The temporal order of the reference views is shown in the first row. Ground-truth poses are visualized in black, and predicted poses are aligned to the ground truth via an optimal similarity transform. E-RayZer outperforms baselines in pose accuracy, demonstrating its grounded 3D understanding. While RayZer [23] typically produces high-quality novel views, it often exhibits grid-like artifacts in low-texture regions (highlighted with red boxes; best viewed when zoomed in), likely due to its latent-rendering formulation.

[15] Yang Fu, Ishan Misra, and Xiaolong Wang. Mononerf: Learning generalizable nerfs from monocular videos without camera poses. In *ICML*, 2023. 3

[16] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *CVPR*, 2020. 3

[17] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. In *CVPR*, 2022. 2, 3

[18] Sunghwan Hong, Jaewoo Jung, Heeseong Shin, Jisang Han, Jiaolong Yang, Chong Luo, and Seungryong Kim. Pf3plat: Pose-free feed-forward 3d gaussian splatting for novel view synthesis. In *ICML*, 2025. 2, 6

[19] Ranran Huang and Krystian Mikolajczyk. No pose at all: Self-supervised pose-free 3d gaussian splatting from sparse views. In *ICCV*, 2025. 2, 6, 4

[20] Varun Jampani, Kevis-Kokitsi Maninis, Andreas Engelhardt, Arjun Karpur, Karen Truong, Kyle Sargent, Stefan Popov, André Araujo, Ricardo Martin Brualla, Kaushal Patel, et al. Navi: Category-agnostic image collections with high-quality 3d shape and pose annotations. In *NeurIPS*, 2023. 7, 8, 1, 2, 3

[21] Hanwen Jiang, Zhenyu Jiang, Kristen Grauman, and Yuke Zhu. Few-view object reconstruction with unknown categories and camera poses. In *3DV*, 2024. 2

[22] Hanwen Jiang, Zhenyu Jiang, Yue Zhao, and Qixing Huang. LEAP: Liberate sparse-view 3d modeling from camera poses. In *ICLR*, 2024. 2

[23] Hanwen Jiang, Hao Tan, Peng Wang, Haian Jin, Yue Zhao, Sai Bi, Kai Zhang, Fujun Luan, Kalyan Sunkavalli, Qixing Huang, et al. Rayzer: A self-supervised large view synthesis model. In *ICCV*, 2025. 1, 2, 3, 4, 6, 7, 8, 5

[24] Lihan Jiang, Yucheng Mao, Linning Xu, Tao Lu, Kerui Ren, Yichen Jin, Xudong Xu, Mulin Yu, Jiangmiao Pang, Feng Zhao, et al. Anysplat: Feed-forward 3d gaussian splatting from unconstrained views. In *ACM SIGGRAPH Asia*, 2025. 2

[25] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *ECCV*, 2016. 3

[26] Gyeongjin Kang, Jisang Yoo, Jihyeon Park, Seungtae Nam, Hyeonsoo Im, Sangheon Shin, Sangpil Kim, and Eunbyung Park. Selfsplat: Pose-free and 3d prior-free generalizable 3d gaussian splatting. In *CVPR*, 2025. 2

[27] Alex Kendall and Roberto Cipolla. Geometric loss functions for camera pose regression with deep learning. In *CVPR*, 2017. 7, 2

[28] Alex Kendall, Matthew Grimes, and Roberto Cipolla. Posenet: A convolutional network for real-time 6-dof camera relocalization. In *ICCV*, 2015. 1

[29] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. In *ACM ToG*, 2023. 2, 3, 4

[30] Zihang Lai, Sifei Liu, Alexei A Efros, and Xiaolong Wang. Video autoencoder: self-supervised disentanglement of static 3d structure and motion. In *ICCV*, 2021. 3

[31] Vincent Leroy, Yohann Cabon, and Jérôme Revaud. Grounding image matching in 3d with mast3r. In *ECCV*, 2024. 6, 4

[32] Amy Lin, Jason Y Zhang, Deva Ramanan, and Shubham Tulsiani. Relpose++: Recovering 6d poses from sparse-view observations. In *3DV*, 2024. 2

[33] Lu Ling, Yichen Sheng, Zhi Tu, Wentian Zhao, Cheng Xin, Kun Wan, Lantao Yu, Qianyu Guo, Zixun Yu, Yawen Lu, et al. Dl3dv-10k: A large-scale scene dataset for deep learning-based 3d vision. In *CVPR*, 2024. 5, 6, 7, 8, 1, 2, 3, 4

[34] Andrew Liu, Richard Tucker, Varun Jampani, Ameesh Makadia, Noah Snavely, and Angjoo Kanazawa. Infinite nature: Perpetual view generation of natural scenes from a single image. In *ICCV*, 2021. 5, 1

[35] Thomas W Mitchel, Hyunwoo Ryu, and Vincent Sitzmann. True self-supervised novel view synthesis is transferable. *arXiv preprint arXiv:2510.13063*, 2025. 3

[36] 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. In *TMLR*, 2024. 2, 5, 7

[37] Julius Plucker. Xvii. on a new geometry of space. In *Philosophical Transactions of the Royal Society of London*, 1865. 3

[38] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *ICML*, 2021. 3

[39] René Ranftl, Alexey Bochkovskiy, and Vladlen Koltun. Vision transformers for dense prediction. In *ICCV*, 2021. 7

[40] Jeremy Reizenstein, Roman Shapovalov, Philipp Henzler, Luca Sbordone, Patrick Labatut, and David Novotny. Common objects in 3d: Large-scale learning and evaluation of real-life 3d category reconstruction. In *ICCV*, 2021. 5, 7, 8, 1, 2, 3, 4

[41] Chris Rockwell, Justin Johnson, and David F Fouhey. The 8-point algorithm as an inductive bias for relative pose prediction by vits. In *3DV*, 2022. 2

[42] Mehdi SM Sajjadi, Aravindh Mahendran, Thomas Kipf, Etienne Pot, Daniel Duckworth, Mario Lučić, and Klaus Greff. Rust: Latent neural scene representations from unposed imagery. In *CVPR*, 2023. 3

[43] Kyle Sargent, Zizhang Li, Tanmay Shah, Charles Herrmann, Hong-Xing Yu, Yunzhi Zhang, Eric Ryan Chan, Dmitry Lagun, Li Fei-Fei, Deqing Sun, et al. Zeronvs: Zero-shot 360-degree view synthesis from a single real image. In *CVPR*, 2024. 5

[44] Johannes L Schonberger and Jan-Michael Frahm. Structure-from-motion revisited. In *CVPR*, 2016. 2

[45] Philipp Schröppel, Jan Bechtold, Artemij Amiranashvili, and Thomas Brox. A benchmark and a baseline for robust multi-view depth estimation. In *3DV*, 2022. 6, 7, 1

[46] Jamie Shotton, Ben Glocker, Christopher Zach, Shahram Izadi, Antonio Criminisi, and Andrew Fitzgibbon. Scene coordinate regression forests for camera relocalization in rgb-d images. In *CVPR*, 2013. 7, 1, 2

[47] Oriane Siméoni, Huy V Vo, Maximilian Seitzer, Federico Baldassarre, Maxime Oquab, Cijo Jose, Vasil Khalidov, Marc Szafraniec, Seungeun Yi, Michaël Ramamonjisoa, et al. Dinov3. *arXiv preprint arXiv:2508.10104*, 2025. 1, 2, 7

[48] Samarth Sinha, Jason Y Zhang, Andrea Tagliasacchi, Igor Gilitschenski, and David B Lindell. Sparsepose: Sparse-view camera pose regression and refinement. In *CVPR*, 2023. 2

[49] Brandon Smart, Chuanxia Zheng, Iro Laina, and Victor Adrian Prisacariu. Splatt3r: Zero-shot gaussian splatting from uncalibrated image pairs, 2024. arXiv preprint arXiv:2408.13912. 2

[50] Zhan Tong, Yibing Song, Jue Wang, and Limin Wang. Videomae: Masked autoencoders are data-efficient learners for self-supervised video pre-training. In *NeurIPS*, 2022. 2, 3

[51] Michael Tschannen, Manoj Kumar, Andreas Steiner, Xiaohua Zhai, Neil Houlsby, and Lucas Beyer. Image captioners are scalable vision learners too. In *NeurIPS*, 2023. 3

[52] Khiem Vuong, Anurag Ghosh, Deva Ramanan, Srinivasa Narasimhan, and Shubham Tulsiani. Aerialmegadepth: Learning aerial-ground reconstruction and view synthesis. In *CVPR*, 2025. 2

[53] Haoru Wang, Kai Ye, Yangyan Li, Wenzheng Chen, and Baoquan Chen. The less you depend, the more you learn: Synthesizing novel views from sparse, unposed images without any 3d knowledge. *arXiv preprint arXiv:2506.09885*, 2025. 3

[54] Jianyuan Wang, Christian Rupprecht, and David Novotny. Posediffusion: Solving pose estimation via diffusion-aided bundle adjustment. In *ICCV*, 2023. 2

[55] Jianyuan Wang, Minghao Chen, Nikita Karaev, Andrea Vedaldi, Christian Rupprecht, and David Novotny. Vggt: Visual geometry grounded transformer. In *CVPR*, 2025. 2, 4, 6, 7, 8, 1, 3

[56] Jiahao Wang, Yufeng Yuan, Rujie Zheng, Youtian Lin, Jian Gao, Lin-Zhuo Chen, Yajie Bao, Yi Zhang, Chang Zeng, Yanxi Zhou, et al. Spatialvid: A large-scale video dataset with spatial annotations, 2025. arXiv preprint arXiv:2509.09676. 3

[57] Limin Wang, Bingkun Huang, Zhiyu Zhao, Zhan Tong, Yinan He, Yi Wang, Yali Wang, and Yu Qiao. Videomae v2: Scaling video masked autoencoders with dual masking. In *CVPR*, 2023. 1, 2, 7, 8

[58] Qianqian Wang, Yifei Zhang, Aleksander Holynski, Alexei A Efros, and Angjoo Kanazawa. Continuous 3d perception model with persistent state. In *CVPR*, 2025. 2

[59] Ruoyu Wang, Yi Ma, and Shenghua Gao. Recollection from pensieve: Novel view synthesis via learning from uncalibrated videos. *arXiv preprint arXiv:2505.13440*, 2025. 3

[60] Shuzhe Wang, Vincent Leroy, Yohann Cabon, Boris Chidlovskii, and Jerome Revaud. Dust3r: Geometric 3d vision made easy. In *CVPR*, 2024. 2, 1

[61] Philippe Weinzaepfel, Vincent Leroy, Thomas Lucas, Romain Brégier, Yohann Cabon, Vaibhav Arora, Leonid Antsfeld, Boris Chidlovskii, Gabriela Csurka, and Jérôme Revaud. Croco: Self-supervised pre-training for 3d vision tasks by cross-view completion. In *NeurIPS*, 2022. 3

[62] Philippe Weinzaepfel, Thomas Lucas, Vincent Leroy, Yohann Cabon, Vaibhav Arora, Romain Brégier, Gabriela Csurka, Leonid Antsfeld, Boris Chidlovskii, and Jérôme Revaud. Croco v2: Improved cross-view completion pre-training for stereo matching and optical flow. In *ICCV*, 2023. 1, 2, 3, 7, 8

[63] Olivia Wiles, Georgia Gkioxari, Richard Szeliski, and Justin Johnson. Synsin: End-to-end view synthesis from a single image. In *CVPR*, 2020. 3

[64] Hongchi Xia, Yang Fu, Sifei Liu, and Xiaolong Wang. Rgbd objects in the wild: Scaling real-world 3d object learning from rgb-d videos. In *CVPR*, 2024. 5, 6, 7, 1, 2

[65] Sang Michael Xie, Hieu Pham, Xuanyi Dong, Nan Du, Hanxiao Liu, Yifeng Lu, Percy S Liang, Quoc V Le, Tengyu Ma, and Adams Wei Yu. Doremi: Optimizing data mixtures speeds up language model pretraining. In *NeurIPS*, 2023. 8

[66] Lihe Yang, Bingyi Kang, Zilong Huang, Xiaogang Xu, Jiashi Feng, and Hengshuang Zhao. Depth anything: Unleashing the power of large-scale unlabeled data. In *CVPR*, 2024. 5

[67] Yao Yao, Zixin Luo, Shiwei Li, Jingyang Zhang, Yufan Ren, Lei Zhou, Tian Fang, and Long Quan. Blendedmvs: A large-scale dataset for generalized multi-view stereo networks. In *CVPR*, 2020. 6, 7, 1, 2

[68] Botao Ye, Sifei Liu, Haofei Xu, Xueting Li, Marc Pollefeys, Ming-Hsuan Yang, and Songyou Peng. No pose, no problem: Surprisingly simple 3d gaussian splats from sparse unposed images. In *ICLR*, 2025. 2

[69] Jiasheng Ye, Peiju Liu, Tianxiang Sun, Jun Zhan, Yunhua Zhou, and Xipeng Qiu. Data mixing laws: Optimizing data mixtures by predicting language modeling performance. *arXiv preprint arXiv:2403.16952*, 2024. 8

[70] Vickie Ye, Ruilong Li, Justin Kerr, Matias Turkulainen, Brent Yi, Zhuoyang Pan, Otto Seiskari, Jianbo Ye, Jeffrey Hu, Matthew Tancik, et al. gsplat: An open-source library for gaussian splatting. In *Journal of Machine Learning Research*, 2025. 4

[71] Chandan Yeshwanth, Yueh-Cheng Liu, Matthias Nießner, and Angela Dai. Scannet++: A high-fidelity dataset of 3d indoor scenes. In *ICCV*, 2023. 6, 7, 8, 1, 2, 3

[72] Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. pixelnerf: Neural radiance fields from one or few images. In *CVPR*, 2021. 2

[73] Xianggang Yu, Mutian Xu, Yidan Zhang, Haolin Liu, Chongjie Ye, Yushuang Wu, Zizheng Yan, Chenming Zhu, Zhangyang Xiong, Tianyou Liang, et al. Mvimgnet: A large-scale dataset of multi-view images. In *CVPR*, 2023. 5, 1

[74] Jason Y Zhang, Deva Ramanan, and Shubham Tulsiani. Relpose: Predicting probabilistic relative rotation for single objects in the wild. In *ECCV*, 2022. 2

[75] Jason Y. Zhang, Amy Lin, Moneish Kumar, Tzu-Hsuan Yang, Deva Ramanan, and Shubham Tulsiani. Cameras as rays: Sparse-view pose estimation via ray diffusion. In *ICLR*, 2024. 2, 3

[76] Kai Zhang, Sai Bi, Hao Tan, Yuanbo Xiangli, Nanxuan Zhao, Kalyan Sunkavalli, and Zexiang Xu. Gs-lrm: Large reconstruction model for 3d gaussian splatting. In *ECCV*, 2024. 2

[77] Yuchen Zhang, Nikhil Keetha, Chenwei Lyu, Bhuvan Jhamb, Yutian Chen, Yuheng Qiu, Jay Karhade, Shreyas Jha, Yaoyu

Hu, Deva Ramanan, et al. Ufm: A simple path towards unified dense correspondence with flow. In *NeurIPS*, 2025. 5, 7

[78] Qitao Zhao and Shubham Tulsiani. Sparse-view pose estimation and reconstruction via analysis by generative synthesis. In *NeurIPS*, 2024. 2

[79] Qitao Zhao, Amy Lin, Jeff Tan, Jason Y Zhang, Deva Ramanan, and Shubham Tulsiani. Diffusionsfm: Predicting structure and motion via ray origin and endpoint diffusion. In *CVPR*, 2025. 2

[80] Tinghui Zhou, Matthew Brown, Noah Snavely, and David G Lowe. Unsupervised learning of depth and ego-motion from video. In *CVPR*, 2017. 3

[81] Tinghui Zhou, Richard Tucker, John Flynn, Graham Fyffe, and Noah Snavely. Stereo magnification: Learning view synthesis using multiplane images. In *ACM SIGGRAPH*, 2018. 5, 6, 7, 8, 1, 2, 3, 4

[82] Zhizhuo Zhou and Shubham Tulsiani. Sparsefusion: Distilling view-conditioned diffusion for 3d reconstruction. In *CVPR*, 2023. 2