



City Research Online

City, University of London Institutional Repository

Citation: Rondao, D., He, L. & Aouf, N. (2024). AI-based monocular pose estimation for autonomous space refuelling. *Acta Astronautica*, 220, pp. 126-140. doi: 10.1016/j.actaastro.2024.04.003

This is the accepted version of the paper.

This version of the publication may differ from the final published version.

Permanent repository link: <https://openaccess.city.ac.uk/id/eprint/32693/>

Link to published version: <https://doi.org/10.1016/j.actaastro.2024.04.003>

Copyright: City Research Online aims to make research outputs of City, University of London available to a wider audience. Copyright and Moral Rights remain with the author(s) and/or copyright holders. URLs from City Research Online may be freely distributed and linked to.

Reuse: Copies of full items can be used for personal research or study, educational, or not-for-profit purposes without prior permission or charge. Provided that the authors, title and full bibliographic details are credited, a hyperlink and/or URL is given for the original metadata page and the content is not changed in any way.

City Research Online:

<http://openaccess.city.ac.uk/>

publications@city.ac.uk

Journal Pre-proof

AI-based monocular pose estimation for autonomous space refuelling

Duarte Rondao, Lei He, Nabil Aouf

PII: S0094-5765(24)00201-7

DOI: <https://doi.org/10.1016/j.actaastro.2024.04.003>

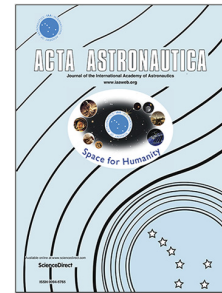
Reference: AA 10385

To appear in: *Acta Astronautica*

Received date: 15 November 2023

Revised date: 29 February 2024

Accepted date: 2 April 2024



Please cite this article as: D. Rondao, L. He and N. Aouf, AI-based monocular pose estimation for autonomous space refuelling, *Acta Astronautica* (2024), doi: <https://doi.org/10.1016/j.actaastro.2024.04.003>.

This is a PDF file of an article that has undergone enhancements after acceptance, such as the addition of a cover page and metadata, and formatting for readability, but it is not yet the definitive version of record. This version will undergo additional copyediting, typesetting and review before it is published in its final form, but we are providing this version to give early visibility of the article. Please note that, during the production process, errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

© 2024 Published by Elsevier Ltd on behalf of IAA.

AI-based Monocular Pose Estimation for Autonomous Space Refuelling

Duarte Rondao^{1,*}, Lei He², Nabil Aouf³

City, University of London, ECV1 0HB London, United Kingdom

Abstract

Cameras are rapidly becoming the choice for on-board sensors towards space rendezvous due to their small form factor and inexpensive power, mass, and volume costs. When it comes to docking, however, they typically serve a secondary role, whereas the main work is done by active sensors such as lidar. This paper documents the development of a proposed AI-based (artificial intelligence) navigation algorithm intending to mature the use of on-board visible wavelength cameras as a main sensor for docking and on-orbit servicing (OOS), reducing the dependency on lidar and greatly reducing costs. Specifically, the use of AI enables the expansion of the relative navigation solution towards multiple classes of scenarios, e.g., in terms of targets or illumination conditions, which would otherwise have to be crafted on a case-by-case manner using classical image processing methods. Multiple convolutional neural network (CNN) backbone architectures are benchmarked on synthetically generated data of docking manoeuvres with the International Space Station (ISS), achieving position and attitude estimates close to 1% range-normalised and 1 deg, respectively, an established rule of thumb for the navigation measurement accuracy during final approach. The

*Corresponding author

¹Postdoctoral Research Fellow with the Department of Electrical and Electronic Engineering (email address: duarte.rondao@city.ac.uk).

²Former Postdoctoral Research Fellow with the Department of Electrical and Electronic Engineering (currently a PhD candidate at Northwestern Polytechnical University, Xi'an, China).

³Professor of Robotics and Autonomous Systems with the Department of Electrical and Electronic Engineering.

integration of the solution with a physical prototype of the refuelling mechanism is validated in laboratory using a robotic arm to simulate a berthing procedure.

Keywords: AI, deep learning, spacecraft, navigation, docking and berthing

List of Acronyms

This document is incomplete. The external file associated with the glossary ‘acronym’ (which should be called `aa10385-ftft.acr`) hasn’t been created.

Check the contents of the file `aa10385-ftft.acn`. If it’s empty, that means
 5 you haven’t indexed any of your entries in this glossary (using commands like `\gls` or `\glsadd`) so this list can’t be generated. If the file isn’t empty, the document build process hasn’t been completed.

Try one of the following:

- Add `automake` to your package option list when you load `glossaries-extra.sty`.

10 For example:

```
\usepackage[automake]{glossaries-extra}
```

- Run the external (Lua) application:

```
makeglossaries-lite.lua "aa10385-ftft"
```

- Run the external (Perl) application:

15 `makeglossaries "aa10385-ftft"`

Then rerun \LaTeX on this document.

This message will be removed once the problem has been fixed.

1. Introduction

For the majority of the 64-year history of space launches, satellites have been
 20 seen as an expendable medium: once the propellant is depleted, the mission is ended. Northrop Grumman’s [Mission Extension Vehicle \(MEV\)](#) programme has recently challenged this paradigm by achieving the first teleoperated [on-orbit](#)

servicing (OOS) to reposition existing spacecraft. This has opened up a new market segment where the spacecraft's life cycle can be extended beyond its original planning, avoiding the costs of launching, manufacturing, and keeping a new one. The success of the mission has attracted the attention of the United States Department of Defence, which has awarded the company a contract to study the possibility of servicing commercial and government satellites using robotics technology [7].

Despite its breakthrough, it is incontrovertible that the MEV was built on the shoulders of previous demonstrators: for example, the Kepler Automated Transfer Vehicle's (ATV) first refuelling operation of the International Space Station (ISS) to supply the station's thrusters in 2011 [3], or NASA's Robotic Refuelling Mission which demonstrated the technology to refuel satellites in orbit by robotic means in 2014 [20]. More recently, ESA has also recognised the potential for this new market segment by opening up calls for ideas related to OOS, having previously invested M€50 in support for research and development of relevant technologies [1]. Overall, OOS and manufacturing alone is projected to have a cumulative global market size of over B\$4.4 by 2030 [2]. Still, existing satellites were, and still are, built without thinking of their serviceability and, more specifically, refuelling, which is a fulcral part of the servicing operations of these assets and represents a significant cost saving measure.

Paramount to the safe accomplishment of refuelling, and OOS in general, is the estimation of the relative states between the service vehicle (SV) and the client or target vehicle (TV) during docking or berthing. At such small distances, this entails the estimation of the six degree-of-freedom (DOF) relative pose, which is typically achieved with two types of optical sensors: lidar and camera sensors [36]. However, current flight-proven solutions using either sensor require optical corner-cube reflectors to be mounted on the TV [9]. The viability of large-scale OOS involves rendezvous and docking or berthing (RVD/B) via autonomous navigation with minimal or no human input, and the cost associated with active sensors can hinder the massification of orbital servicers. Indeed, cameras are already rapidly becoming the choice for on-board sensors towards

spacecraft rendezvous (RV) due to their small form factor and inexpensive power,
 55 mass, and volume costs. The past three years have witnessed a consolidation of
 AI-based (artificial intelligence) techniques for RV, particularly through deep
 learning (DL), using monocular cameras which do not make assumptions about
 the level of cooperation of the target [29]. However, this has not yet been
 established as a navigation approach for docking.

60 2. Related Work

The use of vision-based sensors (VBSs) for RVD/B has traditionally consisted of
 establishing geometric relationships derived from the laws of imaging on the focal
 plane of a lens [9]: the SV illuminates retro-reflectors on the TV-side, which are
 configured according to a known pattern and are then imaged by the on-board
 65 camera. By knowing this configuration (i.e., the relative distances between the
 pattern markers), and the intrinsic parameters of a calibrated VBS (i.e., the
 field of view [FOV] and focal length), information on the SV-TV range, line of
 sight direction, and relative attitude can be computed by detecting said markers
 on each image. In the context of RVD/B, navigation requires the estimation
 70 of said quantities, which make up the 6-DOF relative pose T_{bt} mapping the
 target vehicle frame of reference \mathcal{F}_t to the service vehicle frame \mathcal{F}_b (Fig. 1). This
 entails the need for relative navigation sensors which, in the case of a VBS, define
 two extra frames: the physical camera frame \mathcal{F}_c (which can often be assumed
 coincident with \mathcal{F}_b without loss of generality) and the image plane frame \mathcal{F}_Π
 75 containing the image of the TV and where the image processing (IP) tasks occur.

In the computer vision literature, this problem is called the perspective-*n*-
 point (PnP [31, 10]) problem. PnP can be used with only four markers to retrieve
 the full relative pose, although additional markers can be used to robustify the
 solution. Traditionally, cooperative spacecraft markers have been designed with
 80 concentric patterns of varying sizes to cover different operational sub-ranges
 and maintain acceptable resolution. While fiducial marker designs can be small,
 for long-range targets, multi-reflector spots are sometimes needed due to the

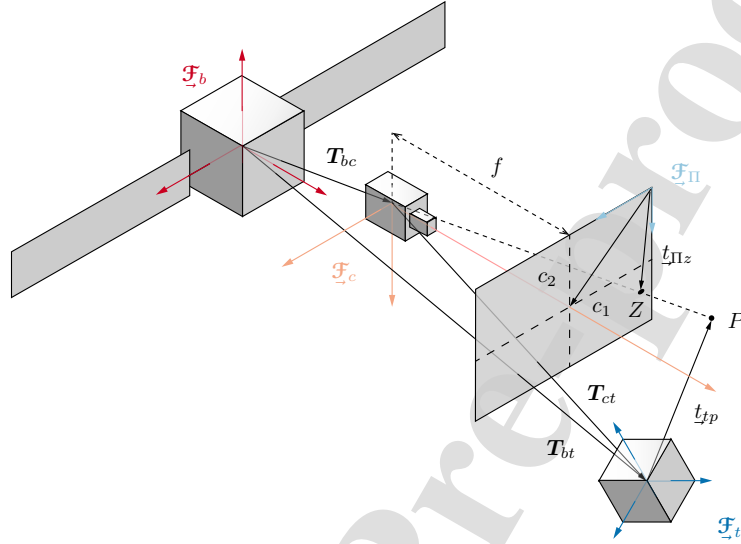


Figure 1: Frames of reference involved in the relative pose estimation problem for RVD/B [25].

low power density of the return signal [9]. Recently, there has been renewed interest in cooperative pose estimation literature, driven by the anticipation of autonomously approached spacecraft in future OOS missions. ArUco markers, popular in robot navigation and augmented reality, have been proposed as an alternative for docking operations [32]. These have the advantage a single marker being sufficient in providing the number of correspondences needed to retrieve the relative pose. Despite their simplicity, though, multiple ArUcos of different sizes are often used for robustness, as a reliable marker detection pipeline is crucial for accurate PnP-based pose estimation [34]. Furthermore, cooperative markers cannot be applied to existing satellites under missions which had not originally been designed with servicing in mind.

In contrast, modern AI techniques, namely through the advent of DL and deep neural networks (DNNs), have gained resurgence in the field of computer vision from the beginning of the previous decade onward, due to advances in

commercial-off-the-shelf [graphics processing units \(GPUs\)](#) and accessibility to large-scale datasets such as ImageNet [8]. The eruption in popularity of [DL](#) arguably occurred in 2012 with AlexNet [18], a [convolutional neural network \(CNN\)](#) which won the ImageNet Large Scale Visual Recognition Challenge with a top-5 classification error more than almost 11 percent points lower than the runner-up; the novel use of [GPU](#)-based training massively accelerated the process, enabling deep learning to be competitive. [CNNs](#), i.e., [DNNs](#) tailored to process image inputs through efficient convolution kernels, later became the norm, as new designs have competed in the challenge each year, resulting in exponential classification score improvements. Two notable examples are GoogLeNet [30], marked not only by a very deep architecture, but also by the implementation of parallel layers to extract multi-scale features; and ResNet [12], which introduced residual connections allowing the breakthrough to even deeper architectures. Most of these state-of-the-art [CNNs](#) have been open-sourced and model weights made available from pre-training on image classification tasks using the ImageNet dataset, which has since contributed to the swift evolution of [DL](#) in general and in the adoption of such models as [CNN](#) front-ends.

It took more than five years for the popularity of [CNNs](#) to migrate onto the domain of spacecraft relative pose estimation for [RV](#). In 2019, ESA Kelvins' [Spacecraft Pose Estimation Challenge \(SPEC\)](#) benchmarked estimation errors obtained on assorted image inputs taken with on-board [VBS](#) from a simulated [RV](#) with the Tango spacecraft under random poses [17]. Although it did not tackle docking or berthing, it did demonstrate the good performance of [AI](#)-based approaches in vision-based navigation for space. The proposed techniques are model-based since they operate under the assumption that a priori knowledge of the target's aspect and structure is available in the form of training images. Since the target does not feature any physical fiducial or retro-reflective markers, only "natural" features derived from shape and texture, they become appropriate for uncooperative scenarios.

The relative pose can be estimated in an end-to-end fashion by following a direct approach (Fig. 2, top). This involves designing a [DNN](#) which generates a

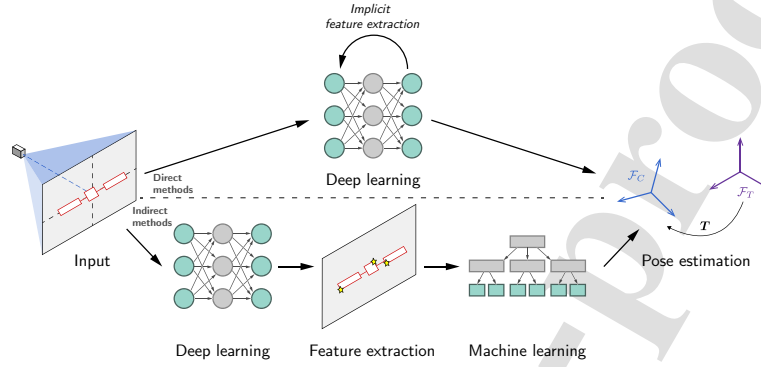


Figure 2: Direct versus indirect methods for DL-based pose estimation [29].

pose directly from image inputs, and has been tested for close-range rendezvous on SPEC [22]. It has the advantage of not relying on additional machine learning (ML) pipelines to generate the desired solution. Moreover, it entails a natural framework for the implementation of DL-based temporal modelling, which has been shown to improve the solution for time-series where there is a correlation between successive relative poses [27].

Interestingly, most of the highest scoring SPEC competitors [5] followed a so-called indirect approach towards DL-based pose estimation (Fig. 2, bottom): this entails relaying the use of a CNN entirely to the task of extracting features (usually 2D keypoints), and then using a PnP solver, or optimizer techniques such as Levenberg-Marquardt, to retrieve the solution from their correspondences to 3D control points. These are defined by the user at the pre-training stage and can be selected from the surface of a computer aided design (CAD) model of the target. As outlined by Kisantal et al. [17], the optimization step applied to the 2D-3D correspondences output by these networks enables a refinement of the solution leading in turn to a more accurate pose estimate, explaining their increased performance on the SPEC dataset. As such, indirect methods have since become the preferred choice for DL-based pose estimation in RV. The potential of cameras as cost-effective navigation sensors for uncooperative RV

has actually been demonstrated previously through this prism using traditional IP techniques under both model-based [26] and model-free [4] assumptions, as well as in combination with other economical hardware such as range finders [35]. However, DNNs in general benefit from representation learning, thus eliminating the need for manually selecting the most appropriate type of feature, the optimality of which would potentially be linked to each specific scenario or dataset.

To the best of our knowledge, neither technique has been evaluated on docking scenarios, thus making the **Orbital AI-based Autonomous Refuelling (OIBAR)** project the first to demonstrate the viability of deep neural networks using vision-only inputs for pose estimation of a docking structure. When designing the proposed navigation pipeline, a direct approach was chosen over an indirect one for two main reasons. Firstly, as the docking phase requires continuous estimation of relative states, there is motivation to explore the influence of temporal modelling using a direct framework. While a similar argument could be made for the rendezvous stage, the incentive for single-image pose estimation—often studied with indirect methods in the literature—is greater in this context than in the docking stage. Secondly, it can also be argued that this avenue can lead to a stronger decoupling from the camera intrinsics. Indirect methods may experience difficulties regressing on precise keypoint locations at larger distances if the camera resolution is not sufficiently high, and at shorter distances a subset of these may fall outside the FOV. On the other hand, direct methods are end-to-end and can thus rely on different elements of the input image in such situations.

3. Methodology and Design

This section details the approach followed for the execution of the **OIBAR** project. The design of the docking mechanism has originally been detailed by He et al. [13]; a rundown of the refuelling operations is reiterated below for context. Then, the AI-based navigator introduced in this paper is described.

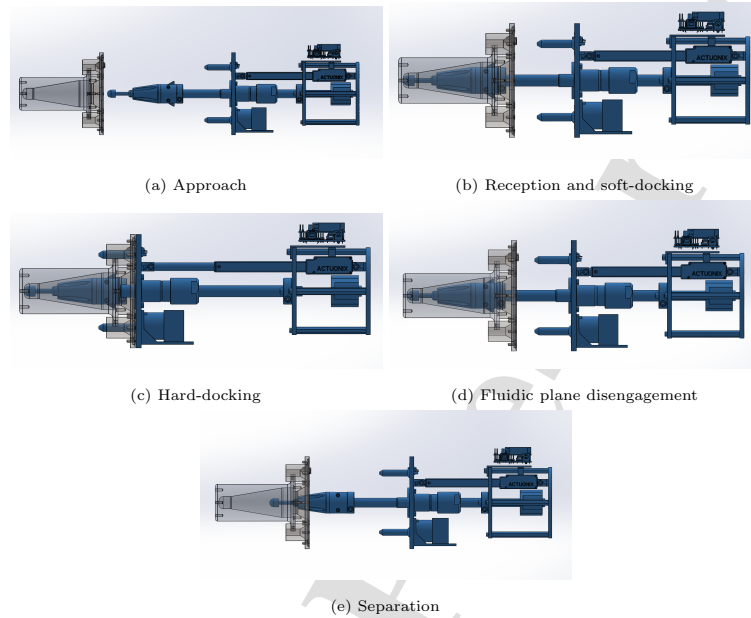


Figure 3: Refuelling operations between TV and SV using the designed docking mechanism for OIBAR [13].

3.1. Refuelling Operations

The refuelling procedure in OIBAR can be broken down into the following sequence of operations, which are also illustrated in Figure 3.

Initially, the AI-guided SV approaches the TV, adjusting velocity and aligning docking interfaces (Fig. 3a). The reception phase then begins with the end-effector probe entering the drogue, utilising a spring-damper to absorb shock and provide retraction force. Soft-docking is facilitated by two spring-loaded latches, preventing accidental detachment during probe entry into the drogue's varying diameter cavity (Fig. 3b). After soft-docking, the end-effector becomes restricted for hard-docking, which in turn ensures precise alignment for fluidic plane connection via two alignment pins engaged in the guide cavities on the berthing fixture side. (Fig. 3c). Fuel transfer is initiated, and after reaching the target

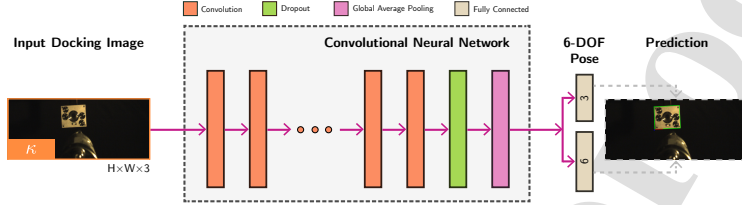


Figure 4: OibarNet base DNN architecture for AI-based docking navigation.

pressure, valves are closed. Electrical connectors transmit control commands and pressure data. Post-refuelling, the end-effector is released, retracting the
 190 fluidic plane and automatically extracting the probe (Fig. 3d). The latch stepper motor resets the latch for the next refuelling process (Fig. 3e).

3.2. Software Development

3.2.1. Architecture

Figure 4 illustrates the base architecture of the AI-based navigation algorithm
 195 for OIBAR, termed OibarNet.

The proposed network is a direct (i.e., end-to-end) DNN taking red-green-blue
 (RGB) images of the TV's berthing fixture at time-step $\tau = k$ and outputting
 the corresponding 6-DOF pose relative to the SV. The front-end and backbone
 of OibarNet is the CNN that processes the image inputs (Fig. 4, in orange).
 200 Multiple CNN model candidates are considered and evaluated within the project
 (see Section 4); however, two architectural aspects are kept constant. The first
 is a dropout layer (Fig. 4, in green) added after the last convolutional layer to
 prevent overfitting [14]. The second is a global average pooling (GAP) layer
 (Fig. 4, in pink). GAP converts the CNN output to a fixed-dimension vector
 205 dependant only on the number of output channels, regardless of the image inputs
 spatial dimensions. This allows OibarNet to work with large input resolutions
 without needing to increase the network depth, and to potentially train it on
 different datasets without needing to alter the architecture. The CNN-processed
 features are then subject to a fully connected (FC) layer back-end which estimates

210 the pose via regression.

As part of the project, the temporal modelling of the CNN features is also considered. This is investigated via the inclusion of a recurrent neural network (RNN) between the GAP output and the FC input.

3.2.2. Relative Pose Representation

215 The first FC block head maps the CNN output to a 3-vector estimate of the position, whereas the second head maps it to the 6-dimensional estimate of the attitude formulated by Zhou et al. [37]. This representation, denoted by \mathbf{r} , is obtained from the direction cosine matrix \mathbf{R} through the mapping:

$$\text{SO}(3) \rightarrow \mathbb{R}^6$$

$$\mathbf{R} \mapsto \left[-\mathbf{R}_{:,1}^\top, -\mathbf{R}_{:,2}^\top \right] = \mathbf{r}^\top, \quad (1)$$

i.e., discarding the third column of \mathbf{R} and stacking the result. The transformation $\mathbf{r} \mapsto \mathbf{R}$ involves in turn reshaping \mathbf{r} into a 3×2 matrix followed by a Gram-Schmidt orthogonalisation. The attitude representation warrants special attention, as the 4-dimensional quaternion, \mathbf{q} , is normally used to represent the attitude of a spacecraft due to its low dimensionality and lack of singularities. However, its antipodal ambiguity-induced discontinuities (i.e., $\mathbf{q} = -\mathbf{q}$) have 225 been shown to yield sub-optimal results in a deep learning environment compared to the 6D representation — further details are given in Ref. [37]. For post-processing or error quantification, \mathbf{q} can then be obtained from \mathbf{R} using well-known isomorphisms [19].

3.2.3. Loss Function

230 The combination of predicted position and attitude quantities in a single loss function requires the incorporation of a scaling factor since these two quantities normally deal in different magnitudes [22]. Typically, this scaling factor has been considered a hyperparameter part of the DNN's tuning process, which is sub-optimal.

235 In contrast, OibarNet follows the approach of Cipolla et al. [6] and attributes one weight each to the position and attitude, σ_t and σ_r , respectively, which become learnables and converge during the training process. The weights represent the task-specific variances of two Gaussian distributions, yielding a combined L^2 norm loss:

$$\mathcal{L} = \mathcal{L}_r \exp(-2\hat{\sigma}_r) + \mathcal{L}_t \exp(-2\hat{\sigma}_t) + 2(\hat{\sigma}_r + \hat{\sigma}_t), \quad (2)$$

240 where

$$\mathcal{L}_r = \sum_{i=1}^B \|\hat{\mathbf{r}}^{(i)} - \mathbf{r}^{(i)}\|, \quad \mathcal{L}_t = \sum_{i=1}^B \|\hat{\mathbf{t}}^{(i)} - \mathbf{t}^{(i)}\|. \quad (3)$$

Here, $\hat{\mathbf{r}}, \hat{\mathbf{t}}$ are the 6D attitude and 3D position estimated by the network, respectively; \mathbf{r}, \mathbf{t} are the corresponding ground truths; $\|\cdot\|$ denotes the L^2 norm; and B is the batch size.

4. Demonstration and Testing

245 In this section, the methodology adopted for demonstrating the reliability of the developed AI-based solution for space docking and refuelling under OIBAR. The validation tests are divided into three fronts: hardware validation, software validation, and integration validation. The former has been reported by He et al. [13], the latter two are introduced herein.

250 4.1. Software Validation

Supervised ML algorithms require labelled and well-structured datasets not only for evaluation, but also for training. This is especially true, and even more relevant, for DL-based methods, which require large and diverse batches of data to learn how to generalise towards unseen scenarios due to the very large number

255 of parameters at play.

However, labelled datasets for spacecraft pose estimation are scarce and expensive to obtain due to the intricate environmental conditions that must be

emulated. As such, the first step in the software validation campaign for OIBAR is to create a framework that allows for the generation of synthetic data: in contrast to real world sets, synthetic data is generally inexpensive and allows the possibility of having virtually unlimited samples.

Once this simulation environment is defined, the next step entails using it to produce synthetic docking trajectories. Lastly, an architecture selection round is performed based on the estimation performance on the data.

4.1.1. Simulation Environment

The simulation environment for OIBAR is composed of two main components: a simulator designed in MATLAB/Simulink to replicate the orbital motion of the SV and TV under the influence of Earth; and an interface with the open-source 3D modelling software Blender⁴ for the purpose of generating synthetic but realistic imagery of the TV as viewed from the SV on-board VBS based on the states computed by the simulator.

The MATLAB/Simulink orbital simulator propagates 6-DOF pose of a body orbiting Earth based on an initial state at a given date, which is obtained from two-line element (TLE) sets. The effects of aerodynamic drag and the nonspherical mass distribution of Earth (J_2 zonal coefficient) are implemented as acting force perturbations. Aerodynamic and gravity gradient torques are implemented as acting torque perturbations. Using the planetary ephemerides blocks available in Simulink, the relative states of Earth and the Sun are also computed.

Figure 5 illustrates the developed orbital trajectory simulator environment. The approach to building the simulator followed a modular design (Fig. 5a) which aimed to create basic blocks with fundamental functions, such as linear algebra, attitude manipulation and kinematics, and orbital mechanics, in order to facilitate any needed changes or customisation to the environment. The front-end (Fig. 5b) allows for a simple configuration of initial states, including

⁴<https://www.blender.org>.

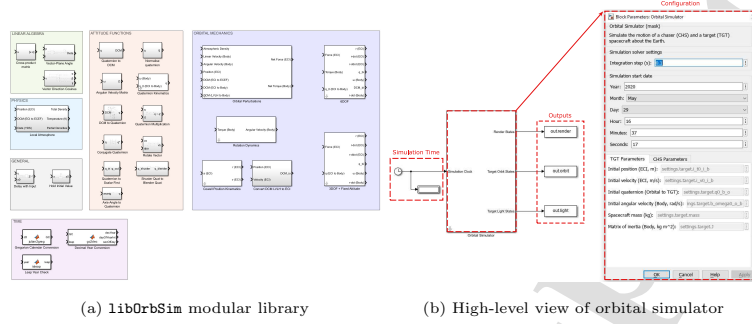


Figure 5: MATLAB/Simulink orbital trajectory simulator environment.

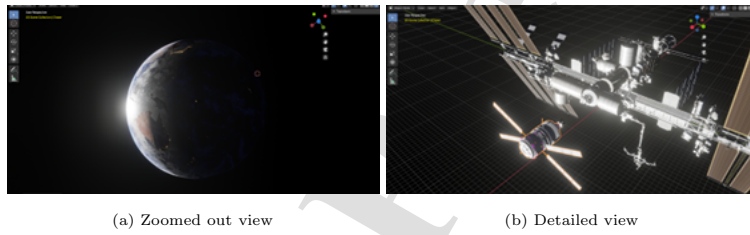


Figure 6: Blender 3D rendering environment.

the date and **TV** pose, permitting configurations for different scenarios.

The propagated states are then saved and interfaced with Blender, allowing the recreation of realistic images of the **RVD/B** dynamics of an Earth-orbiting **SV** and **TV** to be used for algorithm testing. Figure 6 illustrates an example of a simulated scene within Blender. Besides the **TV** and **SV**, the Sun and Earth states have both been imported as well, ensuring a correct correspondence to the simulated time of day and solar phase angles.

4.1.2. Synthetic Dataset Generation

Using the simulator presented in Subsection 4.1.1, a synthetic **VBS** docking dataset was generated to validate the **AI**-based navigation system. The simulated scenario was chosen to be a refuelling of the **ISS** by the Kepler **ATV**. The **CAD** model of the **OIBAR** docking mechanism was imported into the Blender

environment and attached to the vehicles: the end-effector replaced the *ATV*'s existing *RVD/B* system, and the berthing fixture replaced the docking ports on the *ISS*. Due to a size mismatch, the original *OIBAR CADs* were scaled up to accommodate the vehicles within the simulation. In total, six different docking port locations on the *ISS* were considered; the inclusion of several docking ports was deemed beneficial as it grants diversity in terms of backgrounds, approach vectors (i.e., R-bar and V-bar) and illumination conditions.

The MATLAB/Simulink simulator is used to propagate the trajectory of the *ISS* from an initial *TLE* set, as well as the Earth and Sun states. The *ATV SV* guidance trajectories are generated directly relative to the *TV* frame; in particular, relative to the docking port considered for docking. Each *SV* trajectory begins at a relative shaft-plane-berthing-plane distance of 10 m and consists of three parts:

- 1) **Acquisition.** This stage is characterised by a large cross-track motion ($x - y$ plane) whereby the *VBS* is acquiring the target, prior to the end-effector and berthing fixture axes coinciding, but keeping the relative attitudes aligned. The *SV* translates between two randomly generated waypoints at radial distances between 1–2 m from the alignment axis, before reducing this distance to zero. The linear velocity in this stage varies from 0.09–0.12 m s⁻¹.
- 2) **Forced translation.** Once the previous stage is concluded, the *SV* end-effector and *RV* berthing fixture are aligned in terms of a common along-track axis (z -axis), still at a relative distance of 10 m. The *SV* then translates along this axis at a nominal velocity of 0.03 m s⁻¹ to close the distance until 3 m. To add variations amongst sequences, small perturbations are randomly generated and added to the translational and rotational motions; this is achieved by modelling a simple **proportional integral (PI)** controller and generating the next pose state from the feedback error. The magnitudes of the allowed perturbations vary from ± 0.002 m s⁻¹ for the along-track velocity, ± 0.01 m for the cross-track position, and

± 0.1 deg for the attitude; all with a probability of occurrence of 10 %.

3) **Alignment and soft-docking.** This final stage begins with the position-
 330 and attitude-wise alignment of the fluidic and berthing planes from what-
 ever misalignment state the previous stage may have ended in. The **SV**
 then translates towards the **TV**, while keeping an aligned attitude, until
 the spring protruding pin enters the central cavity, the latches are engaged,
 and soft-docking is achieved.

335 The average sequence duration is ~ 5 min, and the synthetic dataset comprises
 12 sequences in total. Furthermore, the synthetic dataset is composed of two
 subsets. The first one, **synthetic/iss**, consists of nominal scenario conditions
 where the docking mechanism is mounted on one of the docking ports of the **ISS**.
 The second subset, **synthetic/perlin**, models similar relative trajectories, but
 340 removes all meshes except for the docking mechanism, replacing the background
 with randomised Perlin noise. The objective of **synthetic/perlin** is to comple-
 ment the **synthetic/iss** subset to provide additional training data and to help
 OibarNet focus on extracting features of the target and ignore the background
 and environment, such as the appearance of Earth behind the target. Table
 345 **1** summarises the characteristics of the generated dataset. The Sun elevation
 angle, i.e., the angle between the Sun direction vector and the orbital velocity
 vector, was selected to ensure daylight conditions for each specific docking port,
 whereby variance on the visual docking conditions was introduced by selecting
 values close to sunrise or sunset periods. Figure 7 illustrates a few sample frames
 350 from the dataset.

The synthetic dataset emulates the **VBS** used in the integration testing (see
 Section 4.4, Table 3); images are generated at a resolution of 744×480 px and a
 framerate of 10 Hz.

4.1.3. Training

355 From Table 1, sequences 1 and 8 were selected exclusively for testing, whereas
 the remaining sequences were used for training.

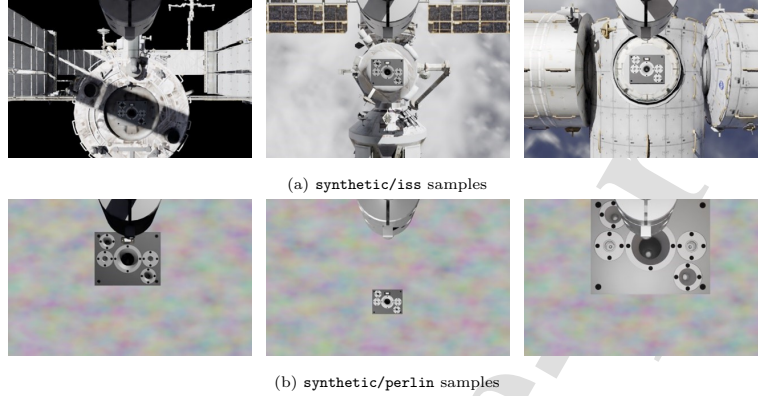


Figure 7: Sample frames from OIBAR's synthetic dataset.

Further, the latter were divided according to a 80%/20% split to include validation data as well and allow the benchmarking of different models. Due to the data consisting of long temporal sequences, these were partitioned into smaller ones to guarantee the proportions of the intended split in an unbiased way. To achieve this, each sub-sequence i 's length in seconds defined by randomly sampling a power of two, n_i , according to the following formula:

$$\text{len}(i) = 2^{n_i}. \quad (4)$$

The values for n_i were sampled from the set of integers $\{6, \dots, 10\}$, resulting in sub-sequence lengths belonging to the set of $\{64, \dots, 1024\}$ seconds.

Image augmentation was performed online on the training data to prevent overfitting (i.e., to boost generalisation performance during inference with new, unseen data). This was applied in two fronts. The first one is related to transform operations in the IP domain: randomised changes in terms of brightness, contrast, colour, Gaussian noise and blur, for example, are generated to robustify the network against potentially unpredicted imaging conditions during deployment. The second one is related to operations in the pose domain, whereby randomly generated perspective transforms are applied to images in the sequence to

Table 1: Synthetic dataset characteristics used for training, validation, and testing of OibarNet.

Sequence	Docking port	Approach axis		Sun elevation angle (deg)	Background		Duration (s)
		V-bar	R-bar		iss	perlin	
1	1	+		37	×		332
2	1	+		75		×	319
3	2		−	56	×		358
4	2		−	146		×	336
5	3	−		127	×		348
6	3	−		165	×		327
7	4		−	56		×	333
8	4		−	146	×		329
9	5	+		56	×		333
10	5		+	146		×	342
11	6		+	56		×	323
12	6		+	146	×		355

simulate deviations in the trajectory (i.e., translation shifts, in-plane rotations, homography-induced off-plane rotations). The latter is of particular importance since, despite the sequence partitioning for training, the forced translation phase dominates each sequence, generating an imbalance on the distribution of position states. The reader is directed towards Rondao [25] for illustrations of the implemented augmentation techniques. Table 2 compiles the quantitative parameters of the randomised transformations used during training.

OibarNet is implemented in MATLAB R2021b using a custom-developed library. Models are trained for 100 epochs with a cyclical learning rate decay of 5 cycles [28]. The Adam optimiser [16] is used. A dropout probability of 0.2 is used. Training is performed on City, University of London’s high performance computing facility Hyperion using one NVIDIA® Quadro RTX™ 8000 GPU with 48 GB VRAM.

4.1.4. Testing

The test results are presented in terms of the position and attitude error metrics, respectively:

Table 2: Image augmentation randomisation parameters used in the training process of OibarNet.

Transform	Parameter		Unit	Description
	Minimum	Maximum		
Channel shift	-20	20	-	Pixel intensity shift value
Gaussian blur	7	13	px	Kernel size
Gaussian noise	3×10^{-3}	1×10^{-2}	-	Variance
JPEG compression	2	8	-	Intensity
Median blur	7	13	px	Kernel size
Patch dropout	10	10	%	Mean image area covered
	3	5	%	Patch size relative to smallest image dimension
Brightness	-0.2	0.2	-	Intensity
Contrast	0.8	1.2	-	Intensity
CLAHE	2	6	-	Pixel intensity clip range
	8	8	-	Number of tiles
Gamma	0.35	1.50	-	Factor
Camera rotation (3 DOF)	-5	5	deg	Magnitude per axis
Image rotation (in-plane)	-5	5	deg	Magnitude
Image translation	-150	150	px	Magnitude

$$\delta \tilde{t} := \|\hat{\mathbf{t}} - \mathbf{t}\|, \quad (5)$$

$$\delta \tilde{q} := 2 \arccos (\hat{\mathbf{q}}^{-1} \otimes \mathbf{q})_4, \quad (6)$$

where the subscript “4” denotes the scalar element of the resulting quaternion.

390 Additionally, the position error is also assessed in terms of the relative range:

$$\delta \tilde{t}_r := \frac{\delta \tilde{t}}{\|\mathbf{t}\|}. \quad (7)$$

It is noted that the pose estimation accuracy requirements for OIBAR have been defined by He et al. [13] as 5% of range for the position and 5 deg for attitude. This is driven by the fact that the functional testing setup is configured to use a robotic arm, i.e., more closely simulating a berthing operation. In such cases, pose estimation accuracy becomes less critical, and the “1%-1 deg” rule

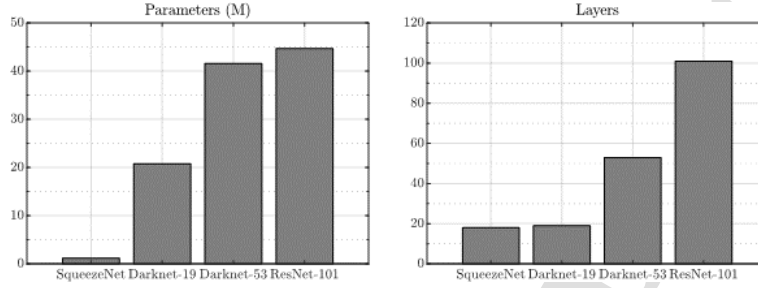


Figure 8: Characteristics of the four different benchmarked CNN models.

of thumb normally applied for the final approach can be relaxed by a factor of 5 [9].

4.1.5. CNN Architecture Selection

From the relative pose estimation point of view, a few differences may be expected between an RV manoeuvre and a docking sequence. Firstly, a reduced variation in the attitude is expected during docking since the SV is expected to be inside the cone-shaped approach corridor of the TV [9]; in opposition, the target may be tumbling during RV. Secondly, an increased apparent variation in the position can be expected during docking, as due to the reduced relative distance any small shift will result in a large displacement of the TV berthing fixture in the FOV.

To better assess the influence of these factors, multiple CNN architectures are benchmarked. The baseline is Darknet-19 [23], which has successfully been applied in the past to the problem of pose estimation in RV [27]. To analyse the effect of increasing the capacity of the model, Darknet-53 [24] and ResNet-101 [12] are included. Lastly, it is also important to verify the change in performance when reducing the capacity, and SqueezeNet [15] is thus included in the benchmark as well.

Figure 8 summarises the number of parameters, in millions, and number of layers of the four different CNN models considered for benchmarking.

4.2. Integration Validation Methodology

The goal of the integration testing is to validate the combination of the hardware and software blocks outlined above. In this setup, the navigation VBS is incorporated into the hardware setup [13] to acquire a stream of images to be processed by the navigation algorithm (Sec. 4.1) during the docking manoeuvre emulated by the robotic setup.

Figure 9 illustrates the integration validation setup. A blackout backdrop is placed behind the target berthing fixture to simulate the imaging conditions of a featureless deep space background. The target itself is illuminated by a single 400 W halogen directional floodlight. The distance at which this illumination source was placed from the target was adjusted to simulate an irradiance of approximately 1361 W m^{-2} , typical for low Earth orbit, under the following rationale. The light is modelled as a point source illuminating a cone with an apex angle θ and generatrix r . These represent, respectively, the beam spread and the distance between the apex (or origin) and the target being illuminated. Let $\theta = 60 \text{ deg}$, which is a typical beam spread for wide flood halogen lamps, such as the one used herein. Given θ , the corresponding solid angle Ω can be calculated as:

$$\Omega = 2\pi \left(1 - \cos\left(\frac{\theta}{2}\right) \right) = 0.8418 \text{ sr.} \quad (8)$$

The solid angle Ω , the irradiance E , the power P , and the distance r are related by the formula:

$$E = \frac{P}{\Omega \cdot r^2}. \quad (9)$$

Solving the above for r and plugging in the remaining values yields a distance of approximately 0.6 m at which the lamp should be positioned relative to the berthing fixture in the experiment. This distance also allows for the floodlight to be conveniently positioned outside of the imaging sensor's FOV mounted on the robotic arm manipulator.

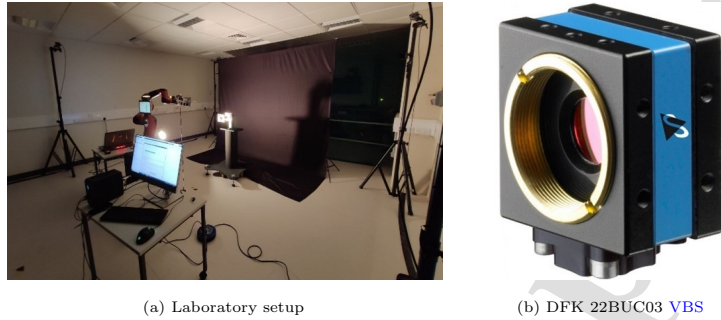


Figure 9: Integration validation setup at City, University of London's [ASML](#).

The [SV](#) and [TV](#) are placed inside the capture volume of an OptiTrack⁵ motion capture system for recording the ground truth measuring approximately $5 \times 5 \times 3$ m. OptiTrack can record 6-[DOF](#) pose data of rigid and flexible bodies by detecting, tracking, and triangulating passive near infrared markers placed

445 on targets. The data can be saved or stream over a local network in real-time.

The OptiTrack setup at City consists in six PrimeX 13 cameras with a resolution of 1280×1024 px running at a native framerate of 240 Hz, capable of achieving positional errors less than ± 0.20 mm and rotational errors less than 0.5 deg.

450 The used [VBS](#) is the Imaging Source DFK 22BUC03 colour camera with a $1/3$ inch format CMOS sensor (Onsemi MT9V024) and a native resolution of 744×480 px, fitted with a Kowa LM4NCL 3.5 mm focal length lens. Table 3 summarises the technical characteristics of the [VBS](#).

The workstation consists of an Intel[®] NUC 9 Pro with an NVIDIA[®]

455 RTX[™] 3060 Ti Mini GPU with 8 GB VRAM. The workstation is used for both experimental data offline validation of OibarNet and real-time online testing of the network, at a framerate of 10 Hz.

⁵<https://optitrack.com>.

Table 3: Technical data – DFK 22BUC03.

Parameter	Units	Value
Resolution	px	744 × 480
Maximum frame rate	Hz	76
Focal length	mm	3.5
Horizontal FOV	deg	65.6
Vertical FOV	deg	44.7

4.2.1. Experimental Dataset Generation

The docking imaging sequences acquired with the experimental setup follow the same structure as the synthetic dataset (Sec. 4.1.2) albeit with two key differences. The first one is that all experimental sequences feature the same type of background (black, deep space). The second is that, rather than implementing PI-induced pose perturbations during the forced translation (Phase 2), a static misalignment of the pose is randomly introduced in each sequence at the beginning of the phase, which is then corrected at the beginning of the final one.

In total, 12 experimental trajectories are collected, whereby the angle of illumination alternates between port and starboard. The average sequence duration is ~ 3.15 min. The first 10 sequences are used for training and validation of the model according to the methodology of Section 4.1.3. Sequences `experimental/11` and `experimental/12` are used exclusively for testing.

4.2.2. Ground Truth Calibration Toolbox

The OptiTrack system used to record the ground truth measures the poses of rigid bodies equipped with infrared markers. However, it does not directly output the relative pose between the VBS and the TV (as illustrated in Fig. 1, Sec. 2), which is required by the navigation algorithm.

To this end, a toolbox was developed in MATLAB to calibrate the output OptiTrack data and generate the required relative pose, based on the work of

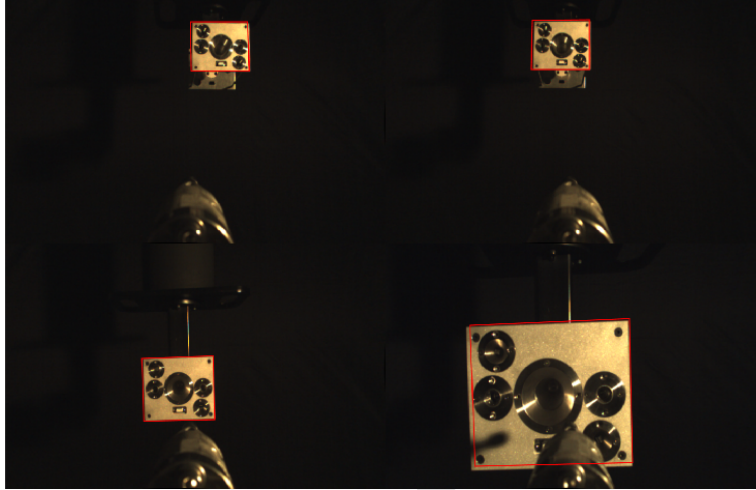


Figure 10: Ground truth calibration toolbox output, visualised on some frames of the experimental dataset by reprojecting the target's CAD model (in red) according to the measured pose.

Valmorbida et al. [33] and Pasqualetto Cassinis et al. [21]. The output of the calibration toolbox are the static transforms T_{ic} , mapping the camera frame \mathcal{F}_c to the frame of reference \mathcal{F}_i defined by the physical markers placed on its housing and tracked by OptiTrack, and T_{sb} , mapping the target's body frame \mathcal{F}_b to the frame of reference \mathcal{F}_s defined by the markers placed on it. These transforms then make possible to map the OptiTrack marker-defined rigid bodies' poses, which are measured relative to \mathcal{F}_o , the system's arbitrary global frame of reference, into a usable ground truth T_{bc} defined in terms of the VBS frame of reference. The overall uncertainty of the framework can be quantified in terms of the total reprojection error between model and image corners and is estimated to be below 20 px for ranges above 2 m [21]. However, this upper bound is expected to be lower for the present application due to the shorter relative ranges considered. Figure 10 illustrates the output of the calibration procedure on select samples of the dataset.

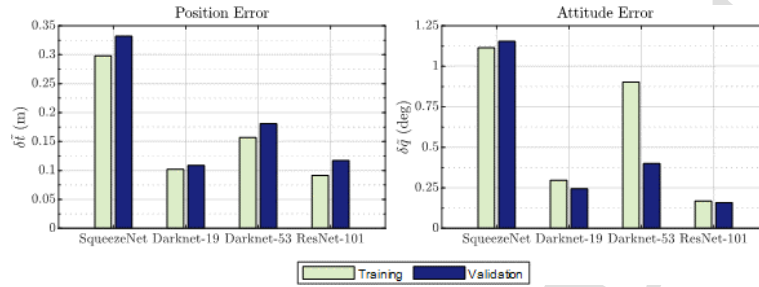


Figure 11: Average training and validation pose estimation errors for different CNN architectures trained on the synthetic dataset.

5. Results

5.1. Software Testing

5.1.1. CNN Model Benchmarking

Figure 11 illustrates the results of the different models trained on the synthetic dataset, presented in terms of mean position and attitude errors averaged per trajectory. It can be seen that the performance of SqueezeNet is considerably worse than the baseline Darknet-19, yielding errors twice as large for both position and attitude. Both networks have a very similar number of layers, but Darknet-19 has substantially more learnable parameters (as indicated in Fig. 8); the reduced attitude variance in the docking manoeuvres is thus shown not to justify a decrease in parameters.

Interestingly, the error for Darknet-53 actually increases with respect to the baseline. Once the capacity of the CNN is further increased with ResNet-101, though, the error decreases again, making the network the best performing model (except on position validation error, which is slightly larger than Darknet-19's).

The results of Figure 11 are presented with the caveat that they represent average errors per trajectory, but where the data is not composed of random images but time sequences. As such, while a histogram visualisation is useful for a first analysis of each model's performance, it is also important to look at how these perform in specific, individual situations. For example, Figure

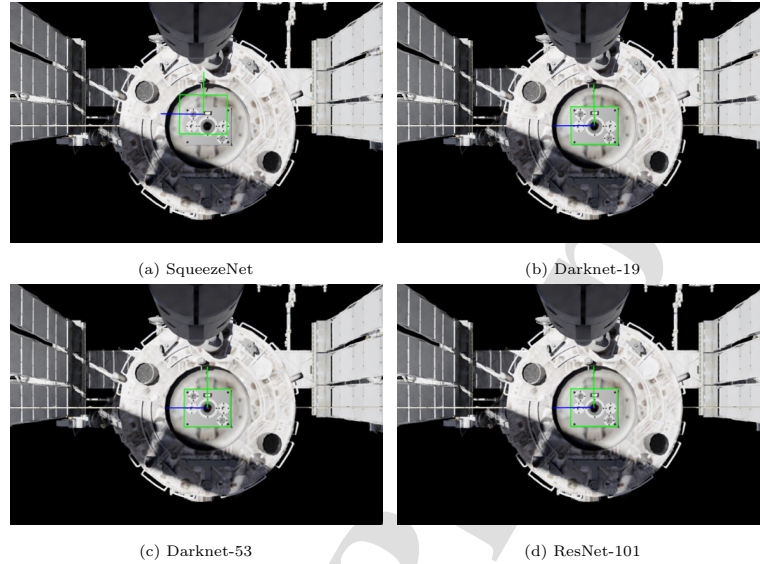


Figure 12: Qualitative pose estimation performance on a validation sequence of the synthetic dataset for different CNN models.

12 represents the qualitative performance of each model on a single frame of one of the synthetic validation sequences; the rectangular boundary of the berthing fixture is reprojected in green using the predicted pose, and the axes of the estimated frame \mathcal{F}_t are also shown. The results show that SqueezeNet is overfitting at least on the position state, as it expects the berthing fixture to be located in the centre of the FOV, when in reality the SV end-effector is still misaligned. The other three models with increased capacity demonstrate no issues in estimating the correct relative position.

520 Consider now, however, the performance on one training sequence, as illustrated in Figure 13: SqueezeNet (a) is shown to be underfitting, but so is Darknet-19 (b). This suggests that increasing the capacity would benefit OibarNet, as confirmed by the frame output by Darknet-53 (c) showing a better fit, despite the summary metrics in Figure 12. The performance with ResNet-101

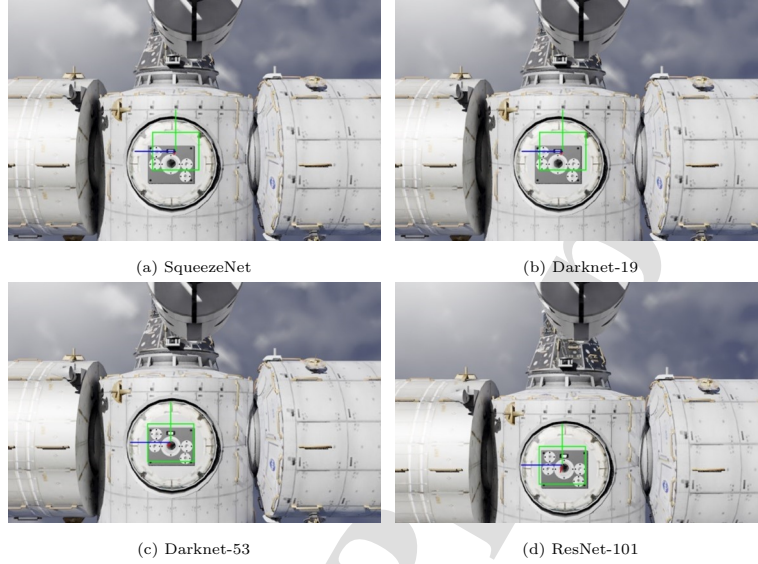


Figure 13: Qualitative pose estimation performance on a training sequence of the synthetic dataset for different CNN models.

525 (d) is slightly better even, confirming it as the choice for the final CNN model in OibarNet.

5.1.2. Performance Evaluation on Synthetic Dataset

In this subsection, the performance of the navigation algorithm is evaluated on the test sequences `synthetic/01` and `synthetic/08`, as outlined in Section 4, and according to the selected ResNet-101 CNN architecture for OibarNet.

530 Figure 14 showcases the attained pose estimation errors for each sequence using the final OibarNet model; the position errors are normalised as a percentage of range. Table 4 summarises these statistics.

The figures demonstrate that, for `synthetic/01`, OibarNet fulfils the 5% maximum range-normalised position error requirement (defined in Ref. [13]) for most of the trajectory. The exception is a segment corresponding to phase 1 (acquisition) whereby the SV moves to a waypoint representing a large displace-

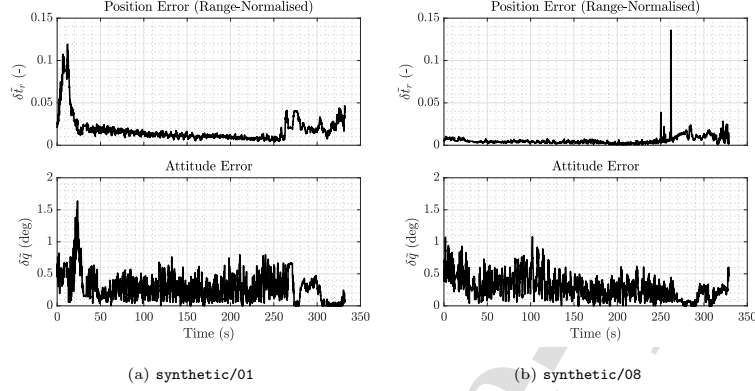


Figure 14: Estimated position and attitude errors over time on the two test sequences of the synthetic dataset.

Table 4: Summary performance statistics on the two test sequences of the synthetic dataset. Position errors are range-normalised. “Std.” denotes standard deviation.

Sequence	Position Error (%)			Attitude Error (deg)			Requirement Compliance (%)	
	Mean	Median	Std.	Mean	Median	Std.	Position	Attitude
synthetic/01	1.81	1.39	1.53	0.29	0.26	0.20	96.33	100
synthetic/08	0.56	0.43	0.50	0.28	0.26	0.17	99.91	100

ment relative to the alignment axis, representing about 3.7% of the sequence’s duration. After this period, the error converges to values below 2.5% of range, further decreasing as the SV closes in on the TV, until the beginning of phase 3 (alignment and soft-docking), where the very short range causes the error to rise, but not above the requirement threshold. The attitude estimation performance is shown to fully comply with the 5 deg maximum error requirement.

The position estimation performance of the navigation algorithm on synthetic/08 is observed to be better than the previous sequence, as an improvement of 1.25 percent points on the mean value and 0.96 percent points on the median value are achieved. Furthermore, the position estimate is virtually fully compliant with the defined requirement, save for a singular spike (less than 0.1% of the

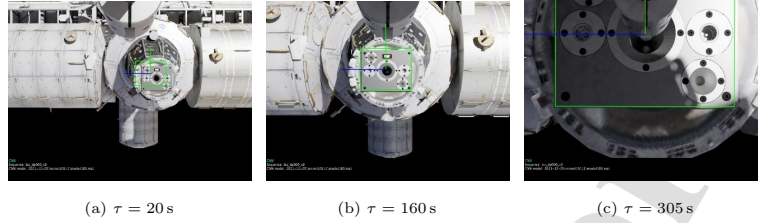


Figure 15: Qualitative pose estimation performance on the `synthetic/01` test sequence.

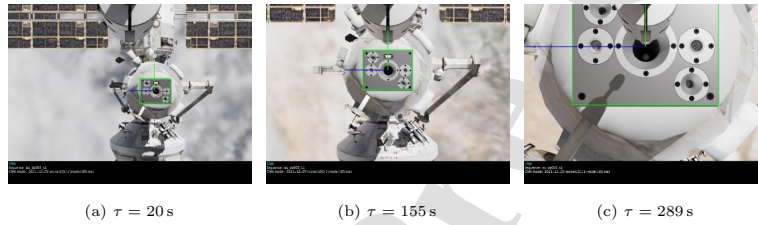


Figure 16: Qualitative pose estimation performance on the `synthetic/08` test sequence.

trajectory). The attitude estimation is again entirely compliant and practically
 550 does not surpass 1 deg in error.

Figure 15 and Figure 16 exhibit some frames from each sequence with the respective qualitative pose estimation fit overlaid. On `synthetic/01` the TV structure surrounding the berthing fixture is quite complex, which from the IP point of view represents a more challenging background than the case
 555 of `synthetic/08`, despite it being an R-bar trajectory which includes Earth. Additionally, the illumination conditions on the former appear to make the berthing fixture harder to distinguish from the ISS structure relative to the latter. Both aspects could provide an explanation to the increased position error seen in the beginning of `synthetic/01`.

5.1.3. Effect of Temporal Modelling

560 The designed OibarNet pipeline uses a CNN front-end to process incoming images and extract features. However, these images are processed individually,

whereby the data as a whole represents a time sequence depicting a docking manoeuvre, which implies that each sample is time correlated. Modifications
565 to CNN architectures have been proposed in the past to account for this correlation and shown to improve the relative pose estimation error for rendezvous. Specifically, deep recurrent convolutional neural networks (DRCNNs) include a recurrent sub-network as the back-end of the pipeline that models the features extracted by the CNN [27].

570 This test investigates the effect of applying a DRCNN to the problem of relative pose estimation for docking. To this end, the trained CNN model was appended with a recurrent model, further trained on the output of the CNN for the same dataset, consisting of bi-directional long-short-term memory (BiLSTM) cells [11]. Contrary to regular long-short-term memories (LSTMs), BiLSTMs
575 run sequence inputs in two directions: one from past to future, and the other from future to past, thus preserving information from both past and future. This feature can be beneficial for RVD/B pose estimation problems since trajectories are continuous, meaning that not only do the previous states influence the present, but states in the future provide context to the preceding ones.

580 The results of the benchmark are illustrated in Figure 17. It can be seen that the addition of a recurrent layer degrades not only the validation performance, but also the training performance; this is witnessed both in terms of position and attitude estimation. Adding more recurrent layers lowers the error on the attitude estimate, but even with three layers this is still higher than that obtained
585 for the CNN alone. Furthermore, the position error is shown not to decrease.

This study represents an interesting result since it is seemingly counter-intuitive and diverges from the findings reported by Rondao et al. [27]. However, whereas the apparent relative motion during an RV is typically smooth and predictable (e.g., SV at a hold point observing the TV tumbling), the docking
590 trajectories modelled within the scope of OIBAR are actually more dynamic and include higher stochasticity due to the random perturbations added during the approach phase. As such, one explanation towards the poor performance of the DRCNN in this case could be the failure in modelling these high-frequency,

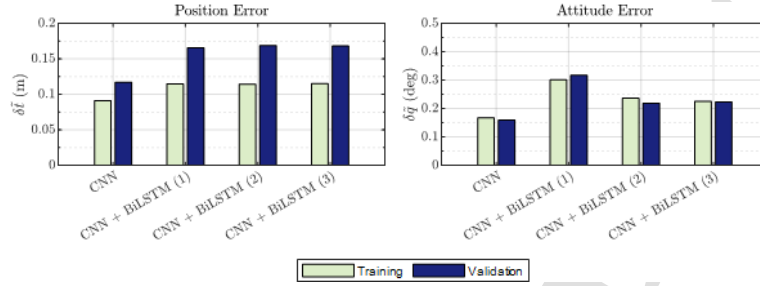


Figure 17: Effect of adding recurrent neural layers on the training and validation performance of the synthetic dataset. The numbers in parenthesis denote the number of recurrent layers.

random changes in motion, in which the CNN indeed an advantage as it is
 595 processing each time-step individually.

Further avenues of research could still be pursued, however. For example, the
 inclusion of attention-based mechanisms remains to be investigated for RVD/B,
 where the network would be capable of self-learning weights to be attributed
 to each time-step in the sequence, thus becoming able to let certain segments
 600 influence the estimate more than others (e.g., placing less attention on the
 immediate perturbations and more on the overall along-track motion).

Due to the attained results, the OibarNet architecture was not altered for
 the integration tests.

5.2. Integration Testing

605 5.2.1. Performance Evaluation on Experimental Dataset

This section is analogous to Subsection 5.1.2 with the difference that the
 selected OibarNet CNN architecture is evaluated and tested on experimental
 data collected in laboratory. As outlined in Section 4, the performance of the
 navigation algorithm is evaluated on the test sequences `experimental/11` and
 610 `experimental/12`.

Figure 18 displays the attained pose estimation errors for both trajecto-
 ries. Table 5 summarises the performance metric statistics. Lastly, Figure
 19 and Figure 20 illustrate qualitative estimation results for a few frames of

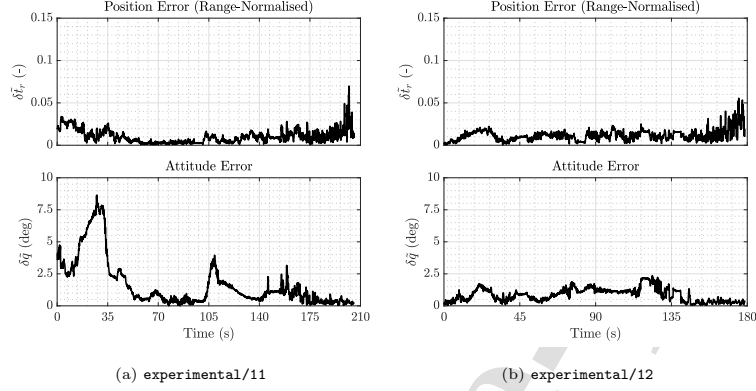


Figure 18: Estimated position and attitude errors over time on the two test sequences of the experimental dataset.

Table 5: Summary performance statistics on the two test sequences of the experimental dataset. Position errors are range-normalised. “Std.” denotes standard deviation.

Sequence	Position Error (%)			Attitude Error (deg)			Requirement Compliance (%)	
	Mean	Median	Std.	Mean	Median	Std.	Position	Attitude
experimental/11	1.02	0.84	0.80	1.65	0.91	1.85	99.71	91.20
experimental/12	1.17	1.08	0.67	0.86	0.87	0.52	99.72	100.00

the `experimental/11` and `experimental/12` sequences, respectively. The relative position estimation error follows a similar trend to the synthetic dataset case: lower during the approach phase and increasing in the final alignment and soft-docking phase. Overall, the curves oscillate more in amplitude for both trajectories; this is a possible by-product of using real-data which can be contaminated with random errors (e.g., sensor noise) and systematic errors (e.g., errors in the motion capture system calibration), which are not seen in the ideal development conditions of synthetic datasets. The reduced number of training samples relative to the synthetic case also affects the solution (i.e., the experimental trajectories are shorter). Nevertheless, the requirement compliance is virtually 100 % for both trajectories.

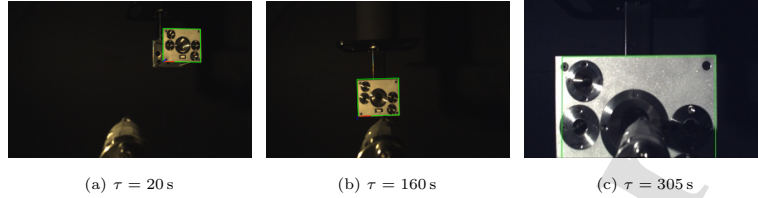


Figure 19: Qualitative pose estimation performance on the `experimental/11` test sequence.

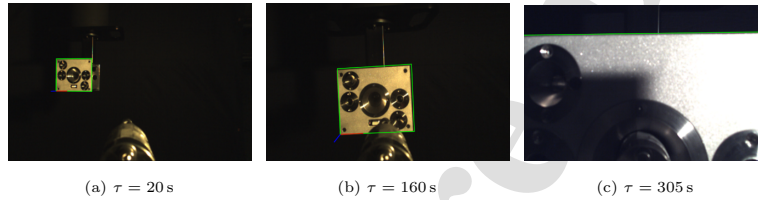


Figure 20: Qualitative pose estimation performance on the `experimental/12` test sequence.

625 The experimental evaluation demonstrates, on average, a higher attitude
 error than the synthetic evaluation case. In particular, for `experimental/11`,
 a spike in the initial 35 seconds of the sequence cause the error to surpass
 7.5 deg which brings down the requirement compliance to 91.2%. This is due
 to the `SV` travelling to a waypoint during the acquisition phase that is quite
 630 distinctive from the others present in the training data, making the berthing
 fixture appear in the top right corner of the `FOV` close to the image edge
 (Figure 27 a). However, the proposed training scheme which includes image
 augmentation prevents the error from diverging, and the estimate begins to
 recover after $\tau = 35$ s, reaching minimum values during and immediately before
 635 the final phase. In `experimental/12`, the attitude estimation error is bounded
 at 2.5 deg.

6. Conclusions

OOS is now becoming increasingly important and represents a significant
 cost saving measure, opening up a new global market. Whereas the latest

640 OOS initiatives and demonstrators have focused on clear near-term commercial opportunities, such as life extension and end-of-life, longer-term OOS segments expected to emerge this decade such as refuelling are set to unravel novel and wider business opportunities, and have the potential to unlock new orbital ecosystems.

645 On this basis, City, University of London have developed OIBAR, a novel AI-based solution for space docking and refuelling applications consisting of the combination of two major components: a vision-based orbital relative navigation algorithm to safely approach and dock to the target vehicle; and an intelligent hardware mechanism achieving the mechanical docking and refuelling operation
650 of the target. The present document reported the development and achievements of the OIBAR project, namely the design procedure of its key features adopted to tackle the problem, the modelling of the mechanism and software architecture, and the validation of the combined solution. Functional testing of the prototype was performed in laboratory using a 7-DOF robotic manipulator to simulate
655 docking/berthing trajectories and a state-of-the-art Optitrack ground truth measurement system to assess the quality of the navigation solution.

A CNN-based direct VBS navigation algorithm was proposed to estimate the relative states between SV and TV to achieve docking. A MATLAB/Simulink simulator was developed to generate synthetic data intended to train and evaluate
660 the solution. A benchmarking campaign was performed to assess the best architecture candidate. The final model reported average errors per trajectory of 1.19% and 0.29 deg for range-normalised position and for attitude, respectively, with accompanying average standard deviations of 1.14% and 0.19 deg. The performance requirements were satisfied for nearly the whole length of the test
665 sequences. The inclusion of BiLSTM-based recurrent layers was analysed but found not to improve the base CNN model.

Lastly, the combined solution was assessed through an integration testing campaign. The navigator was trained and tested on experimental data collected in laboratory using the mechanical docking prototype. Similarly to the synthetic
670 dataset results, these have achieved near-complete compliance with the proposed

accuracy requirements, thus validating the findings. Some differences between the two sets have been observed, however, namely in terms of an overall increase in the mean attitude error which can be attributed to an increased variation in the possible attitude states induced by the waypoint programming on the robotic manipulator. An enlargement of the training dataset poses is expected to further reduce the error.

Acknowledgements

This research has been supported by the UK Robotics and Artificial Intelligence (RAI) hub in Future AI and Robotics for Space (FAIR-SPACE: EP/R026092).

References

- [1] European Space Agency. ESA invites ideas to open up in-orbit servicing market, 2021-04-01. URL https://www.esa.int/Safety_Security/Clean_Space/ESA_invites_ideas_to_open_up_in-orbit_servicing_market.
- [2] Astroscale, Fair-Space, and Catapult Space Applications. Uk in-orbit servicing capability. Technical report, UK Space Agency, Swindon, UK, 2021-05. URL <https://sa.catapult.org.uk/wp-content/uploads/2021/05/Catapult-Astroscale-Fairspace-Platform-for-Growth-report-final-27-05-21.pdf>.
- [3] Patrice Benarroche, Martial Vanhove, and Mauro Augelli. ATV Operations: from Demo Flight to Human Spaceflight Partner. In *SpaceOps 2014 Conference*. American Institute of Aeronautics and Astronautics, 05 2014. doi:10.2514/6.2014-1665.
- [4] Vincenzo Capuano, Kyunam Kim, Alexei Harvard, and Soon-Jo Chung. Monocular-based pose determination of uncooperative space objects. *Acta Astronautica*, 166:493–506, January 2020. doi:10.1016/j.actaastro.2019.09.027.

- [5] Bo Chen, Jiewei Cao, Alvaro Parra, and Tat-Jun Chin. Satellite Pose Estimation with Deep Landmark Regression and Nonlinear Pose Refinement. In *2019 IEEE/CVF International Conference on Computer Vision Workshop (ICCVW)*. IEEE, oct 2019. doi:[10.1109/iccvw.2019.00343](https://doi.org/10.1109/iccvw.2019.00343).
700
- [6] Roberto Cipolla, Yarin Gal, and Alex Kendall. Multi-task Learning Using Uncertainty to Weigh Losses for Scene Geometry and Semantics. In *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*. IEEE, 6 2018. doi:[10.1109/cvpr.2018.00781](https://doi.org/10.1109/cvpr.2018.00781).
- [7] Vicki Cox. Northrop Grumman’s Wholly Owned Subsidiary, SpaceLogistics, Selected by DARPA as Commercial Partner for Robotic Servicing Mission, 2020-03-04. URL <https://news.northropgrumman.com/news/releases/northrop-grummans-wholly-owned-subsidiary-spacelogistics-selected-by-darpa-as-commercial-partner-for-robotic-servicing-mission>.
705
710
- [8] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. ImageNet: A large-scale hierarchical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*. IEEE, 6 2009. doi:[10.1109/cvpr.2009.5206848](https://doi.org/10.1109/cvpr.2009.5206848).
- [9] Wigbert Fehse. *Sensors for rendezvous navigation*, page 218–282. Cambridge Aerospace Series. Cambridge University Press, 2003. doi:[10.1017/CBO9780511543388.008](https://doi.org/10.1017/CBO9780511543388.008).
715
- [10] Martin A. Fischler and Robert C. Bolles. Random sample consensus: a paradigm for model fitting with applications to image analysis and automated cartography. *Communications of the ACM*, 24(6):381–395, June 1981. ISSN 1557-7317. doi:[10.1145/358669.358692](https://doi.org/10.1145/358669.358692).
720
- [11] Alex Graves, Santiago Fernández, and Jürgen Schmidhuber. Bidirectional LSTM Networks for Improved Phoneme Classification and Recognition. In Włodzisław Duch, Janusz Kacprzyk, Erkki Oja, and Sławomir Zadrozny,

- 725 editors, *Artificial Neural Networks: Formal Models and Their Applications*
– *ICANN 2005*, pages 799–804, Berlin, Heidelberg, 2005. Springer Berlin
Heidelberg. ISBN 978-3-540-28756-8. doi:[10.1007/11550907_126](https://doi.org/10.1007/11550907_126).
- [12] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual
learning for image recognition, 2015.
- 730 [13] Lei He, Duarte Ronda, and Nabil Aouf. A Novel Mechanism for Orbital
AI-based Autonomous Refuelling. In *2023 AIAA Guidance, Navigation, and*
Control Conference. American Institute of Aeronautics and Astronautics, 1
2023.
- [14] Geoffrey E. Hinton, Nitish Srivastava, Alex Krizhevsky, Ilya Sutskever,
735 and Ruslan R. Salakhutdinov. Improving neural networks by preventing
co-adaptation of feature detectors, 2012.
- [15] Forrest N. Iandola, Song Han, Matthew W. Moskewicz, Khalid Ashraf,
William J. Dally, and Kurt Keutzer. SqueezeNet: AlexNet-level accuracy
with 50x fewer parameters and 0.5 MB model size, 2016.
- 740 [16] Diederik P. Kingma and Jimmy Ba. Adam: A Method for Stochastic
Optimization, 2014.
- [17] Mate Kisantal, Sumant Sharma, Tae Ha Park, Dario Izzo, Marcus Martens,
and Simone D'Amico. Satellite pose estimation challenge: Dataset, compe-
tition design, and results. *IEEE Transactions on Aerospace and Electronic*
745 *Systems*, 56(5):4083–4098, 10 2020. doi:[10.1109/taes.2020.2989063](https://doi.org/10.1109/taes.2020.2989063).
- [18] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. ImageNet classi-
fication with deep convolutional neural networks. *Communications of the*
ACM, 60(6):84–90, 5 2017. doi:[10.1145/3065386](https://doi.org/10.1145/3065386).
- [19] F. Landis Markley and John L. Crassidis. *Fundamentals of Spacecraft Atti-*
750 *tude Determination and Control*. Springer New York, 2014. doi:[10.1007/978-1-4939-0802-8](https://doi.org/10.1007/978-1-4939-0802-8).

- [20] Laurie Metcalfe and Tara Hillebrandt. Robotic Refuelling Mission: Demonstrating Satellite Refuelling Technology on Board the ISS. In *Proceedings of 12th International Symposium on Artificial Intelligence, Robotics and Automation in Space (i-SAIRAS)*, Montréal, Canada, 2014. European Space Agency.
- [21] Lorenzo Pasqualetto Cassinis, Alessandra Menicucci, Eberhard Gill, Ingo Ahrens, and Manuel Sanchez-Gestido. On-ground validation of a cnn-based monocular pose estimation system for uncooperative spacecraft: Bridging domain shift in rendezvous scenarios. *Acta Astronautica*, 196:123–138, July 2022. ISSN 0094-5765. doi:10.1016/j.actaastro.2022.04.002.
- [22] 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.
- [23] Joseph Redmon and Ali Farhadi. YOLO9000: Better, Faster, Stronger. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, 7 2017. doi:10.1109/cvpr.2017.690.
- [24] Joseph Redmon and Ali Farhadi. YOLOv3: An Incremental Improvement, 2018.
- [25] Duarte Rondao. *Multimodal Navigation for Accurate Rendezvous Missions*. PhD thesis, Cranfield University, 2021. Available from: <https://ethos.bl.uk/OrderDetails.do?uin=uk.bl.ethos.883372>.
- [26] Duarte Rondao, Nabil Aouf, Mark A. Richardson, and Vincent Dubanchet. Robust On-Manifold Optimization for Uncooperative Space Relative Navigation with a Single Camera. *Journal of Guidance, Control, and Dynamics*, 44(6):1157–1182, 6 2021. doi:10.2514/1.g004794.
- [27] Duarte Rondao, Nabil Aouf, and Mark A. Richardson. ChiNet: Deep Recurrent Convolutional Learning for Multimodal Spacecraft Pose Estimation. *IEEE Transactions on Aerospace and Electronic Systems*, pages 1–13, 2022.

- 780 doi:10.1109/taes.2022.3193085. URL <https://doi.org/10.1109/2Ftaes.2022.3193085>.
- [28] Leslie N. Smith. Cyclical Learning Rates for Training Neural Networks. In *2017 IEEE Winter Conference on Applications of Computer Vision (WACV)*. IEEE, 3 2017. doi:10.1109/wacv.2017.58.
- 785 [29] Jianing Song, Duarte Rondao, and Nabil Aouf. Deep learning-based spacecraft relative navigation methods: A survey. *Acta Astronautica*, 191:22–40, 2 2022. doi:10.1016/j.actaastro.2021.10.025.
- [30] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew
790 Rabinovich. Going deeper with convolutions. In *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, jun 2015. doi:10.1109/cvpr.2015.7298594.
- [31] Richard Szeliski. *Structure from motion and SLAM*, pages 543–594. Springer International Publishing, Cham, 2022. doi:10.1007/978-3-030-34372-9_11.
- 795 [32] Katsuyoshi Tsujita and Asumi Nishimura. Development of an autonomous soft docking system of small spacecraft using visual guidance. In *2022 IEEE/SICE International Symposium on System Integration (SII)*. IEEE, January 2022. doi:10.1109/sii52469.2022.9708765.
- [33] Andrea Valmorbida, Mattia Mazzucato, and Marco Pertile. Calibration
800 procedures of a vision-based system for relative motion estimation between satellites flying in proximity. *Measurement*, 151:107161, 2 2020. doi:10.1016/j.measurement.2019.107161.
- [34] Claudio Vela, Giancarmine Fasano, and Roberto Opromolla. Pose determination of passively cooperative spacecraft in close proximity using a monocular
805 camera and aruco markers. *Acta Astronautica*, 201:22–38, December 2022. ISSN 0094-5765. doi:10.1016/j.actaastro.2022.08.024.

- [35] Renato Volpe, Giovanni B. Palmerini, and Marco Sabatini. A passive camera based determination of a non-cooperative and unknown satellite's pose and shape. *Acta Astronautica*, 151:805–817, October 2018. doi:[10.1016/j.actaastro.2018.06.061](https://doi.org/10.1016/j.actaastro.2018.06.061).
810
- [36] Bong Wie, Vaios Lappas, and Jesús Gil-Fernández. *Attitude and Orbit Control Systems*, pages 323–369. Springer Berlin Heidelberg, Berlin, Heidelberg, 2014. ISBN 978-3-642-41101-4. doi:[10.1007/978-3-642-41101-4_12](https://doi.org/10.1007/978-3-642-41101-4_12).
- [37] Yi Zhou, Connelly Barnes, Jingwan Lu, Jimei Yang, and Hao Li. On the continuity of rotation representations in neural networks, 2018.
815

- A novel artificial intelligence-based solution for spacecraft docking and refuelling applications is presented.
- The viability of deep neural networks using vision-only inputs for pose estimation of a docking structure is demonstrated for the first time.
- Multiple convolutional neural network backbones are benchmarked on a photorealistic dataset of a refuelling scenario with the International Space Station.
- A ResNet-101 architecture reports average errors per trajectory of 1.19% for range-normalised position and 0.29 degrees for attitude.
- The end-effector and berthing fixture prototypes are manufactured and validated in the laboratory with a 7 degree-of-freedom robotic manipulator.

Declaration of interests

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

The authors declare the following financial interests/personal relationships which may be considered as potential competing interests:

Duarte Rondao reports financial support was provided by UK Research and Innovation. If there are other authors, they declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.