

LEVERAGING NEURAL RADIANCE FIELDS FOR POSE ESTIMATION OF AN UNKNOWN SPACE OBJECT DURING PROXIMITY OPERATIONS

Antoine Legrand ^{1,2,4*}, Renaud Detry ^{2,3}, and Christophe De Vleeschouwer ¹

¹Department of Electrical Engineering (ELEN), ICTEAM, UCLouvain

²Department of Electrical Engineering (ESAT), KU Leuven

³Department of Mechanical Engineering (MECH), KU Leuven

⁴AerospaceLab

ABSTRACT

We address the estimation of the 6D pose of an unknown target spacecraft relative to a monocular camera, a key step towards the autonomous rendezvous and proximity operations required by future Active Debris Removal missions. We present a novel method that enables an "off-the-shelf" spacecraft pose estimator, which is supposed to know the target CAD model, to be applied on an unknown target. Our method relies on an in-the-wild NeRF, *i.e.* a Neural Radiance Field that employs learnable appearance embeddings to represent varying illumination conditions found in natural scenes. We train the NeRF model using a sparse collection of images that depict the target, and in turn generate a large dataset that is diverse both in terms of viewpoint and illumination. This dataset is then used to train the pose estimation network. We validate our method on the Hardware-In-the-Loop images of SPEED+ [1] that emulate lighting conditions close to those encountered on orbit. We demonstrate that our method successfully enables the training of an off-the-shelf spacecraft pose estimation network from a sparse set of images. Furthermore, we show that a network trained using our method performs similarly to a model trained on synthetic images generated using the CAD model of the target.

1 INTRODUCTION

With an ever growing number of satellites in orbit, the risk of collision between a satellite and space debris, *e.g.*, rocket bodies, defunct satellites or pieces from a previous collision, is steadily rising. Such a collision would not only cause the destruction of a functional satellite but also dramatically increase the number of space debris [2], thereby further increasing the risk of such a collision [3]. As a result, private companies and space agencies are working on Active Debris Removal (ADR) [4, 5, 6] missions that aims at de-orbiting space debris. These ADR missions require to perform Rendezvous and Proximity Operations (RPO) with a non-cooperative target, *i.e.*, a chaser spacecraft must operate close to, or even dock with, a target spacecraft which was not designed to support RPOs. Due to the risk of human failures implied by tele-operated operations, those RPOs should be carried autonomously by the chaser spacecraft.

A key capability to perform autonomous RPOs is the on-board estimation of the relative pose, *i.e.*, position and orientation of the target spacecraft relative to the chaser. Due to their low cost, low mass and compactness, monocular cameras are considered for this task [7, 8]. Although the vision-based estimation of the relative pose of a non-cooperative spacecraft has been studied in depth in the literature [9, 10, 11], current solutions assume the knowledge of the CAD model of the target spacecraft that enables the generation of large synthetic training sets. In the case of Active Debris Removal, this assumption does not hold since little information is known about the debris. This work aims at leveraging Neural Radiance Fields (NeRFs) models [12] to extend the scope of existing pose estimation methods to unknown targets, *i.e.*, targets for which the CAD model is not available.

For this purpose, we consider a three-steps approach, as illustrated in Figure 1. Firstly, the chaser spacecraft is tele-operated to approach the target up to a safe distance. During the approach,

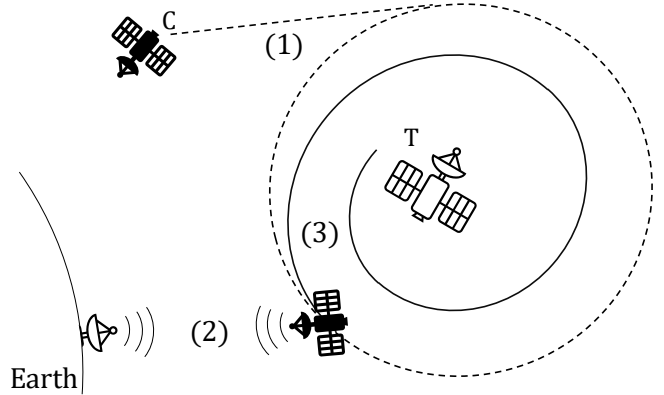


Figure 1: Overview of the considered three-steps operation. **(1)** The chaser spacecraft (C) approaches the target (T) while taking pictures. **(2)** The images are downloaded on ground and processed to train a Spacecraft Pose Estimation network, whose weights are uploaded on the chaser. **(3)** The chaser finishes the operation autonomously, by relying on the trained network.

the chaser acquires images of the target and transmit them to a ground station. Then, those few images are processed on-ground to synthesize additional views of the target under varying illuminations, so as to build a sufficiently rich set of images to train an "off-the-shelf" pose estimation network, *i.e.* an existing neural network that only requires to be trained on a new set depicting the target. Finally, the model weights are uploaded on the spacecraft which autonomously performs the final approach. The on-ground processing step enables the use of the virtually infinite computing resources available on ground, in contrast with the low-power on-board hardware. Furthermore, even if the chaser spacecraft requires ground support in this scenario, it operates autonomously during the critical, *i.e.*, close-range, phase of the operation.

SoA SPE methods	Model-agnostic	Validation on realistic images
known target [9, 10, 11]		✓
unknown target [15, 16, 17]	✓	
Ours	✓	✓

Table 1: Overview of the State-of-The-Art in the field of Spacecraft Pose Estimation. Our method is the first model-agnostic one to be validated on realistic images.

To demonstrate the feasibility of such a 3-steps procedure in terms of image analysis requirements, this paper aims at studying the performance obtained when training a pose estimator from a few images. Hence, it focuses on the on-ground processing step, *i.e.* on the training of a spacecraft pose estimation model from a small number of spaceborne images depicting that target spacecraft. For this purpose, our method resorts to an implicit representation of the target, under the form of an "in-the-wild" Neural Radiance Field [12, 13, 14]. This NeRF is then used to generate a sufficiently large training set which captures the diversity of both the pose distribution and the illumination conditions encountered in orbit. Finally, an off-the-shelf Spacecraft Pose Estimation (SPE) network is trained on this set. As pointed out in Table 1, our work is the first to validate a vision-based spacecraft pose estimation method for unknown targets on realistic images, *i.e.*, the Hardware-in-the-loop images of SPEED+ [1].

2 SPACECRAFT POSE ESTIMATION

The existing methods for spacecraft pose estimation are either model-agnostic or model-based. While the first ones do not assume any prior on the target spacecraft, the latter ones exploit the CAD model of the target.

Regarding the model-agnostic methods, most works rely on (i) sensors such as stereo [18, 19] or Time-of-Flight [20] cameras or (ii) on multiple chaser spacecraft [21]. A few works [15, 16] tackled the problem using a monocular camera while relying on keypoint detection and description [15, 22, 23, 24] or edge detection [16] to estimate the pose from the observed image. Apart from those keypoint/edge detection methods, Park *et al.* [17] proposed a CNN-based method to recover both the shape and the pose of an unknown spacecraft from a single image. However, those works were only validated on synthetic images. Unlike them, our method, which extends the applicability of model-based methods to an unknown target observed from a sparse set of viewpoints, is successfully demonstrated on Hardware-In-the-Loop (HIL) images from the SPEED+ dataset [1].

Model-based methods received a strong interest in recent years, notably through the Spacecraft Pose Estimation Challenges (SPECs) [25, 26]. State-of-the-Art solutions rely on Convolutional Neural Networks (CNN) to either directly estimate the spacecraft pose from the image [27, 11, 9] or predict keypoint coordinates [10, 28, 29] which are then used to recover the pose by solving the Perspective-n-Points (PnP) problem [30]. While most works [10, 28, 29] rely on an intermediate step of spacecraft detection to identify a Region-of-Interest to crop and process through the CNN, some works [11, 9] directly perform the estimation on the whole image. In all cases, those methods

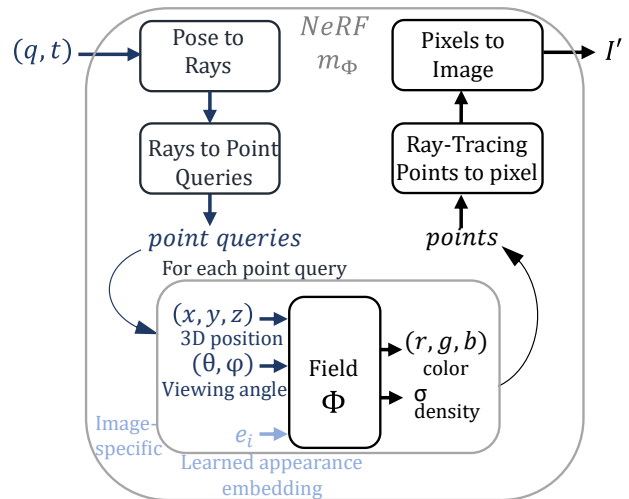


Figure 2: Overview of the generation of an image $I' = m_\Phi(q, t)$ through an in-the-wild Neural Radiance Field [14, 13] m_Φ for a given pose (q, t) . As explained in section 3, given a pose, *i.e.* rotation q and translation t , the NeRF renders an image I' by querying multiple times a MLP, which approximates the radiance field, and by aggregating the predicted color and density of the points through ray-tracing techniques. Unlike most NeRFs where the MLP takes as input a 3D position and 2 viewing angles to output the color and density of each point, In-the-wild NeRFs [14] also feeds the MLP with a learnable appearance embedding which is specific to each image. This offers the possibility to change the illumination conditions when synthesizing new images (see section 4.4).

always rely on a large synthetic training set generated using the CAD model of the target spacecraft and therefore implicitly assume the knowledge of this CAD model, making them unsuitable in front of unknown targets.

Our work leverages a Neural radiance Field [12] (NeRF) trained on a few spaceborne images to generate the large training set required by model-based Spacecraft Pose Estimation (SPE) pipelines. Thereby, it makes an arbitrary off-the-shelf, model-based, SPE network relevant to infer the pose of an unknown target, only characterized by a few image samples.

3 NEURAL RADIANCE FIELDS

Neural Radiance Fields [12] (NeRF) were developed as a tool to render novel views of a scene. In a NeRF, the scene is implicitly represented by a learned neural network. To render a novel view, the image is processed pixel per pixel. For each pixel, a ray is projected in the scene and several points are sampled along that ray. Each point is determined by a 3D position and 2 viewing angles which are fed in a MLP that outputs an rgb triplet and a volume density σ . The value of the pixel is recovered through differentiable ray-tracing techniques that aggregate all the points along the ray. To train this implicit representation, a reconstruction loss, computed between the rendered image and the corresponding ground-truth, is back-propagated through the neural network.

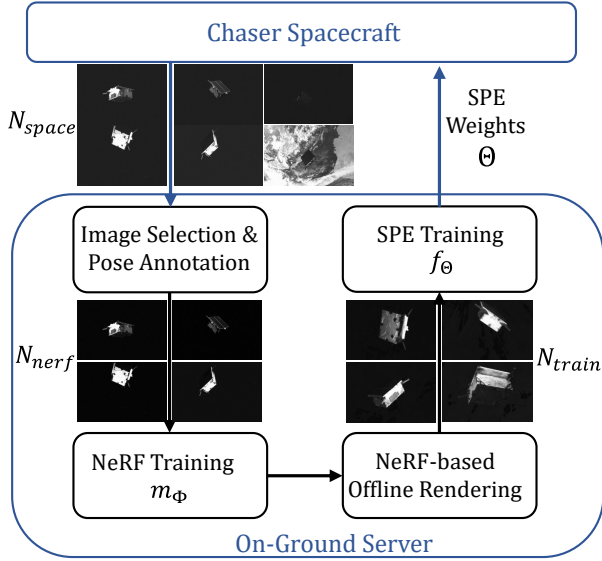


Figure 3: On-ground processing pipeline. A subset of the images downloaded from the spacecraft is annotated and used to train a NeRF [12], m_ϕ . This radiance field is then used to render a large training set, which is exploited to train an off-the-shelf SPE network f_Θ . Finally, the network weights, Θ , are uploaded on the spacecraft. Variables are defined in the text.

Improvements to the original NeRF method [12] have been proposed regarding the training time [31] or the ability to represent real scenes [14]. NeRF in the Wild [14] has introduced learnable appearance embeddings that are associated to each training image. As depicted in Figure 2, the embedding is given to the MLP along with the input coordinates, thereby allowing the NeRF to render the same view under different appearances through the appearance embeddings. In this work, as both the fast training and the ability to represent a real scene are necessary, *K*-Planes [13] is used as it combines an efficient encoding with learnable appearance embeddings.

Neural Radiance Fields are used in many fields, *e.g.* robotics [32, 33], virtual reality [34], and used for diverse tasks such as novel views rendering [12], visual scene understanding [35] or pose estimation [36]. In the aerospace field, NeRFs have been used to render novel views of Mars [37], recover the 3D shape of space objects [38, 39] or estimate the state of a spacecraft [40]. However, they have never been used as a tool enabling the training of a Spacecraft Pose Estimation network on an unknown target.

4 METHOD

As stated in the introduction, this paper provides a method that enables the use of an off-the-shelf, model-based, Spacecraft Pose Estimation (SPE) network on an unknown target in the context of a semi-autonomous Rendezvous and Proximity Operation (RPO).

4.1 Overview

As illustrated in Figure 1, the considered RPO is made of 3 steps. First, the chaser spacecraft is tele-operated to approach the target and take pictures that are transferred to a ground station. On ground, the images are processed to train a SPE network whose weights are then uploaded on the chaser spacecraft. Finally, the chaser performs the final approach autonomously by exploiting the so trained pose estimation network.

This section describes the on-ground processing required to train an off-the-shelf spacecraft pose estimation model from a sparse set of spaceborne images. As depicted in Figure 3, N_{space} images are downloaded from the chaser spacecraft. From this set, N_{nerf} high-quality images, *i.e.* with good illumination conditions, are selected and their pose is annotated. They are then used to train a Neural Radiance Field (NeRF) [12] m_ϕ that learns an implicit representation of the target spacecraft. This radiance field is then used to generate a training set made of N_{train} images which is used to train an off-the-shelf SPE network f_Θ whose weights Θ are finally uploaded on the chaser spacecraft. Those steps are detailed in the following sections.

4.2 Images Selection and Pose Annotation

Due to the harsh lighting conditions encountered in orbit, some of the downloaded images can be over-exposed or under-exposed. As those images contain little information and would act as a noisy and misleading supervision in the NeRF training, they are discarded. Similarly, all the images where the Earth appears in the background are removed. Indeed, in a field aligned with the target, the Earth is a transient object which can not be explained by the NeRF. Since exploiting those images to train the NeRF would introduce significant artifacts, they are simply discarded. Finally, each image is annotated with pose information.

4.3 NeRF Training

Using 90% of the N_{nerf} images, we train an "in-the-wild" NeRF [14, 13] m_ϕ , *i.e.*, a Neural Radiance Field that contains learnable appearance embeddings as illustrated in Figure 2. Those embeddings enable the network to capture illumination conditions that are specific to each image and therefore render images with a larger illumination diversity.

4.4 Offline Image Rendering

Training a SPE network requires a large number of images in order to capture the diversity of the pose distribution as well as of the illumination conditions. To generate this large training set, the learned NeRF, m_ϕ , is used to render N_{train} images with pose labels randomly sampled in $SE(3)$, *i.e.* the set of rigid-body transformation in 3D space. As introduced in [14], for each image, an appearance embedding is generated by interpolating two random appearance embeddings from the NeRF training set, *i.e.*, let α be a random scalar between 0 and 1, and let e_i and e_j be two randomly picked appearance embeddings from the NeRF training images, the interpolated appearance embedding e is computed as:

$$e = e_i + \alpha(e_j - e_i) \quad (1)$$

Figure 4 depicts several images generated using this appearance interpolation strategy.

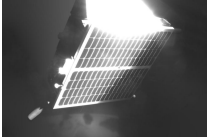


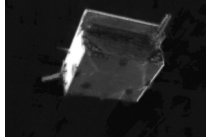
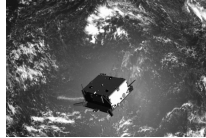

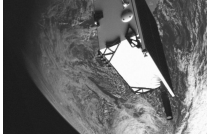


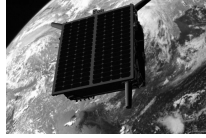

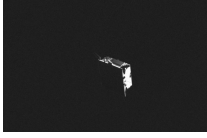


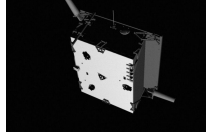


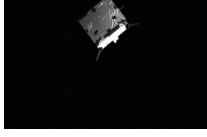
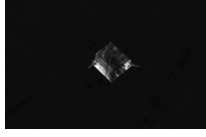
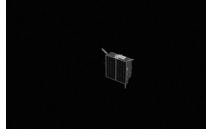
	<i>Sunlamp</i> [1]	<i>Lightbox</i> [1]		our training set	synthetic training set [1]
		<i>Lightbox*</i>	<i>Lightbox-500</i>		
Usage Domain	SPE test HIL	SPE test HIL	NeRF-training HIL	SPE train NeRF-based	SPE train Synthetic
Illumination	Direct	Diffuse	Diffuse	Randomized	Synthetic
# images	2791	6240	500	48,000	48,000
Examples					
					
					
					

Table 2: Overview of the sets used in this paper. *Sunlamp* and *Lightbox* are two Hardware-In-the-Loop (HIL) sets [1]. *Lightbox* is divided into *Lightbox** and *Lightbox-500*. *Lightbox-500* contains 500 images, selected as explained in section 4.2, while *Lightbox** contains all the remaining ones. *Sunlamp* and *Lightbox** are only used for testing. Because of the direct illumination of the Sun, which causes over-exposed images, *Sunlamp* is the most challenging test set. *Lightbox-500* contains the images used for training the NeRF. Our training set contains images generated through the trained NeRF and uses the same ground-truth labels as the synthetic training set of SPEED+ [1]. Our purpose is to demonstrate that the NeRF-based training set, generated using few spaceborne images, results in a SPE network that achieves an accuracy close to the one it achieves with the large synthetic set, even if the latter approach requires the knowledge of the CAD model to generate the synthetic set while our method is model-agnostic.

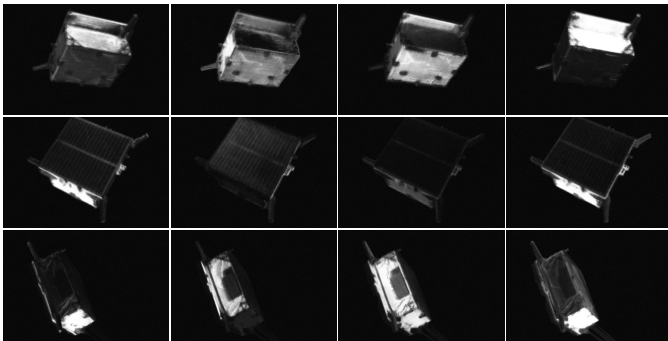


Figure 4: Images rendered by the same NeRF using 4 different appearance embeddings (columns 1-4). Randomizing the appearance embeddings enables the generation of images with diverse illumination conditions, regarding both the intensity and the orientation of the illumination source.

4.5 SPE Network Training

The off-the-shelf Spacecraft Pose Estimation (SPE) network, f_{Θ} , is trained on the generated N_{train} images. Its trained weights, Θ , are finally uploaded on the chaser spacecraft.

5 EXPERIMENTS

This section demonstrates how our method successfully enables the training of a Spacecraft Pose Estimation (SPE) network from a sparse set of real images depicting an unknown target spacecraft. Implementation details are provided in section 5.1. section 5.2 discusses the accuracy of the proposed approach while section 5.3 provides an ablation study evaluating the impact of the in-the-wild abilities of the NeRF on the SPE accuracy.

5.1 Implementation Details

Dataset

Our method is validated on the SPEED+ dataset [1], used in the Spacecraft Pose Estimation Challenge (SPEC) 2021 [26], co-organized by the European Space Agency and Stanford University. SPEED+ contains synthetic and Hardware-In-the-Loop (HIL) images of the TANGO spacecraft from the PRISMA mission [41] and the corresponding pose labels. The spacecraft CAD model was not released. SPEED+ contains 59,960 grayscale synthetic images of resolution 1920x1200. The interspacecraft distance is between 2.2 and 10 meters. In the following experiments, these synthetic images are divided into a training set (80%) and a validation set (20%).

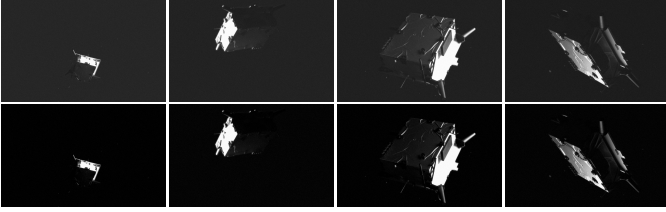


Figure 5: **(Top)** Original images. **(Bottom)** Images after pre-processing. For each image, the average intensity of the background is subtracted and the range is extended so that the maximal intensity of both images is unchanged.

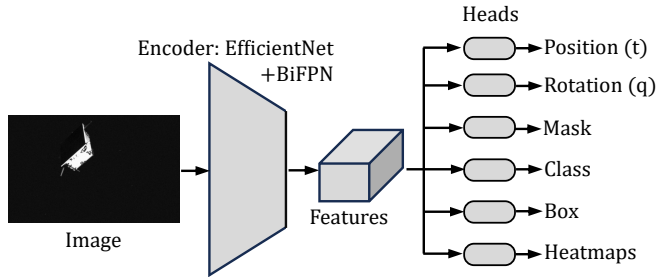


Figure 6: Overview of the off-the-shelf Spacecraft Pose Estimation network used in this paper, *i.e.* SPNv2 [9]. The network consists in a shared encoder followed by six heads inspired from EfficientPose [43]. The encoder is made of an EfficientNet backbone [44] and a BiFPN [45]. Two heads are used to estimate the relative pose, *i.e.* position and orientation while the four last heads are used to enhance the generalization abilities of the network through Multi-Task Learning. The auxiliary tasks consist in classification, bounding box prediction, foreground segmentation and heatmap-based keypoint regression. The $K = 11$ keypoints correspond to the 8 corners of the spacecraft body and the top of 3 antennas.

The Hardware-In-the-Loop (HIL) images were taken in the Testbed for Rendezvous and Optical Navigation (TRON) [42] facility, which emulates the lighting conditions encountered in Low Earth Orbit. HIL images are split in two domains, *lightbox* and *sunlamp* that aim at replicating different illumination conditions. The *lightbox* domain contains 6740 images where the Earth albedo is simulated using light boxes, while the *sunlamp* domain contains 2791 images where the direct illumination of the Sun is simulated using a metal halide lamp. In both domains, the inter-spacecraft distance is between 2.5 and 9.5 meters.

In this paper, similar to [9, 26], the *lightbox* set is assumed to accurately render the lighting conditions encountered in Low Earth Orbit, and will thus be used to provide training and testing images that have to be representative of real images captured in orbit. From this set, all the images where (i) the inter-spacecraft distance is below five meters, or (ii) the Earth appears in the background, or (iii) the spacecraft is under-exposed or over-exposed, are removed, as explained in section 4.2. 500 images are randomly picked from the remaining ones and the average intensity of their background is removed from each image, as illustrated in Figure 5, to ensure a smooth training. The resulting set is named *lightbox-500*. It is used only for training the NeRF while the rest of *lightbox*, dubbed as *lightbox**, is kept for testing. The Neural Radiance Field, m_ϕ , is trained on N_{nerf} images,

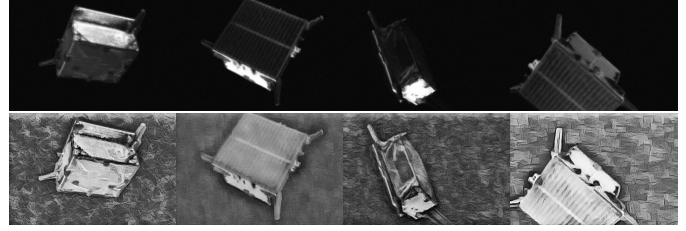


Figure 7: Examples of Neural Style Augmentation [46]. **(Top)** Original images **(Bottom)** Corresponding images after Style Augmentation. The spacecraft shape is always preserved by the transformation but its texture is augmented.

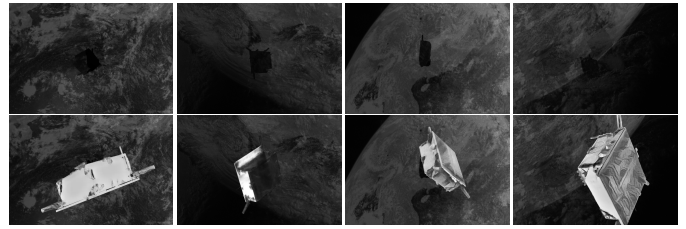


Figure 8: **(Top)** Examples of background images recovered by merging 4 randomly picked background images from the SPEED+ [1] synthetic set. **(Bottom)** Examples of images augmented with these backgrounds. The two first ones are performed from the NeRF output while the two last ones first went through the style augmentation stage.

resized to 960x600 pixels, randomly picked from *lightbox-500*. Although the pose labels associated to the N_{nerf} images could be annotated through point matching across the images, we use the pose labels provided by the SPEED+ dataset [1]. The NeRF is then used to generate a set of $N_{train} = 48,000$ images to train the pose estimation model f_Θ . Those N_{train} images share the same pose labels as the synthetic train set of SPEED+. Finally, the SPE network, f_Θ , is trained on the rendered set and the SPE metrics are computed on the *sunlamp* and *lightbox** sets. Table 2 presents an overview of the sets used in this paper.

In-the-wild NeRF m_ϕ

In this paper, the in-the-wild Neural Radiance Field, m_ϕ , follows the K -Planes [13] implementation because it features both learnable appearance embeddings and efficient encoding, thereby reducing the training time. m_ϕ uses the default configuration of K -Planes except that it does not use a linear decoder but a MLP [13]. The appearance embedding dimension is kept at 32. The network is trained for 30,000 steps of 4096 points per batch. The NeRF training takes 40 minutes on a NVIDIA RTX3090. Figure 4 highlights the ability of K -Planes to render not only the pose distribution, but also variable illumination conditions thanks to its appearance embedding. m_ϕ is then used to render 48,000 images that correspond to the pose labels from the synthetic training set of SPEED+ [1] used with randomly interpolated appearance embeddings, as explained in section 4.4. The rendering of the whole set takes 16 GPU-hours and is performed on a cluster made of 4 NVIDIA RTX3090 and 10 NVIDIA TeslaA100 GPUs.

Training Strategy	N	No CAD required	<i>Sunlamp</i>				<i>Lightbox*</i>			
			S_{Pose}^* [/]	E_R [°]	E_T [m]	\bar{E}_T [%]	S_{Pose}^* [/]	E_R [°]	E_T [m]	\bar{E}_T [%]
Synthetic (w. CAD)	/		0.322	15.7	0.30	0.05	0.203	9.5	0.25	0.04
Baseline	500	✓	2.005	99.4	1.20	0.27	1.96	95.0	1.31	0.30
NeRF (Ours)	50	✓	0.502	24.1	0.57	0.08	0.213	9.9	0.24	0.04
NeRF (Ours)	100	✓	0.491	23.6	0.54	0.08	0.161	7.4	0.19	0.03
NeRF (Ours)	200	✓	0.375	15.3	0.37	0.06	0.158	7.2	0.20	0.03
NeRF (Ours)	500	✓	0.341	16.9	0.29	0.04	0.158	7.2	0.20	0.03

Table 3: Comparison of three SPE training strategies, *i.e.*, (i) trained on 48,000 synthetic images generated using the CAD model of the target, (ii) trained on N real images (baseline), or (iii) trained on 48,000 images generated by a NeRF trained on N real images, and tested on two test sets, *i.e.*, *Sunlamp* and *Lightbox**. Our method significantly outperforms the baseline and reaches similar performance as the synthetic approach while requiring no CAD model of the target.

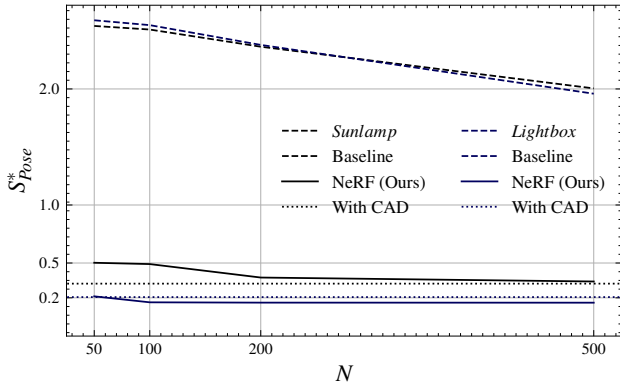


Figure 9: Relationship between the number of real images, N and the SPEED+ score of the SPE when trained on (i) N real images (dashed lines), (ii) 48,000 images generated by a NeRF trained on N images (solid lines) or (iii) 48,000 synthetic images rendered with the CAD model of the target (dotted lines), for the *Sunlamp* (blue) and *Lightbox** (black) sets. Our NeRF-based approach always outperforms the baseline by a large margin. When trained with enough images, *e.g.* 500, it performs similarly as the CAD-based approach while requiring no CAD model.

Off-the-shelf SPE network f_{Θ}

In our experiments, the Spacecraft Pose Estimation network, f_{Θ} , consists in a SPNv2 [9] with scaling coefficient $\phi = 3$ and Batch-Normalization layers [47]. Figure 6 illustrates the SPNv2 architecture. The training is conducted using the same procedure as in the original paper [9], including the offline neural Style Augmentation [46], as illustrated in Figure 7. In addition, we add a custom background augmentation, illustrated in Figure 8, which is applied on half of the training images to add the Earth in background. The SPE is trained for 30,000 steps on batches of 32 images on a NVIDIA TeslaA100 GPU.

Metrics

Our method is evaluated by the SPEED+ score, S_{Pose}^* , introduced in SPEC2021 [26], which averages the pose scores of the N images of the test set. For each image, the translation error, e_t , is computed as the norm of the error between the predicted position, \hat{t} , and the ground-truth position, t , *i.e.*

$$e_t = \|\hat{t} - t\|, \quad (2)$$

while the normalized translation error, \bar{e}_t , is equal to the ratio between the translation error and the norm of the ground-truth

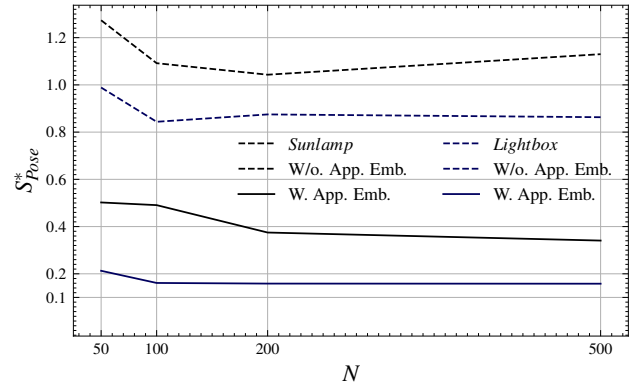


Figure 10: Relationship between the number of real images, N , and the SPEED+ score of the SPE network when trained on 48,000 images generated using either (i) the average appearance embedding (dashed lines), or (ii) interpolated appearance embeddings (solid lines), trained on N real images and tested on the *Sunlamp* (blue) and *Lightbox** (black) sets. The SPE network trained on images rendered with interpolated appearance embeddings always outperforms by a large margin a SPE network trained on images generated with the average appearance embedding.

position, *i.e.*

$$\bar{e}_t = \frac{e_t}{\|t\|}. \quad (3)$$

The rotation error, e_q , is computed as the angular error between the predicted quaternion, \hat{q} , and the ground-truth quaternion, q , *i.e.*

$$e_q = 2\arccos\left(\left|\hat{q}q^T\right|\right). \quad (4)$$

As the ground-truth positions and rotations of the HIL domains were measured with some uncertainty [1], the predictions are considered as perfect if they are lower than the calibration error, *i.e.*

$$e_q^* = \begin{cases} 0, & \text{if } e_q < 0.00295\text{rad} \\ e_q, & \text{otherwise,} \end{cases} \quad (5)$$

and

$$\bar{e}_t^* = \begin{cases} 0, & \text{if } \bar{e}_t < 2.173\text{mm/m} \\ \bar{e}_t, & \text{otherwise} \end{cases} \quad (6)$$

The average normalized translation, $E_{T,N}$, translation, E_T , and rotation errors, E_R , are computed as the average of the corre-

Offline Rendering	<i>Sunlamp</i>				<i>Lightbox*</i>			
	S_{Pose}^* [°]	E_R [°]	E_T [m]	\bar{E}_T [°]	S_{Pose}^* [°]	E_R [°]	E_T [m]	\bar{E}_T [°]
Average Appearance Embedding	1.130	60.2	0.45	0.08	0.863	39.6	0.95	0.18
Interpolated Appearance Embedding	0.341	16.9	0.29	0.04	0.158	7.2	0.20	0.03

Table 4: Comparison of the SPE metrics when trained on 48,000 images generated by a NeRF using either (i) the average appearance embedding, or (ii) using interpolated appearance embeddings. In both cases, the NeRF is trained on 500 real images. Generating the SPE training set with interpolated appearance embeddings decreases the SPE errors.

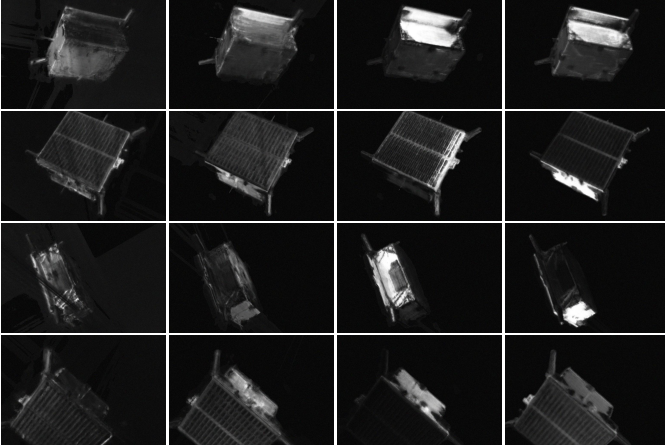


Figure 11: Images generated by NeRFs trained on 50, 100, 200 or 500 images (left to right). Even if 50 images are sufficient to recover a fair representation of the target spacecraft, increasing the number of training images reduces the artifacts and enables the recovery of finer details.

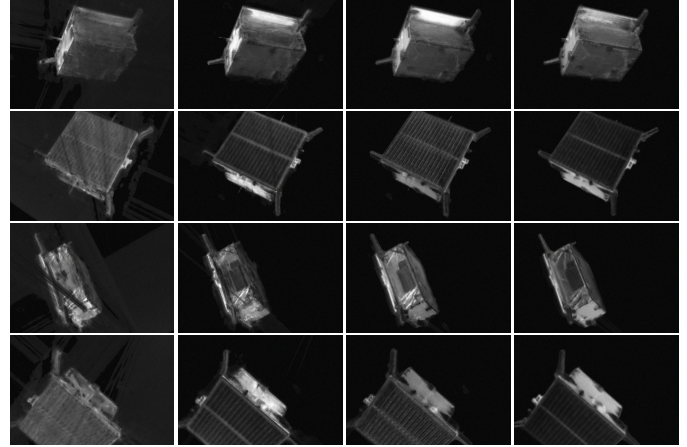


Figure 12: Images generated by NeRFs trained on 50, 100, 200 or 500 images (left to right) and rendered using average appearance embeddings. As for the appearance interpolation strategy (see fig. 11), 50 images are sufficient to recover a fair representation of the target spacecraft while finer details can be recovered by increasing the number of training images. Unlike the images generated using randomly interpolated appearance embeddings, the rendered images do not capture the illumination diversity depicted in Figure 11.

sponding values over the set, *i.e.*

$$\bar{E}_T = \frac{1}{N} \sum_{i=1}^N \bar{e}_t^{(i)}, \quad (7)$$

$$E_T = \frac{1}{N} \sum_{i=1}^N e_t^{(i)}, \quad E_R = \frac{1}{N} \sum_{i=1}^N e_q^{(i)}. \quad (8)$$

Finally, the SPEED+-score, S_{Pose}^* , computes the average of the sum of the normalized translation and rotation errors over the set, *i.e.*

$$S_{Pose}^* = \frac{1}{N} \sum_{i=1}^N (\bar{e}_t^{*(i)} + e_q^{*(i)}). \quad (9)$$

5.2 Evaluation

Figure 9 compares the SPEED+ score, *i.e.* the sum of the average rotation and normalized translation errors (see section 5.1), achieved by the SPE network trained either on N spaceborne images directly or on 48,000 images generated by a NeRF trained on those N images. While the baseline, straightforward, approach fails in estimating the pose of the target spacecraft, our approach successfully predicts the pose from a limited set of spaceborne images regardless of the number of images used to train the NeRF. Indeed, as highlighted in table 3, while the baseline approach exhibits errors of 99.4° - 1.20m on *sunlamp*

and 95.0° - 1.31m on *lightbox*, our approach achieves errors of 16.9° - 0.29m and 7.2° - 0.20m on those sets. Figure 9 also depicts the SPEED+ score obtained with the same SPE network but trained on the 48,000 synthetic images of SPEED+, rendered using the CAD model of the target. Our approach, which is model-agnostic, trained with 500 real images, performs as well as, or even better than, the CAD-dependent one, which achieves errors of 15.7° - 0.30m and 9.5° - 0.25m on those sets. Examples of images generated with our method are depicted in Figure 11.

5.3 Ablation Study: Appearance Interpolation Strategy

This section explores the impact of the appearance interpolation strategy on the accuracy of the SPE network. For this purpose, Figure 10 depicts the relationship between the number of images used to train the NeRF and the SPEED+ score of the SPE trained on images rendered using either randomly interpolated appearance embeddings or a single appearance embedding that corresponds to the average of the learned appearance embeddings. Figure 12 depicts images rendered by NeRFs using average appearance embeddings. Even if the NeRF recovers the shape of the spacecraft, it can not render the diversity of the illumination conditions encountered on the train set. On

both *sunlamp* and *lightbox**, the SPE trained using interpolated appearance embeddings outperforms the one trained with the average appearance.

6 CONCLUSIONS

This paper addressed the problem of estimating the pose, *i.e.* position and orientation, of an unknown target spacecraft relative to a chaser, in the context of autonomous Rendezvous and Proximity Operations. For this purpose, we introduced a method that enables the training of an off-the-shelf spacecraft pose estimation network from a limited set of real images depicting the target. To this end, our method resorts to an in-the-wild NeRF, *i.e.* a Neural Radiance Field that uses learnable appearance embeddings to capture the varying appearance conditions encountered on real images, which is trained on this sparse set of images. It is then used to generate a large training set which is finally used to train the pose estimation network.

The method was validated on the Hardware-In-the-Loop images of SPEED+ [1]. We showed that a pose estimation network trained with our approach does not only perform much better than the baseline solution, which consists in training the network directly on the real images, but also achieves a similar accuracy as the same network trained on images generated using the CAD model of the target. In addition, we highlighted the role of the in-the-wild abilities of the Neural Radiance Fields in the efficiency of the proposed method.

ACKNOWLEDGEMENTS

The research was funded by Aerospacelab and the Walloon Region through the Win4Doc program. Christophe De Vleeschouwer is a Research Director of the Fonds de la Recherche Scientifique - FNRS. Computational resources have been provided by the Consortium des Équipements de Calcul Intensif (CÉCI), funded by the Fonds de la Recherche Scientifique de Belgique (F.R.S.-FNRS) under Grant No. 2.5020.11 and by the Walloon Region.

REFERENCES

- [1] Tae Ha Park, Marcus Märtens, Gurvan Lecuyer, Dario Izzo, and Simone D’Amico. Speed+: Next-generation dataset for spacecraft pose estimation across domain gap. In *2022 IEEE Aerospace Conference (AERO)*, pages 1–15. IEEE, 2022.
- [2] A Rossi, L Anselmo, A Cordelli, P Farinella, and C Pardini. Modelling the evolution of the space debris population. *Planetary and Space Science*, 46(11-12):1583–1596, 1998.
- [3] Donald J Kessler, Nicholas L Johnson, JC Liou, and Mark Matney. The kessler syndrome: implications to future space operations. *Advances in the Astronautical Sciences*, 137(8):2010, 2010.
- [4] Jason L Forshaw, Guglielmo S Aglietti, Simon Fellowes, Thierry Salmon, Ingo Retat, Alexander Hall, Thomas Chabot, Aurélien Pisseloup, Daniel Tye, Cesar Bernal, et al. The active space debris removal mission removedebris. part 1: From concept to launch. *Acta Astronautica*, 168:293–309, 2020.
- [5] Guglielmo S Aglietti, Ben Taylor, Simon Fellowes, Thierry Salmon, Ingo Retat, Alexander Hall, Thomas Chabot, Aurélien Pisseloup, Christopher Cox, A Mafficini, et al. The active space debris removal mission removedebris. part 2: In orbit operations. *Acta Astronautica*, 168:310–322, 2020.
- [6] Mithun Poozhivil, Manu H Nair, Mini C Rai, Alexander Hall, Connor Meringolo, Mark Shilton, Steven Kay, Danilo Forte, Martin Sweeting, Nikki Antoniou, et al. Active debris removal: A review and case study on leopard phase 0-a mission. *Advances in Space Research*, 2023.
- [7] Simone D’Amico, Mathias Benn, and John L Jørgensen. Pose estimation of an uncooperative spacecraft from actual space imagery. *International Journal of Space Science and Engineering* 5, 2(2):171–189, 2014.
- [8] Leo Pauly, Wassim Rharbaoui, Carl Shneider, Arunkumar Rathinam, Vincent Gaudillère, and Djamilia Aouada. A survey on deep learning-based monocular spacecraft pose estimation: Current state, limitations and prospects. *Acta Astronautica*, 2023.
- [9] Tae Ha Park and Simone D’Amico. Robust multi-task learning and online refinement for spacecraft pose estimation across domain gap. *Advances in Space Research*, 2023.
- [10] Bo Chen, Jiewei Cao, Alvaro Parra, and Tat-Jun Chin. Satellite pose estimation with deep landmark regression and nonlinear pose refinement. In *Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops*, pages 0–0, 2019.
- [11] Pedro F Proença and Yang Gao. Deep learning for spacecraft pose estimation from photorealistic rendering. In *2020 IEEE International Conference on Robotics and Automation (ICRA)*, pages 6007–6013. IEEE, 2020.
- [12] 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. *Communications of the ACM*, 65(1):99–106, 2021.
- [13] Sara Fridovich-Keil, Giacomo Meanti, Frederik Rahbæk Warburg, Benjamin Recht, and Angjoo Kanazawa. K-planes: Explicit radiance fields in space, time, and appearance. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12479–12488, 2023.
- [14] Ricardo Martin-Brualla, Noha Radwan, Mehdi SM Sajjadi, Jonathan T Barron, Alexey Dosovitskiy, and Daniel Duckworth. Nerf in the wild: Neural radiance fields for unconstrained photo collections. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7210–7219, 2021.
- [15] Yefei Huang, Zexu Zhang, Hutao Cui, and Liang Zhang. A low-dimensional binary-based descriptor for unknown satellite relative pose estimation. *Acta Astronautica*, 181: 427–438, 2021.
- [16] Xiaoyuan Ren, Libing Jiang, and Zhuang Wang. Pose estimation of uncooperative unknown space objects from a single image. *International Journal of Aerospace Engineering*, 2020:1–9, 2020.

- [17] Tae Ha Park and Simone D'Amico. Rapid abstraction of spacecraft 3d structure from single 2d image. In *AIAA SCITECH 2024 Forum*, page 2768, 2024.
- [18] Vincenzo Pesce, Michèle Lavagna, and Riccardo Bevilacqua. Stereovision-based pose and inertia estimation of unknown and uncooperative space objects. *Advances in Space Research*, 59(1):236–251, 2017.
- [19] Qian Feng, Zheng H Zhu, Quan Pan, and Yong Liu. Pose and motion estimation of unknown tumbling spacecraft using stereoscopic vision. *Advances in Space Research*, 62(2):359–369, 2018.
- [20] Jiawei Guo, Yucheng He, Xiaozhi Qi, Guangxin Wu, Ying Hu, Bing Li, and Jianwei Zhang. Real-time measurement and estimation of the 3d geometry and motion parameters for spatially unknown moving targets. *Aerospace Science and Technology*, 97:105619, 2020.
- [21] Kai Matsuka, Angel Santamaria-Navarro, Vincenzo Capuano, Alexei Harvard, Amir Rahmani, and Soon-Jo Chung. Collaborative pose estimation of an unknown target using multiple spacecraft. In *2021 IEEE Aerospace Conference (50100)*, pages 1–11. IEEE, 2021.
- [22] Chris Harris, Mike Stephens, et al. A combined corner and edge detector. In *Alvey vision conference*, volume 15, pages 10–5244. Citeseer, 1988.
- [23] David G Lowe. Distinctive image features from scale-invariant keypoints. *International journal of computer vision*, 60:91–110, 2004.
- [24] Ethan Rublee, Vincent Rabaud, Kurt Konolige, and Gary Bradski. Orb: An efficient alternative to sift or surf. In *2011 International conference on computer vision*, pages 2564–2571. Ieee, 2011.
- [25] Mate Kisantal, Sumant Sharma, Tae Ha Park, Dario Izzo, Marcus Märten, and Simone D'Amico. Satellite pose estimation challenge: Dataset, competition design, and results. *IEEE Transactions on Aerospace and Electronic Systems*, 56(5):4083–4098, 2020.
- [26] Tae Ha Park, Marcus Märten, Mohsi Jawaid, Zi Wang, Bo Chen, Tat-Jun Chin, Dario Izzo, and Simone D'Amico. Satellite pose estimation competition 2021: Results and analyses. *Acta Astronautica*, 204:640–665, 2023.
- [27] Sumant Sharma, Connor Beierle, and Simone D'Amico. Pose estimation for non-cooperative spacecraft rendezvous using convolutional neural networks. In *2018 IEEE Aerospace Conference*, pages 1–12. IEEE, 2018.
- [28] Tae Ha Park, Sumant Sharma, and Simone D'Amico. Towards robust learning-based pose estimation of noncooperative spacecraft. *arXiv preprint arXiv:1909.00392*, 2019.
- [29] Antoine Legrand, Renaud Detry, and Christophe De Vleeschouwer. End-to-end neural estimation of spacecraft pose with intermediate detection of keypoints. In *European Conference on Computer Vision*, pages 154–169. Springer, 2022.
- [30] Vincent Lepetit, Francesc Moreno-Noguer, and Pascal Fua. Ep n p: An accurate o (n) solution to the p n p problem. *International journal of computer vision*, 81:155–166, 2009.
- [31] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *ACM Transactions on Graphics (ToG)*, 41(4):1–15, 2022.
- [32] Michal Adamkiewicz, Timothy Chen, Adam Caccavale, Rachel Gardner, Preston Culbertson, Jeannette Bohg, and Mac Schwager. Vision-only robot navigation in a neural radiance world. *IEEE Robotics and Automation Letters*, 7(2):4606–4613, 2022.
- [33] Antoni Rosinol, John J Leonard, and Luca Carlone. Nerf-slam: Real-time dense monocular slam with neural radiance fields. In *2023 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 3437–3444. IEEE, 2023.
- [34] Nianchen Deng, Zhenyi He, Jiannan Ye, Budmonde Duinkharjav, Praneeth Chakravarthula, Xubo Yang, and Qi Sun. Fov-nerf: Foveated neural radiance fields for virtual reality. *IEEE Transactions on Visualization and Computer Graphics*, 28(11):3854–3864, 2022.
- [35] Thang-Anh-Quan Nguyen, Amine Bourki, Mátyás Macduzinski, Anthony Brunel, and Mohammed Bennamoun. Semantically-aware neural radiance fields for visual scene understanding: A comprehensive review. *arXiv preprint arXiv:2402.11141*, 2024.
- [36] Lin Yen-Chen, Pete Florence, Jonathan T Barron, Alberto Rodriguez, Phillip Isola, and Tsung-Yi Lin. inerf: Inverting neural radiance fields for pose estimation. In *2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 1323–1330. IEEE, 2021.
- [37] Lorenzo Giusti, Josue Garcia, Steven Cozine, Darrick Suen, Christina Nguyen, and Ryan Alimo. Marf: Representing mars as neural radiance fields. In *European Conference on Computer Vision*, pages 53–65. Springer, 2022.
- [38] Anne Mergy, Gurvan Lecuyer, Dawa Derksen, and Dario Izzo. Vision-based neural scene representations for spacecraft. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2002–2011, 2021.
- [39] Basilio Caruso, Trupti Mahendrakar, Van Minh Nguyen, Ryan T White, and Todd Steffen. 3d reconstruction of non-cooperative resident space objects using instant ngp-accelerated nerf and d-nerf. *arXiv preprint arXiv:2301.09060*, 2023.
- [40] Aneesh M Heintz and Mason Peck. Spacecraft state estimation using neural radiance fields. *Journal of Guidance, Control, and Dynamics*, pages 1–14, 2023.
- [41] Eberhard Gill, Simone D'Amico, and Oliver Montenbruck. Autonomous formation flying for the prisma mission. *Journal of Spacecraft and Rockets*, 44(3):671–681, 2007.
- [42] Tae Ha Park, Juergen Bosse, and Simone D'Amico. Robotic testbed for rendezvous and optical navigation: Multi-source calibration and machine learning use cases. *arXiv preprint arXiv:2108.05529*, 2021.
- [43] Yannick Bukschat and Marcus Vetter. Efficientpose: An efficient, accurate and scalable end-to-end 6d multi object pose estimation approach. *arXiv preprint arXiv:2011.04307*, 2020.

- [44] Mingxing Tan and Quoc Le. Efficientnet: Rethinking model scaling for convolutional neural networks. In *International conference on machine learning*, pages 6105–6114. PMLR, 2019.
- [45] Mingxing Tan, Ruoming Pang, and Quoc V Le. Efficientdet: Scalable and efficient object detection. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10781–10790, 2020.
- [46] Philip TG Jackson, Amir Atapour Abarghouei, Stephen Bonner, Toby P Breckon, and Boguslaw Obara. Style augmentation: data augmentation via style randomization. In *CVPR workshops*, volume 6, pages 10–11, 2019.
- [47] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International conference on machine learning*, pages 448–456. pmlr, 2015.