



SCHOOL OF COMPUTER AND  
COMMUNICATION SCIENCES

ÉCOLE POLYTECHNIQUE FÉDÉRALE DE LAUSANNE

**Computer Vision Laboratory**  
**Unseen Spacecraft Pose Estimation**

Baseline solution implementation of a deep learning model  
for unseen spacecraft pose estimation

Bachelor's Thesis in Computer Science

Author:	JÉRÉMY CHAVEROT
Supervisor:	Dr. MATHIEU SALZMANN
Advisors:	Dr. ANDREW PRICE, PhD. CHEN ZHAO
Semester:	Fall 2023

I hereby confirm that I am the sole author of the written work here enclosed and that I have compiled it in my own words. Parts excepted are corrections of form and content by the advisors.

Lausanne, Switzerland, 26.01.24

JÉRÉMY CHAVEROT

## **Acknowledgments**

Before delving into the topic, I'd like to express some acknowledgments. First and foremost I must thank my advisors Andrew Price and Chen Zhao for accepting my request to take part in a semester project under their guidance. I am grateful to have been able to practice my skills with them, and can only hope that the feeling is mutual. Moreover I would also like to thoroughly thank my friends and family for supporting me in my academic journey, despite a rather unstable start in my studies.

## **Abstract**

Nowadays, the significant number of objects in Low Earth Orbit (LEO) is becoming a matter of concern, with the potential to create even more space debris. Therefore, it's becoming crucial to develop technologies to tackle this pollution. One can envision a clean-up mission where devices are equipped with quantized deep learning models to estimate the position of these objects for retrieval. Training these models is challenging, as the real image database is quite sparse. To overcome this, high-quality, realistic synthetic images of spacecraft are used, with the ultimate goal of generalizing to unseen objects.

# Contents

<b>Acknowledgments</b>	<b>ii</b>
<b>Abstract</b>	<b>iii</b>
<b>1. Introduction</b>	<b>1</b>
1.1. Problem Formulation . . . . .	1
1.2. Applications . . . . .	1
1.3. The Work Environment: EPFL Scitas Izar . . . . .	2
<b>2. Gen6D: Formal Description</b>	<b>3</b>
2.1. Overview of the Network . . . . .	3
2.2. Detection . . . . .	4
2.3. Viewpoint Selection . . . . .	4
2.4. Pose Refinement . . . . .	5
<b>3. Implementation of the model</b>	<b>6</b>
3.1. Data Loader . . . . .	6
3.2. From Quaternions to Rotation Matrices . . . . .	7
3.3. Inverted Masks . . . . .	9
3.4. Resized Query Images . . . . .	10
<b>4. Experimental Results and Analysis</b>	<b>12</b>
4.1. Evaluation Metrics . . . . .	12
4.2. Vizualisation and Quantitative Evaluation . . . . .	12
<b>5. Ways of improvements</b>	<b>19</b>
5.1. More diverse Reference Images . . . . .	19
5.2. Specialized Spacecraft Training Set . . . . .	19
5.3. Improve Object Detection and Viewpoint Selection Algorithms . . . . .	19
<b>6. Conclusion</b>	<b>21</b>
<b>Abbreviations</b>	<b>22</b>
<b>A. Python Scripts</b>	<b>23</b>
<b>B. Scitas Izar Setup Tutorial</b>	<b>31</b>
<b>Bibliography</b>	<b>33</b>

# 1. Introduction

This project falls within an Unseen 6 Degrees of Freedom (DoF) competition, organized in collaboration with the European Space Agency (ESA) Advanced Concept Team. Essentially, we are dealing with space objects that are unfamiliar to us, and our objective is to accurately predict their 6DoF poses. The action of the Computer Vision Laboratory (CVLab) team is twofold: firstly, we are tasked with creating a challenging dataset featuring multi-object, unseen, and occluded spacecraft scenarios. This involves ensuring a high degree of rendering realism. Secondly, we are focused on developing a baseline solution, which entails implementing a chosen pose estimation model and conducting thorough training and testing on our dataset. My role this semester was primarily concentrated on the latter aspect, specifically on a track that incorporated 3D target models.

## 1.1. Problem Formulation

To put it simply, our task involves analyzing images of space objects with the aim of accurately determining the relative 3D translation and 3D rotation of the spacecraft in relation to the camera's position.

For the training and testing of the selected deep learning architecture, we utilize synthetic images provided by the SPACECRAFT dataset team. This dataset encompasses four distinct models: the Hubble Space Telescope, the James Webb Space Telescope, the Cosmos Link, and the Rocket Body. These images may feature an Earth-rendered background or be without it. In addition to the images, our dataset includes 3D models, masks, segmented images, and camera settings. Furthermore, the ground truth poses of the objects are available in a format that combines quaternions with a translation vector.

## 1.2. Applications

Beyond the scope of the ESA competition, the project has numerous practical applications in the aerospace industry.

**Remediation of Space Debris in LEO** The developed pose estimation models can significantly enhance the identification and tracking of space debris, enabling precise navigation for cleanup missions.

**Space Monitoring** The technology can also be utilized for the constant surveillance and cataloging of artificial and natural objects in space, improving awareness of the space environment and collision avoidance systems.

**Planetary Defense Technology** Finally we could think about the accurate pose estimation capabilities serving as a critical component in planetary defense missions by ensuring precise targeting of potentially hazardous small celestial bodies, for instance the NASA's Double Asteroid Redirection Test (DART) mission in 2021, which was the first-ever asteroid deflection mission through kinetic impact.

## 1.3. The Work Environment: EPFL Scitas Izar

We needed to establish a suitable environment for executing the code: EPFL Scitas Izar servers. They are ideally configured for our task: equipped with two NVIDIA V100 PCIe 32 GB GPUs, the most advanced data center GPUs ever built to accelerate Artificial Intelligence (AI), High Performance Computing (HPC), data science and graphics (very expensive as well). Additionally, the server enables users to execute the code from any location, eliminating the need for a sufficiently powerful hardware configuration.

However, this stage proved to be more time-consuming than expected. It involved software engineering and various technical challenges, including setting up the virtual environment, installing the necessary dependencies in a manner that avoids conflicts with the modules already installed on the server, composing the bash execution script, and, fundamentally, learning the correct way to utilize the server. Special recognition goes to Emily Bourne from the EPFL HPC team for her essential help and support.

For those who wish to run the later described model, instructions for setting up the environment are provided in Appendix B of this report.

## 2. Gen6D: Formal Description

### 2.1. Overview of the Network

After conducting a scientific literature search, we decided on an interesting model known as Gen6D: Generalizable Model-Free 6-DoF Object Pose Estimation from RGB Images [1]. This architecture is fed with reference and query images, and operates in a relatively straightforward manner, comprising three stages shown in Figure 2.1. Initially, there's the detector, which, as its name suggests, detects the object in the query image. Next is the viewpoint selection, which selects reference images that have the closest viewing angle to the detected object in the query image by using a similarity score. Finally, there's the pose refiner, which enhances the accuracy of the estimated pose through a 3D volume-based process. We have compiled a table summarizing the advantages and disadvantages for a clearer understanding:

Table 2.1.: Summary of Gen6D

Pros	Cons
Generalizability	Limited by Reference Image Quality
Model-Free	Everyday Life Objects Training Data
Simple Input Requirements	Difficulty with Symmetric Objects
Robustness to Background Clutter	Dependence on Initial Detection and Selection
Effective in Diverse Environments	Potential Challenges with Severe Occlusions
Competitive Performance	Computationally Intensive

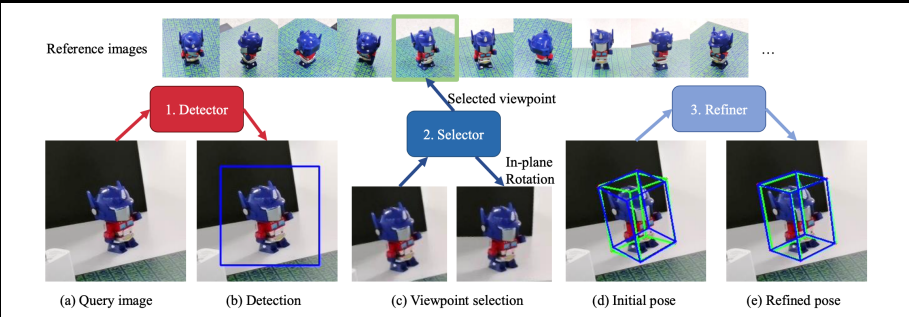


Figure 2.1: Overview of Gen6D, image sourced from Gen6D's research paper.



## 2.2. Detection

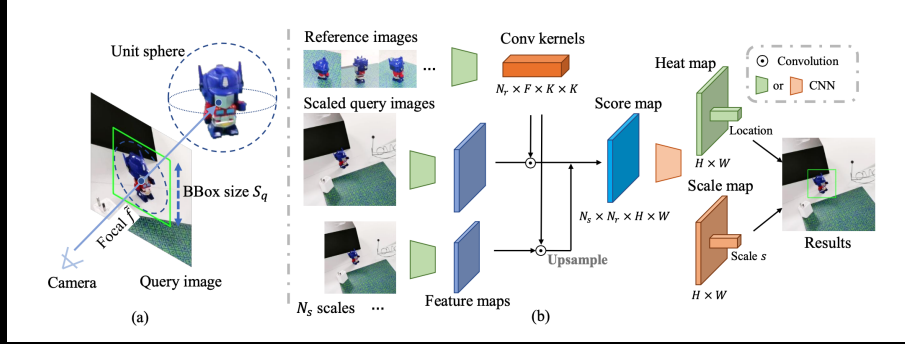


Figure 2.2: Architecture of the detector, image sourced from Gen6D's research paper.

In the detection phase, the initial step involves resizing both reference and query images. The reference images are resized to a dimension of  $128 \times 128$ , with the object centered. The query images are resized at various predefined scales. Subsequently, all the resized images are processed through a Visual Geometry Group Network (VGGNet), which essentially functions as a deep convolutional network for extracting feature maps.

Following this, the feature maps from the reference images are utilized as convolutional kernels to convolve with one of the query images, resulting in the generation of a score map. Leveraging the multi-scale score map, we perform regression to obtain a heat map and a scale map. The final 2D detection is acquired by identifying the maximum value in the heat map, which provides the object's center coordinates, while the scale  $s$  is obtained from the same location in the scale map, by this way determining the size of the 2D bounding box.

## 2.3. Viewpoint Selection

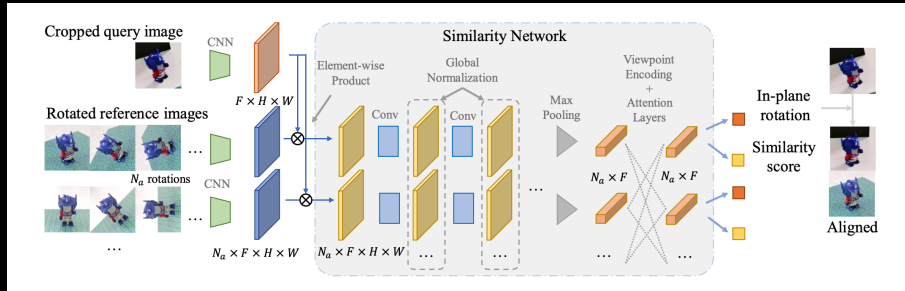


Figure 2.3: Architecture of the viewpoint selector, image sourced from Gen6D's research paper.

During the viewpoint selection phase, we begin by applying predefined rotations to the reference images to accommodate in-plane rotation variations. After that,

we once again extract feature maps from both the query image, which has been cropped based on the detection results, and the rotated reference images.

Next we calculate the element-wise product of the query image’s features with those of each rotated reference image, yielding a correlation score map. These scores are then inputted into a similarity network, which produces a similarity score and the relative in-plane rotation.

It’s worth highlighting that within the similarity network, the authors have implemented global normalization and incorporated a transformer to facilitate information sharing among the reference images.

## 2.4. Pose Refinement

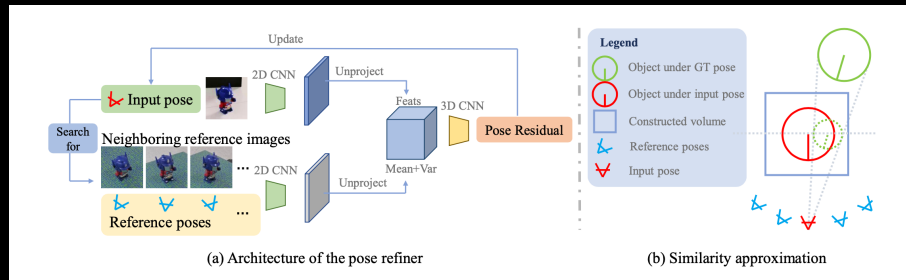


Figure 2.4: Architecture of the pose refiner, image sourced from Gen6D’s research paper.

From the two previous stages, we have an initial coarse object pose. In this refinement process, we one last time extract feature maps from the selected reference images using a 2D Convolutional Neural Network (CNN). Afterwards, these feature maps are unprojected into the 3D volume, and we compute their mean and variance, used as features associated with the volume vertices.

Similarly, for the query image, we perform the same process but utilize the input pose. The unprojected query image features are then concatenated with the previously obtained mean and variance features.

At the end, we have a 3D CNN that operates on the concatenated feature set of the volume, so as to predict a pose residual that updates the input pose. We iterate through the loop three times, which has been determined empirically as the optimal number of steps, in order to enhance the accuracy of the estimated pose.

## 3. Implementation of the model

### 3.1. Data Loader

Now that we have a better understanding of how Gen6D operates, we're interested in testing it with the synthetic images from the SPACECRAFT dataset. To accomplish this, we need to redefine the Data Loader to supply the model with the necessary data in the required format. In the file `database.py`, located in the `Gen6D/dataset` directory, there is an Abstract Base Class (ABC) named `BaseDatabase()` containing all the essential methods, as shown in the code below:

```
1 class BaseDatabase(abc.ABC):
2     def __init__(self, database_name):
3         self.database_name = database_name
4
5     @abc.abstractmethod
6     def get_image(self, img_id):
7         pass
8
9     @abc.abstractmethod
10    def get_K(self, img_id):
11        pass
12
13    @abc.abstractmethod
14    def get_pose(self, img_id):
15        pass
16
17    @abc.abstractmethod
18    def get_img_ids(self):
19        pass
20
21    def get_mask(self, img_id):
22        # dummy mask
23        img = self.get_image(img_id)
24        h, w = img.shape[:2]
25        return np.ones([h,w], np.bool_)
```

Listing 3.1: Python code of ABC `BaseDatabase()`, from file `database.py`.

Our goal, therefore, is to write a new class that inherits from the ABC `BaseDatabase()`, and implement each abstract method in the same manner as it was done for the LINEMOD and GENMOP datasets, for example. We will also adapt the code in various places, especially in the `eval.py` file, to ensure compatibility with the SPACECRAFT dataset. The Python code for the new data loading class `SpaceCraftCVLabDatabase()` can be found in Listing A.1.

## 3.2. From Quaternions to Rotation Matrices

In Gen6D, the model represents the ground truth and estimated poses using the format  $\mathbf{P} = (\mathbf{R}, \mathbf{t})$ . Here,  $\mathbf{R}$  denotes the rotation matrix and  $\mathbf{t}$  is the translation vector. The thing is, in the SPACECRAFT dataset, the poses have the format  $\mathbf{P} = (\underline{\mathbf{q}}, \mathbf{t})$ , where  $\underline{\mathbf{q}}$  is a *quaternion*. We use quaternions for three-dimensional rotation calculations because they offer several benefits over rotation matrices. Notably, quaternions are more compact, requiring only four elements to be stored compared to nine for a matrix. Additionally, they are more efficient when composing rotations thanks to their algebraic properties.

Despite the benefits of quaternions mentioned earlier, other datasets frequently represent poses using a combination of rotation matrices and translation vectors. Therefore, to align with Gen6D's pose format, this section will focus on converting quaternions into rotation matrices [2].

Quaternions were first introduced by the Irish mathematician W. R. Hamilton in 1843 as an extension of the complex numbers [3]. In the first place, we provide the definition of a *quaternion*  $\underline{\mathbf{q}}$ : it is the sum of a scalar  $q_0$  and a vector  $\mathbf{q} = (q_1, q_2, q_3)$ , that is,

$$\underline{\mathbf{q}} \stackrel{\text{def.}}{=} q_0 + \mathbf{q} = q_0 + q_1\mathbf{i} + q_2\mathbf{j} + q_3\mathbf{k}.$$

In the above,  $\mathbf{i}$ ,  $\mathbf{j}$  and  $\mathbf{k}$  denote the three unit vectors of the canonical basis for the set of all ordered triples of real numbers  $\mathbb{R}^3$ . The set of quaternions is denoted by the 4-space  $\mathbb{H}$ .

The quaternion addition is component-wise. Regarding the product of two quaternions, it is essential to first outline the foundational rule established by Hamilton:

$$\mathbf{i}^2 = \mathbf{j}^2 = \mathbf{k}^2 = \mathbf{i}\mathbf{j}\mathbf{k} = -1.$$

We derive the following multiplication table:

$\times$	1	$\mathbf{i}$	$\mathbf{j}$	$\mathbf{k}$
1	1	$\mathbf{i}$	$\mathbf{j}$	$\mathbf{k}$
$\mathbf{i}$	$\mathbf{i}$	-1	$\mathbf{k}$	$-\mathbf{j}$
$\mathbf{j}$	$\mathbf{j}$	$-\mathbf{k}$	-1	$\mathbf{i}$
$\mathbf{k}$	$\mathbf{k}$	$\mathbf{j}$	$-\mathbf{i}$	-1

Let  $(\underline{\mathbf{p}}, \underline{\mathbf{q}}) \in \mathbb{H}^2$ , we are now able to present the multiplication of  $\underline{\mathbf{p}}$  and  $\underline{\mathbf{q}}$ :

$$\underline{\mathbf{p}}\underline{\mathbf{q}} = \underbrace{p_0q_0 - \mathbf{p} \cdot \mathbf{q}}_{\text{scalar part}} + \underbrace{p_0\mathbf{q} + q_0\mathbf{p} + \mathbf{p} \times \mathbf{q}}_{\text{vector part}}.$$

In the preceding, we recall that the binary operation  $\mathbb{R}^3 \times \mathbb{R}^3 \rightarrow \mathbb{R} : (\mathbf{p}, \mathbf{q}) \mapsto \mathbf{p} \cdot \mathbf{q} \stackrel{\text{def.}}{=} \mathbf{p}^\top \mathbf{q}$  is the *dot product*, and the other

$$\mathbb{R}^3 \times \mathbb{R}^3 \rightarrow \mathbb{R}^3 : (\mathbf{p}, \mathbf{q}) \mapsto \mathbf{p} \times \mathbf{q} \stackrel{\text{def.}}{=} \begin{vmatrix} \mathbf{i} & \mathbf{j} & \mathbf{k} \\ p_1 & p_2 & p_3 \\ q_1 & q_2 & q_3 \end{vmatrix} = [\mathbf{p}]_{\times} \mathbf{q}$$

is the *cross product*, and the operator  $[\mathbf{p}]_{\times}$  yields the transformation matrix that when multiplied from the right with a vector  $\mathbf{q}$  gives  $\mathbf{p} \times \mathbf{q}$ .

Several additional definitions are essential before we can begin the conversion of quaternions to rotation matrices.

Let  $\underline{\mathbf{q}} = q_0 + \mathbf{q}$  be a quaternion. The *complex conjugate* of  $\underline{\mathbf{q}}$ , denoted by  $\underline{\mathbf{q}}^*$ , is given by the map

$$\mathbb{H} \rightarrow \mathbb{H} : \underline{\mathbf{q}} \mapsto \underline{\mathbf{q}}^* \stackrel{\text{def.}}{=} q_0 - \mathbf{q}.$$

The *norm* of a quaternion  $\underline{\mathbf{q}}$ , denoted  $|\underline{\mathbf{q}}|$ , is the distance obtained from the map

$$\mathbb{H} \rightarrow \mathbb{R}^+ : \underline{\mathbf{q}} \mapsto |\underline{\mathbf{q}}| \stackrel{\text{def.}}{=} \sqrt{\underline{\mathbf{q}}^* \underline{\mathbf{q}}}.$$

Note that a quaternion whose norm is 1 is referred to as a *unit quaternion*. The *reciprocal* of a quaternion is defined as the map

$$\mathbb{H}^* \rightarrow \mathbb{H}^* : \underline{\mathbf{q}} \mapsto \underline{\mathbf{q}}^{-1} \stackrel{\text{def.}}{=} \frac{\underline{\mathbf{q}}^*}{|\underline{\mathbf{q}}|^2}.$$

We observe that if  $\underline{\mathbf{q}}$  is a unit quaternion, we simply have  $\underline{\mathbf{q}}^{-1} = \underline{\mathbf{q}}^*$ . Furthermore, the subsequent proposition is accepted: if  $\underline{\mathbf{q}}$  is a unit quaternion, there exists a unique  $\theta \in [0, 2\pi]$  such that

$$\underline{\mathbf{q}} = q_0 + \mathbf{q} = \cos \frac{\theta}{2} + \mathbf{u} \sin \frac{\theta}{2},$$

where the unit vector  $\mathbf{u}$  is defined as  $\mathbf{u} \stackrel{\text{def.}}{=} \frac{\mathbf{q}}{|\mathbf{q}|}$ .

**Quaternion Rotation Operator** Let  $\underline{\mathbf{q}} \in \mathbb{H}$  be a unit quaternion, and let  $\mathbf{v} \in \mathbb{R}^3$  be a vector. The action of the  $L_{\underline{\mathbf{q}}}$  function

$$\mathbb{R}^3 \rightarrow \mathbb{R}^3 : \mathbf{v} \mapsto L_{\underline{\mathbf{q}}}(\mathbf{v}) \stackrel{\text{def.}}{=} \underline{\mathbf{q}} \mathbf{v} \underline{\mathbf{q}}^*$$

on  $\mathbf{v}$  is equivalent to a rotation of the vector through an angle  $\theta$  about the axis of rotation  $\mathbf{u}$ .

At last, we can proceed with our conversion task: our aim is to find a  $3 \times 3$  rotation matrix  $\mathbf{R}$ , such that

$$\begin{cases} L_{\mathbf{R}}(\mathbf{v}) \stackrel{\text{def.}}{=} \mathbf{R} \mathbf{v} \\ L_{\mathbf{R}}(\mathbf{v}) = L_{\underline{\mathbf{q}}}(\mathbf{v}), \end{cases}$$

which means we wish to obtain an expression for  $\mathbf{R}$  by manipulating  $L_{\underline{\mathbf{q}}}$ , utilizing principles of linear algebra and vector calculus. We consider the vector  $\mathbf{v}$  as a quaternion with a zero scalar component:

$$\begin{aligned}
 L_{\mathbf{q}}(\mathbf{v}) &= \underline{\mathbf{q}}\mathbf{v}\underline{\mathbf{q}}^* \\
 &= (\mathbf{q}_0 + \mathbf{q})(0 + \mathbf{v})(\mathbf{q}_0 - \mathbf{q}) \\
 &= \underbrace{(-\mathbf{q} \cdot \mathbf{v})}_{\text{scalar part}} + \underbrace{(\mathbf{q}_0\mathbf{v} + \mathbf{q} \times \mathbf{v})}_{\text{vector part}} (\mathbf{q}_0 - \mathbf{q}) \\
 &= \mathbf{q}_0(-\mathbf{q} \cdot \mathbf{v}) - (\mathbf{q}_0\mathbf{v} + \mathbf{q} \times \mathbf{v}) \cdot (-\mathbf{q}) \\
 &\quad + (-\mathbf{q} \cdot \mathbf{v})(-\mathbf{q}) + \mathbf{q}_0(\mathbf{q}_0\mathbf{v} + \mathbf{q} \times \mathbf{v}) \\
 &\quad + (\mathbf{q}_0\mathbf{v} + \mathbf{q} \times \mathbf{v}) \times (-\mathbf{q}) \\
 &= \underbrace{-\mathbf{q}_0(\mathbf{q} \cdot \mathbf{v}) + \mathbf{q}_0(\mathbf{q} \cdot \mathbf{v})}_{\text{scalar part}} \\
 &\quad + \underbrace{\mathbf{q}(\mathbf{q} \cdot \mathbf{v}) + \mathbf{q}_0^2\mathbf{v} + \mathbf{q}_0(\mathbf{q} \times \mathbf{v}) + \mathbf{q} \times (\mathbf{q}_0\mathbf{v} + \mathbf{q} \times \mathbf{v})}_{\text{vector part}} \\
 &= \mathbf{q}(\mathbf{q}^T\mathbf{v}) + \mathbf{q}_0^2\mathbf{v} + \mathbf{q}_0(\mathbf{q} \times \mathbf{v}) + \mathbf{q} \times (\mathbf{q}_0\mathbf{v} + \mathbf{q} \times \mathbf{v}) \\
 &= (\mathbf{q} \otimes \mathbf{q} + \mathbf{q}_0^2 \mathbf{I}_{3 \times 3} + 2\mathbf{q}_0 [\mathbf{q}]_{\times} + [\mathbf{q}]_{\times}^2) \mathbf{v}
 \end{aligned}$$

In the above,  $\otimes$  stands for the *outer product* and  $\mathbf{I}_{3 \times 3}$  is the *identity matrix*.

Since  $L_{\mathbf{R}}(\mathbf{v}) = L_{\underline{\mathbf{q}}}(\mathbf{v})$ , we can identify  $\mathbf{R}$  as  $(\mathbf{q} \otimes \mathbf{q} + \mathbf{q}_0^2 \mathbf{I}_{3 \times 3} + 2\mathbf{q}_0 [\mathbf{q}]_{\times} + [\mathbf{q}]_{\times}^2)$ .

We develop and simplify  $\mathbf{R}$  to find its final expression:

$$\begin{aligned}
 \mathbf{R} &= \mathbf{q} \otimes \mathbf{q} + \mathbf{q}_0^2 \mathbf{I}_{3 \times 3} + 2\mathbf{q}_0 [\mathbf{q}]_{\times} + [\mathbf{q}]_{\times}^2 \\
 &= \begin{bmatrix} \mathbf{q}_1^2 & \mathbf{q}_1\mathbf{q}_2 & \mathbf{q}_1\mathbf{q}_3 \\ \mathbf{q}_2\mathbf{q}_1 & \mathbf{q}_2^2 & \mathbf{q}_2\mathbf{q}_3 \\ \mathbf{q}_3\mathbf{q}_1 & \mathbf{q}_3\mathbf{q}_2 & \mathbf{q}_3^2 \end{bmatrix} + \mathbf{q}_0^2 \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} + 2\mathbf{q}_0 \begin{bmatrix} 0 & -\mathbf{q}_3 & \mathbf{q}_2 \\ \mathbf{q}_3 & 0 & -\mathbf{q}_1 \\ -\mathbf{q}_2 & \mathbf{q}_1 & 0 \end{bmatrix} \\
 &\quad + \begin{bmatrix} 0 & -\mathbf{q}_3 & \mathbf{q}_2 \\ \mathbf{q}_3 & 0 & -\mathbf{q}_1 \\ -\mathbf{q}_2 & \mathbf{q}_1 & 0 \end{bmatrix} \begin{bmatrix} 0 & -\mathbf{q}_3 & \mathbf{q}_2 \\ \mathbf{q}_3 & 0 & -\mathbf{q}_1 \\ -\mathbf{q}_2 & \mathbf{q}_1 & 0 \end{bmatrix} \\
 &= 2 \begin{bmatrix} (\mathbf{q}_0^2 + \mathbf{q}_1^2) - \frac{1}{2} & \mathbf{q}_1\mathbf{q}_2 - \mathbf{q}_0\mathbf{q}_3 & \mathbf{q}_1\mathbf{q}_3 + \mathbf{q}_0\mathbf{q}_2 \\ \mathbf{q}_1\mathbf{q}_2 + \mathbf{q}_0\mathbf{q}_3 & (\mathbf{q}_0^2 + \mathbf{q}_2^2) - \frac{1}{2} & \mathbf{q}_2\mathbf{q}_3 - \mathbf{q}_0\mathbf{q}_1 \\ \mathbf{q}_1\mathbf{q}_3 - \mathbf{q}_0\mathbf{q}_2 & \mathbf{q}_2\mathbf{q}_3 + \mathbf{q}_0\mathbf{q}_1 & (\mathbf{q}_0^2 + \mathbf{q}_3^2) - \frac{1}{2} \end{bmatrix}
 \end{aligned}$$

This completes our derivation of the rotation matrix  $\mathbf{R}$ . The Python code can be found in Listing A.2.

### 3.3. Inverted Masks

During the initial testing of Gen6D on the SPACECRAFT dataset, we noticed that the estimated poses were particularly inaccurate. Upon investigation, we discovered that in the reference images folder, the masks were incorrectly inverted: typically, one expects a white object on a black background, but it was the opposite (see Figure 3.1 below).

After discussing this, we realized that this was actually an error in the current version of the dataset. Therefore, we needed to correct the masks we were using through a Python script utilizing the OpenCV library. The Python script can be found in Listing A.3.

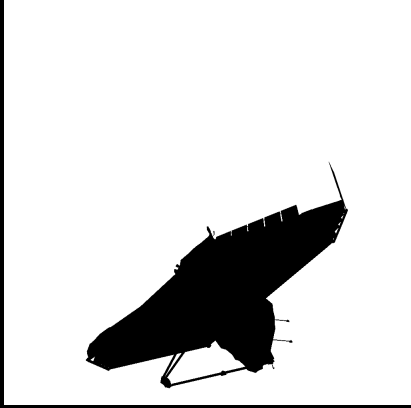


Figure 3.1: James Webb Space Telescope, mask before the correction.

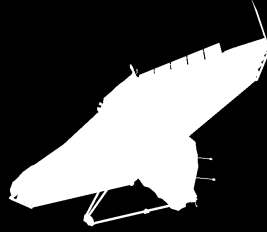


Figure 3.2: James Webb Space Telescope, mask after the correction using a script.

## 3.4. Resized Query Images

During the initial testing of Gen6D, it became apparent that the deep learning model struggled when the object in the query image was either very large or very small (see Figure 3.3 and Figure 3.4 below). Consultation of the GitHub Issues section of the Gen6D repository revealed that this is a recurring problem: since the reference images are resized to  $128 \times 128$ , the model's detector is trained to recognize an object in query images that ranges from half to double the size of the reference images, i.e., from  $64 \times 64$  to  $256 \times 256$ .

Therefore, to improve the detector's results, resizing the query images is necessary. This is accomplished using a Python script that can be found in Listing A.4. It's important to note that this adjustment also requires modifying the camera intrinsic matrix  $\mathcal{K}$  in the data loader, which converts points from the camera coordinate system to the pixel coordinate system.

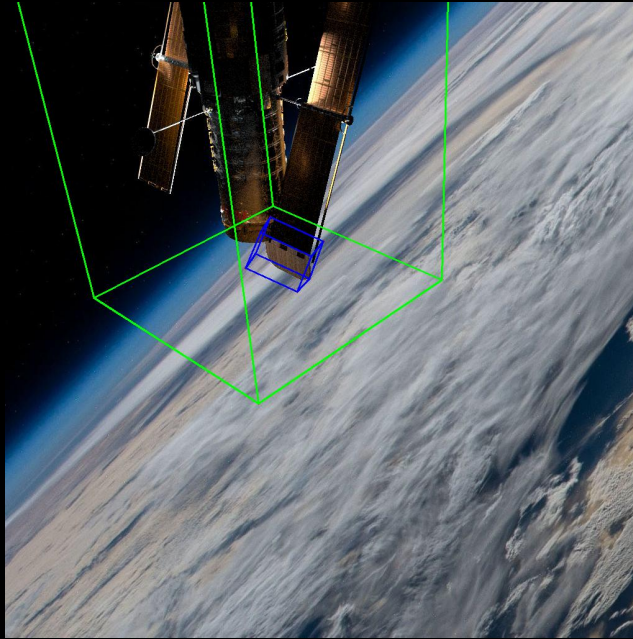


Figure 3.3: Hubble Space Telescope with earth rendered background, 1024x1024 query image number 1, poor pose estimation (blue 3D bounding box) when the object inside the query image is too big.

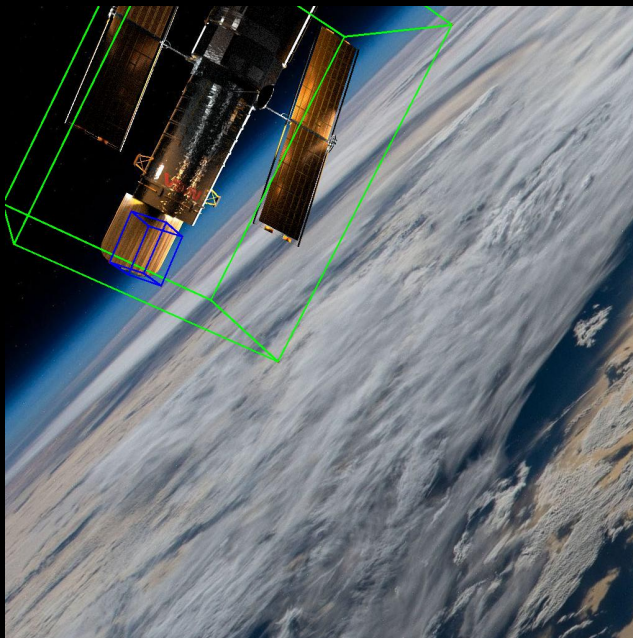


Figure 3.4: Hubble Space Telescope with earth rendered background, 1024x1024 query image number 2, poor pose estimation (blue 3D bounding box) when the object inside the query image is too big.



## 4. Experimental Results and Analysis

### 4.1. Evaluation Metrics

To appreciate the quality of the estimations, the most widely used pose error functions are the Average Distance of Model Points (ADD) and the Average Closest Point Distance (ADD-S) metrics, both introduced by Hinterstoisser et al. [4]. For an object model  $\mathcal{M}$ , we compute the average distance to the corresponding model point. Therefore the error of an estimated pose  $\hat{\mathbf{P}} = (\hat{\mathbf{R}}, \hat{\mathbf{T}})$  w.r.t. the ground truth pose  $\bar{\mathbf{P}} = (\bar{\mathbf{R}}, \bar{\mathbf{T}})$  is calculated as follows:

$$^1e_{\text{ADD}}(\hat{\mathbf{P}}, \bar{\mathbf{P}}, \mathcal{M}) \stackrel{\text{def.}}{=} \text{avg}_{\mathbf{x} \in \mathcal{M}} \left\| \bar{\mathbf{P}}\mathbf{x}^* - \hat{\mathbf{P}}\mathbf{x}^* \right\|_2 \quad (4.1)$$

$$= \text{avg}_{\mathbf{x} \in \mathcal{M}} \left\| (\bar{\mathbf{R}}\mathbf{x} + \bar{\mathbf{T}}) - (\hat{\mathbf{R}}\mathbf{x} + \hat{\mathbf{T}}) \right\|_2 \quad (4.2)$$

When the model  $\mathcal{M}$  has symmetries that leads to indistinguishable views, the error is computed as the average distance to the closest model point:

$$e_{\text{ADD-S}}(\hat{\mathbf{P}}, \bar{\mathbf{P}}, \mathcal{M}) \stackrel{\text{def.}}{=} \text{avg}_{\mathbf{x}_1 \in \mathcal{M}} \min_{\mathbf{x}_2 \in \mathcal{M}} \left\| \bar{\mathbf{P}}\mathbf{x}_1^* - \hat{\mathbf{P}}\mathbf{x}_2^* \right\|_2 \quad (4.3)$$

$$= \text{avg}_{\mathbf{x}_1 \in \mathcal{M}} \min_{\mathbf{x}_2 \in \mathcal{M}} \left\| (\bar{\mathbf{R}}\mathbf{x}_1 + \bar{\mathbf{T}}) - (\hat{\mathbf{R}}\mathbf{x}_2 + \hat{\mathbf{T}}) \right\|_2 \quad (4.4)$$

It's important to point out that  $e_{\text{ADD-S}}$  is more lenient compared to  $e_{\text{ADD}}$ , and should only be applied in cases where there is a definite presence of symmetry in the object and the estimated pose is already notably precise. Otherwise, using  $e_{\text{ADD-S}}$  becomes irrelevant since the estimation would be unfairly advantaged. In the illustrations below, we consistently provide both metrics, however it is up to the reader to assess the relevance of  $e_{\text{ADD-S}}$  based on the observed spacecraft.

### 4.2. Vizualisation and Quantitative Evaluation

With all these tools and the modifications made, we are now in a position to test Gen6D as-is, that is without retraining the model. For each spacecraft, we select approximately 200 reference images from a random image folder, and about ten

<sup>1</sup>In this context, the vector  $\mathbf{x}^*$  represents a vector that has been extended by appending a 1, specifically for the purpose of matrix multiplication.

hand-picked query images to ensure the best possible conditions. Indeed, we aim for the object to be in the foreground, fully within the frame, and under good lighting conditions. Since we are on the ESA competition track with 3D target models included, we skip the Structure From Motion (SFM) step with Colmap. Below are the best pose estimations made by Gen6D for each of the four objects.

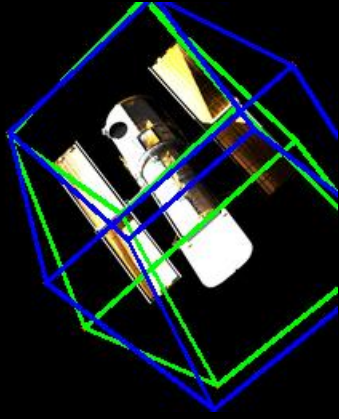


Figure 4.1: Hubble Space Telescope, no background, 256x256 query image,  $e_{\text{ADD}} = 2.925$ ,  $e_{\text{ADD-S}} = 1.183$

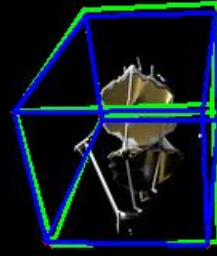


Figure 4.2: James Webb ST, no background, 256x256 query image,  $e_{\text{ADD}} = 1.415$ ,  $e_{\text{ADD-S}} = 0.808$

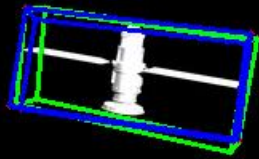


Figure 4.3: Cosmos Link, no background, 256x256 query image,  $e_{\text{ADD}} = 1.718$ ,  $e_{\text{ADD-S}} = 0.383$

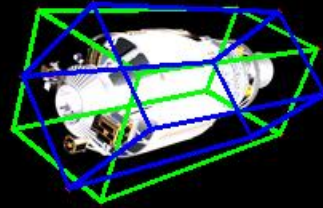


Figure 4.4: Rocket Body, no background, 256x256 query image,  $e_{\text{ADD}} = 1.713$ ,  $e_{\text{ADD-S}} = 0.252$

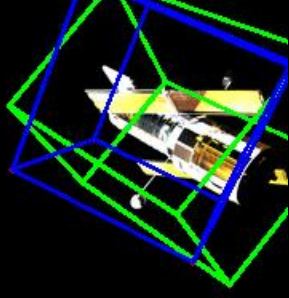


Figure 4.5: Hubble Space Telescope, no background, 256x256 query image,  $e_{\text{ADD}} = 6.514$ ,  $e_{\text{ADD-S}} = 1.571$

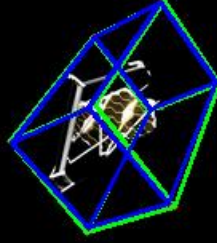


Figure 4.6: James Webb ST, no background, 256x256 query image,  $e_{\text{ADD}} = 2.224$ ,  $e_{\text{ADD-S}} = 1.261$

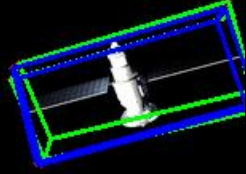


Figure 4.7: Cosmos Link, no background, 256x256 query image,  $e_{\text{ADD}} = 1.925$ ,  $e_{\text{ADD-S}} = 0.377$

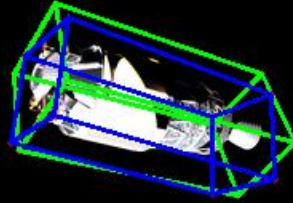


Figure 4.8: Rocket Body, no background, 256x256 query image,  $e_{\text{ADD}} = 1.982$ ,  $e_{\text{ADD-S}} = 0.501$

From our observations, it's notable that the results are particularly less accurate with the Hubble Space Telescope. This seems to stem from errors in the ground truth poses, as pointed out by Haochen, another student working on the project, particularly an axis inversion. Also, we often notice that the estimations on the Rocket Body are rotated around the object's axis of symmetry compared to the ground truth pose. Thus, it is relevant to consider the  $e_{\text{ADD-S}}$  metric in this case.

While the above results are great, the issue is that in many instances, the estimation is rather inaccurate. For completeness, we are now studying cases where Gen6D doesn't perform well, and we are trying to identify the reasons behind these shortcomings. The following text is provided by the authors to offer a clearer understanding of the intermediary results presented below: "[The column of vertical images] shows the viewpoint selection results. The first row shows the input image to the selector. The second row shows the input images rotated by the estimated in-plane rotation (left column) or the ground-truth in-plane rotation (right column). Subsequent 5 rows show the predicted (left) or ground-truth (right) 5 reference images with nearest viewpoints to the input image."



Figure 4.9: Hubble Space Telescope, no background, intermediary result of a poor estimation,  $e_{\text{ADD}} = 9.577$ ,  $e_{\text{ADD-S}} = 5.196$

In Figure 4.9, we observe good detection of the object, however, the selection of the nearest viewpoint is lacking, ultimately leading to an inaccurate pose estimation. This supports the earlier explanations regarding errors in the ground truth poses for the Hubble Space Telescope.

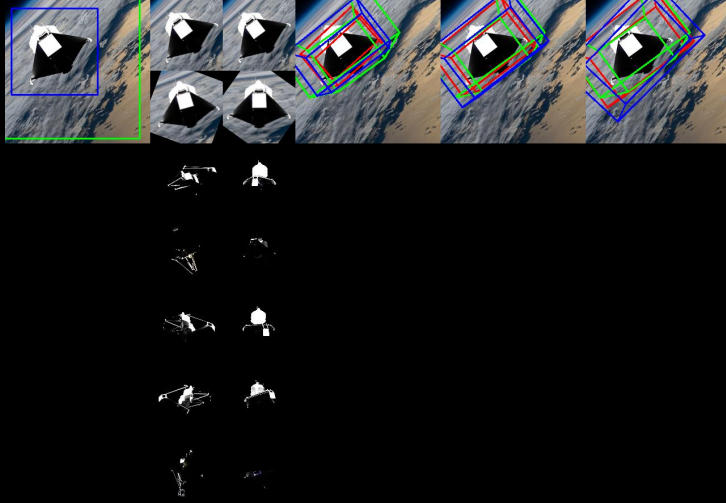


Figure 4.10: James Webb Space Telescope, with earth rendered background, intermediary result of a poor estimation,  $e_{\text{ADD}} = 10.934$ ,  $e_{\text{ADD-S}} = 4.317$

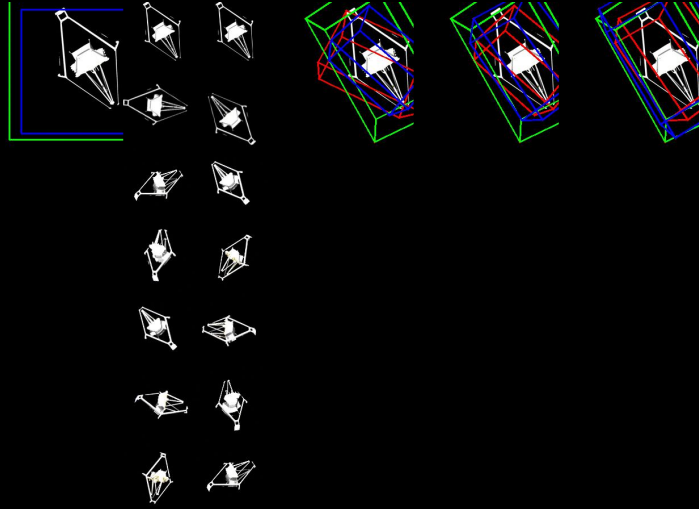


Figure 4.11: James Webb Space Telescope, no background, intermediary result,  $e_{\text{ADD}} = 1.060$ ,  $e_{\text{ADD-S}} = 0.556$

In Figure 4.10, we tested the model using reference images without a background and query images featuring the James Webb Telescope in front of the blue planet. While the object detection appears reasonably accurate, there is a significant difference in lighting conditions between the reference and query images. On the positive side, the refiner seems to work effectively as we observe an improvement in pose accuracy over the course of the three steps. Figure 4.11 is another example of the refiner doing great.

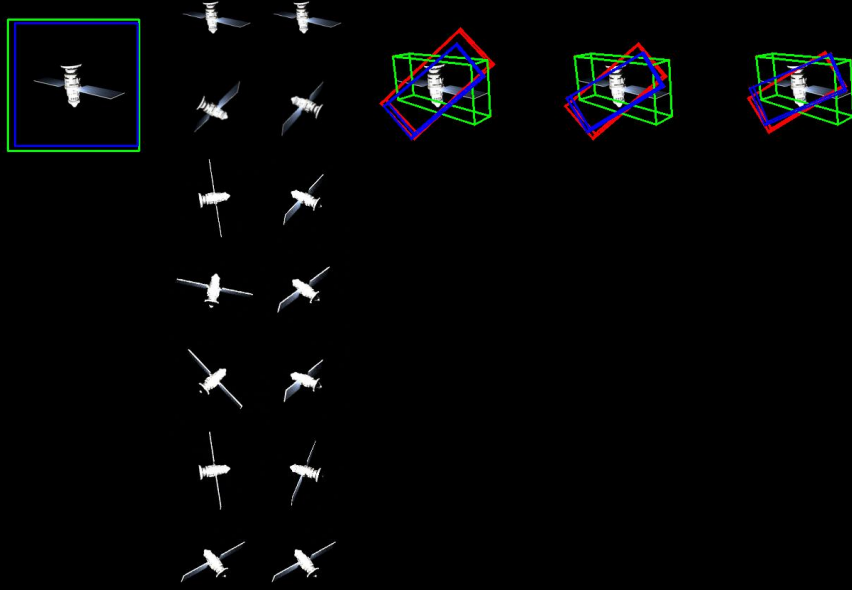


Figure 4.12: Cosmos Link, no background, intermediary result of a poor estimation,  $e_{\text{ADD}} = 11.094$ ,  $e_{\text{ADD-S}} = 6.127$

In Figure 4.12, we have a perfect detection, reasonably favorable viewpoints and yet a poor estimated pose at the end. However, it appears that increasing the number of refinement steps could lead to an improvement.

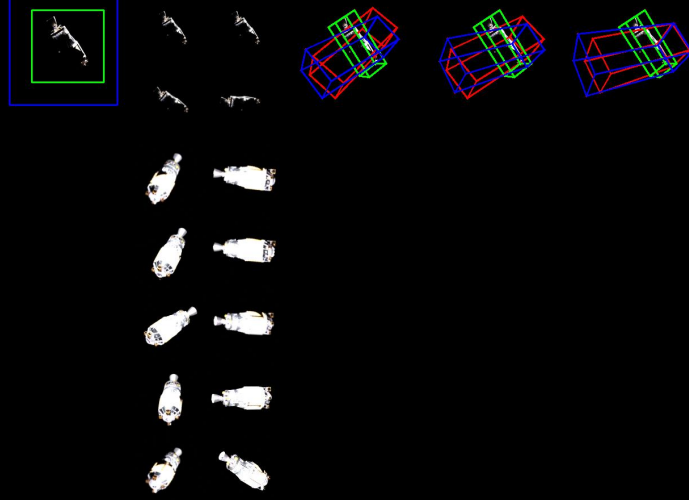


Figure 4.13: Rocket Body, no background, intermediary result of a poor estimation,  $e_{\text{ADD}} = 29.335$ ,  $e_{\text{ADD-S}} = 17.743$

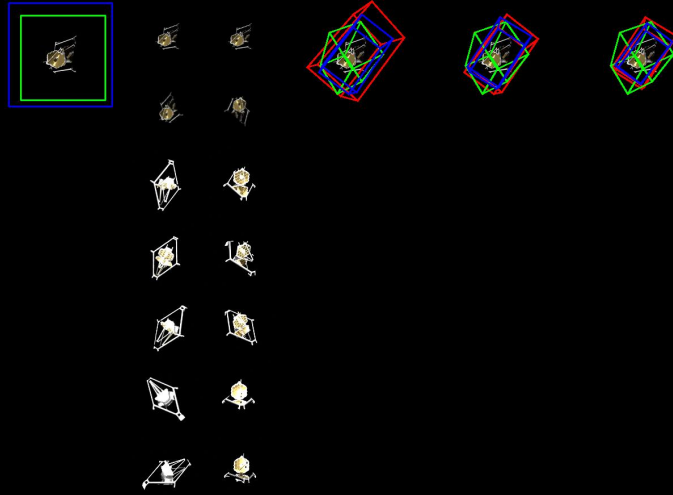


Figure 4.14: James Webb Space Telescope, no background, intermediary result of a poor estimation,  $e_{\text{ADD}} = 21.983$ ,  $e_{\text{ADD-S}} = 12.358$

In Figure 4.13 and Figure 4.14, we are encountering large differences in terms of lighting conditions. Specifically, for the James Webb Space Telescope, the query image reveals the interior of the antenna, which is not present in any of the selected viewpoints. Therefore, we conclude that Gen6D struggles with accurately estimating poses when there are substantial variations in lighting conditions.

## 5. Ways of improvements

As we observed in the previous chapter, Gen6D does not perform very well when generalizing to space objects. In this section, we will attempt to provide some tracks for improving the results of pose estimations.

### 5.1. More diverse Reference Images

Given the time constraints, we evaluated Gen6D under basic conditions, by using reference images from a folder of random synthetic images. To enhance the model's performance, one approach would be to diversify as much as possible the external conditions affecting the object's image through the camera, particularly lighting. As we observed, lighting plays a significant role during viewpoint selection. Having a more varied set of reference images would enable the model to better adapt to new conditions. Additionally, it's essential to ensure that the reference images are distributed as homogeneously as possible around the object.

### 5.2. Specialized Spacecraft Training Set

To assess the ability of Gen6D to generalize, we opted not to retrain the model but instead utilized the pretrained model made available by the authors. This approach may introduce a challenge, as highlighted in Table 2.1, where one of Gen6D's limitations is its training on everyday objects under favorable lighting conditions.

A viable solution would be to retrain Gen6D, focusing on specific components rather than the entire model. Specifically, retraining only the detector and the viewpoint selector, while keeping the refiner as is, could be effective. Indeed, retraining the detector would allow for the addition of new detecting scales for the object in query images, improving the accuracy of the estimated 2D bounding box. Additionally, retraining the viewpoint selector could prove beneficial for making pose estimations under poor lighting conditions. As for the refiner, as previously observed, it seems to perform pretty well and may not require retraining.

### 5.3. Improve Object Detection and Viewpoint Selection Algorithms

A limitation of the current architecture of Gen6D is its detector: it requires predefined object scales in the query images for detection. This is a significant constraint, especially in competitive scenarios where the object could be very close or very far away. Additionally, the current necessity to resize query images leads to information loss and adds a manual step to the process.

To address this, the implementation of a Feature Pyramid Network (FPN) could be a substantial improvement [5]. FPNs are known for their efficiency in detecting



objects at various scales through the utilization of a pyramidal hierarchy of deep convolutional networks. Incorporating a FPN into Gen6D's architecture would enhance its ability to detect objects regardless of their distance or size in the query images, thus potentially improving overall performance.

During the detection and viewpoint selection phase, we could also think about relying more on the 3D model (for now it is only used to get the estimated 3D bounding box size) and the segmented images (currently unused), to get better results.

## 6. Conclusion

In summary, after overcoming several software engineering challenges, we successfully reimplemented Gen6D for the *SPACECRAFT* dataset. While the current results are somewhat disappointing, we have identified several promising tracks for future exploration to enhance the pose estimation accuracy.

Knowing that the dataset team has developed new renderings this semester with dozens of new models, there is ample opportunity to further test Gen6D. More importantly, this development opens up the potential for retraining the model.

From a personal standpoint, this project was my first experience working as part of a significant and motivated team, and I thoroughly enjoyed the experience. It taught me to better organize my work, to operate on a remote server like Scitas Izar, to understand the functioning of a deep learning architecture in a practical, applied setting, to write Python scripts, and finally to improve my English skills through Zoom meetings and the writing of this very report.

# Abbreviations

**ABC** Abstract Base Class

**ADD** Average Distance of Model Points

**ADD-S** Average Closest Point Distance

**AI** Artificial Intelligence

**CNN** Convolutional Neural Network

**CVLab** Computer Vision Laboratory

**DART** Double Asteroid Redirection Test

**DoF** Degrees of Freedom

**ESA** European Space Agency

**FPN** Feature Pyramid Network

**HPC** High Performance Computing

**LEO** Low Earth Orbit

**SFM** Structure From Motion

**VGGNet** Visual Geometry Group Network

## A. Python Scripts

```
1 SPACECRAFT_ROOT='data/SpaceCraft'
2 class SpaceCraftCVLabDatabase(BaseDatabase):
3     def __init__(self, database_name):
4         super().__init__(database_name)
5         _, self.model_name = database_name.split('/')
6         self.img_ids = [str(k) for k in range(len(os.listdir(f'{SPACECRAFT_ROOT}
7             {self.model_name}/images')))]
8         self.model = self.get_ply_model().astype(np.float32)
9         self.object_center = np.zeros(3, dtype=np.float32)
10        self.object_vert = np.asarray([0,0,1], np.float32)
11        self.img_id2depth_range = {}
12        self.img_id2pose = {}
13        _, self.database_type = self.model_name.split('-')
14        # Track 1: SpaceCraft dataset camera intrinsic matrix
15        self.K=np.array([[1406.708374, 0.000000, 512.000000], [0.000000,
16            1406.708374, 512.000000], [0.000000, 0.000000, 1.000000]], dtype=np.float32
17        )
18        # Track 2: SwissCube dataset camera intrinsic matrix
19        # self.K=np.array([[607.5730322397038, 0.000000, 512.000000], [0.000000,
20            607.5730322397038, 512.000000], [0.000000, 0.000000, 1.000000]], dtype=np.
21            float32)
22        if self.database_type == 'test':
23            resize_factor = np.loadtxt(f"{SPACECRAFT_ROOT}/{self.model_name}/
24                resize_factor.txt")
25            self.K *= resize_factor
26            self.K[2, 2] = 1
27            # Note : This following portion of code is to uncomment when the ground
28            # truth poses are store in a JSON file,
29            # like for the SwissCube dataset.
30            # with open(f'{SPACECRAFT_ROOT}/{self.model_name}/scene_gt.json', 'r')
31            as file:
32                self.pose = json.load(file)
33
34        def get_ply_model(self):
35            fn = Path(f'{SPACECRAFT_ROOT}/{self.model_name}/{self.model_name}.pkl')
36            if fn.exists(): return read_pickle(str(fn))
37            ply = plyfile.PlyData.read(f'{SPACECRAFT_ROOT}/{self.model_name}/{self.
38                model_name}.ply')
39            data = ply.elements[0].data
40            x = data['x']
41            y = data['y']
42            z = data['z']
43            model = np.stack([x, y, z], axis=-1)
44            if model.shape[0]>4096:
45                idxs = np.arange(model.shape[0])
46                np.random.shuffle(idxs)
47                model = model[idxs[:4096]]
48            save_pickle(model, str(fn))
49            return model
50
51        def get_image(self, img_id):
52            # Track 1: SpaceCraft dataset image format
```

```

44         return imread(f'{SPACECRAFT_ROOT}/{self.model_name}/images/{int(img_id)
45         :04}.png')
46         # Track 2: SwissCube dataset image format
47         # return imread(f'{SPACECRAFT_ROOT}/{self.model_name}/images/{int(img_id)
48         :06}.jpg')
49
50     def get_K(self, img_id):
51         return np.copy(self.K)
52
53     def get_pose(self, img_id):
54         # Track 1: Gen6D's pose format
55         if img_id in self.img_id2pose:
56             return self.img_id2pose[img_id]
57         else:
58             pose = np.load(f'{SPACECRAFT_ROOT}/{self.model_name}/pose/pose{int(
59             img_id)}.npy')
60             self.img_id2pose[img_id] = pose
61             return pose
62         # Note : This following portion of code is to uncomment when the ground
63         # truth poses are store in a JSON file,
64         #         like for the SwissCube dataset. Do not forget to comment track
65         # 1 above.
66         # Track 2: SwissCube dataset json pose format
67         # raw_pose = (self.pose[f'{int(img_id) + 1}'])[0]
68         # rot_mat = np.asarray(raw_pose['cam_R_m2c'])
69         # rot_mat = np.reshape(rot_mat, (3, 3))
70         # transl_vec = np.asarray(raw_pose['cam_t_m2c'])
71         # transl_vec = np.reshape(transl_vec, (3, 1))
72         # pose = np.concatenate((rot_mat, transl_vec), axis=1)
73         # return pose
74
75     def get_img_ids(self):
76         return self.img_ids.copy()
77
78     def get_mask(self, img_id):
79         # Track 1: SpaceCraft dataset image format
80         return np.sum(imread(f'{SPACECRAFT_ROOT}/{self.model_name}/mask/{int(
81         img_id):04}.png'),-1)>0
82         # Track 2: SwissCube dataset image format
83         #return np.sum(imread(f'{SPACECRAFT_ROOT}/{self.model_name}/mask/{int(
84         img_id):06}.png'),-1)>0

```

Listing A.1: Python code of the data loading class SpaceCraftCVLabDatabase() from file database.py.

```
1 """
2 Author:      Jeremy Chaverot
3 Date:        November 20, 2023
4 Description: Transform a txt file with quaternions and the translation vector
               into multiple npy files containing the rotation matrix concatenated with
               the translation vector.
5 """
6
7 import numpy as np
8 import sys
9 import os
10
11
12 def quaternion_to_matrix(Q, translation):
13     """
14         Covert a quaternion and translation into a full three-dimensional
15         augmented rotation matrix.
16
17         Input
18         :param Q: A 4 element array representing the quaternion (q0, q1, q2, q3)
19         .
20         :param translation: A 3 element array representing the translation (x, y
21         , z).
22
23         Output
24         :return: A 3x4 element matrix representing the 3D rotation matrix
25         concatenated with the translation vector.
26     """
27
28     # Extract the values from arguments
29     q0 = Q[0]
30     q1 = Q[1]
31     q2 = Q[2]
32     q3 = Q[3]
33
34     x = translation[0]
35     y = translation[1]
36     z = translation[2]
37
38     # Compute the rotation matrix
39     r00 = 2 * (q0 * q0 + q1 * q1) - 1
40     r01 = 2 * (q1 * q2 - q0 * q3)
41     r02 = 2 * (q1 * q3 + q0 * q2)
42
43     r10 = 2 * (q1 * q2 + q0 * q3)
44     r11 = 2 * (q0 * q0 + q2 * q2) - 1
45     r12 = 2 * (q2 * q3 - q0 * q1)
46
47     r20 = 2 * (q1 * q3 - q0 * q2)
48     r21 = 2 * (q2 * q3 + q0 * q1)
49     r22 = 2 * (q0 * q0 + q3 * q3) - 1
50
51     # 3x3 rotation matrix concatenated with the 3x1 translation vector
52     matrix = np.array([[r00, r01, r02, x],
53                        [r10, r11, r12, y],
54                        [r20, r21, r22, z]])
55
56     return matrix
57
58 if __name__ == "__main__":
```

```
56
57 # Check if the correct number of arguments is provided
58 if len(sys.argv) != 3:
59     print("Usage: python quaternion_to_matrix.py </path/to/your/text/file>
60         </path/to/the/pose/folder>")
61     sys.exit(1)
62
63 file_path = sys.argv[1]
64 pose_folder_path = sys.argv[2]
65 file_content = None
66
67 try:
68     with open(file_path, 'r') as file:
69         file_content = file.read()
70 except FileNotFoundError:
71     print(f"The file {file_path} was not found.")
72     sys.exit(1)
73 except Exception as e:
74     print(f"An error occurred: {e}")
75     sys.exit(1)
76
77 poses = file_content.split('\n')[:-1]
78
79 # Iterate through each pose, apply the transformation and save it
80 for pose in poses:
81     image_id, obj_id, q0, q1, q2, q3, x, y, z = pose.split(',')
82     Q = np.array([q0, q1, q2, q3], dtype=np.float32)
83     translation = np.array([x, y, z], dtype=np.float32)
84     matrix = quaternion_to_matrix(Q, translation)
85     np.save(pose_folder_path + '/pose' + str(int(image_id)), matrix)
86
87 print(f"Number of transformation processed: {len(poses)}.")
```

Listing A.2: Python script quaternion\_to\_matrix.py to transform a txt file with quaternions and the translation vector into multiple npy files containing the rotation matrix augmented with the translation vector.

```
1 """
2 Author:      Jeremy Chaverot
3 Date:        December 10, 2023
4 Description: Invert the masks from a given folder.
5 """
6
7 import cv2
8 import os
9 import sys
10
11
12 def inverse_masks_in_folder(folder_path):
13     # Iterate through the list of files at the specified path
14     for filename in os.listdir(folder_path):
15         # Filter to include only png image files and exclude MacOS temporary files
16         if filename.endswith(".png") and not filename.startswith('._'):
17             mask_path = os.path.join(folder_path, filename)
18             try:
19
20                 mask = cv2.imread(mask_path, cv2.IMREAD_GRAYSCALE)
21                 if mask is None:
22                     print(f"Failed to read image: {mask_path}")
23                     continue
24                 inverted_mask = cv2.bitwise_not(mask)
25                 temp_path = os.path.join(folder_path, "temp_" + filename)
26                 cv2.imwrite(temp_path, inverted_mask)
27
28                 # Delete the original wrong mask
29                 os.remove(mask_path)
30
31                 # Rename the inverted mask
32                 os.rename(temp_path, mask_path)
33                 print(f"Inverted and replaced mask for: {mask_path}")
34             except Exception as e:
35                 print(f"Error processing {mask_path}: {e}")
36
37
38 if __name__ == "__main__":
39
40     # Check if the correct number of arguments is provided
41     if len(sys.argv) != 2:
42         print("Usage: python invert_mask.py <folder_path>")
43         sys.exit(1)
44
45     folder_path = sys.argv[1]
46     inverse_masks_in_folder(folder_path)
```

Listing A.3: Python script invert\_mask.py to invert the masks from a specified folder. We aim to have a white object set against a black background.



```
1 """
2 Author:      Jeremy Chaverot
3 Date:        January 01, 2024
4 Description: Resize the images from a given folder.
5 """
6
7 import os
8 import sys
9 from PIL import Image
10
11
12 def resize_images(folder_path, resize_factor):
13     # Iterate through the list of files at the specified path
14     for filename in os.listdir(folder_path):
15         # Filter to include only png image files and exclude MacOS temporary files
16         if filename.endswith(".png") and not filename.startswith('._'):
17             img_path = os.path.join(folder_path, filename)
18             with Image.open(img_path) as img:
19                 new_size = tuple([int(dim * resize_factor) for dim in img.size])
20                 resized_img = img.resize(new_size, Image.ANTIALIAS)
21                 temp_path = os.path.join(folder_path, "temp_" + filename)
22                 resized_img.save(temp_path)
23
24             # Delete the original image
25             os.remove(img_path)
26
27             # Rename the resized image
28             os.rename(temp_path, img_path)
29
30
31 if __name__ == "__main__":
32
33     # Check if the correct number of arguments is provided
34     if len(sys.argv) != 3:
35         print("Usage: resize.py <folder_path> <resize_factor>")
36         sys.exit(1)
37
38     folder_path = sys.argv[1]
39     factor = int(sys.argv[2])
40
41     resize_images(folder_path, factor)
```

Listing A.4: Python script `resize.py` designed to alter an image's size with respect to a specified resize factor.

```
1 """
2 Author:      Jeremy Chaverot
3 Date:        November 20, 2023
4 Description: Transform every images of a folder into jpg format.
5 """
6
7 import os
8 import sys
9 from PIL import Image
10
11
12 def transform_image(image_path):
13     img = Image.open(image_path)
14     new_image_path = image_path.split('.')[0] + '.jpg'
15     img.save(new_image_path)
16
17
18 if __name__ == "__main__":
19
20     # Check if the correct number of arguments is provided
21     if len(sys.argv) != 2:
22         print("Usage: python to_jpg.py </path/to/your/images>")
23         sys.exit(1)
24
25     folder_path = sys.argv[1]
26
27     # List of files in the given folder
28     all_files = os.listdir(folder_path)
29
30     # Filter the list to select only images and exclude MacOS temporary files
31     image_files = [file for file in all_files if file.lower().endswith(('.png',
32     '.jpg', '.jpeg', '.gif', '.bmp')) and not file.startswith('.')]]
33
34     num_images = len(image_files)
35
36     # Transform each image
37     for image_file in image_files:
38         image_path = os.path.join(folder_path, image_file)
39         transform_image(image_path)
40         os.remove(image_path)
41
42     print(f"Number of images transformed into .jpg: {num_images}")
```

Listing A.5: Python script to\_jpg.py to transform every images of a specified folder into jpg format.

```
1 """
2 Author:      Jeremy Chaverot
3 Date:        November 29, 2023
4 Description: Create the files val.txt, train.txt and test.txt according to a
               test percentage.
5 """
6
7 import os
8 import sys
9 import random
10
11
12 if __name__ == "__main__":
13
14     # Check if the correct number of arguments is provided
15     if len(sys.argv) != 3:
16         print("Usage: python format.py <object_name> <test_percentage>")
17         sys.exit(1)
18
19     object = sys.argv[1]
20     test_percentage = float(sys.argv[2])
21
22     if (test_percentage < 0 or 1 < test_percentage):
23         print("Wrong value for the variable <test_percentage>. Should be between
24               0 and 1 included.")
25         sys.exit(1)
26
27     # List of files in the given folder
28     all_files = os.listdir(f'data/SpaceCraft/{object}/images')
29
30     # Filter the list to select only images and exclude MacOS temporary files
31     image_files = [file for file in all_files if file.lower().endswith(('.jpg'))
32                   and not file.startswith('.') ]
33
34     num_images = len(image_files)
35
36     # Transform each image
37     with open(f'data/SpaceCraft/{object}/train.txt', 'w') as train, open(f'data/
38     SpaceCraft/{object}/test.txt', 'w') as test:
39         for image_file in image_files:
40             rand = random.random()
41             image_path = 'SpaceCraft/hubble/images/' + image_file
42             if (rand < test_percentage):
43                 test.write(image_path + '\n')
44             else: train.write(image_path + '\n')
45
46     print(f"Done splitting {num_images} images in train.txt and test.txt")
```

Listing A.6: Python script format.py to randomly generate the training set and the test set based on a specified probability. This script is used when creating both the reference set and the query set from a common set of images. Should be run from Gen6D's root folder.

## B. Scitas Izar Setup Tutorial

This tutorial aspires to guide you step by step in setting up the Scitas Izar environment to run Gen6D.

Step 1: Request permission via email to use the Scitas servers. Unfortunately, I can't provide more details, as I had the privilege of attending an EPFL course that automatically granted me access.

Step 2: Download the following GitHub repository:  
<https://github.com/JCHAVEROT/Gen6D/>

Step 3: Follow the instructions provided by the authors of Gen6D to set up the model (specifically, download the Gen6D pretrained model and place it in the correct location).

Step 3: Ensure you're on the EPFL internet network in person or connect to the EPFL VPN.

Step 4: Transfer the Gen6D files to your Scitas Izar session using on your computer the command:

```
1 rsync -azP <your/path/to/Gen6D> <username>@izar.epfl.ch:~/
```

Step 5: Create a virtual environment on your Scitas Izar Session using the instructions at: <https://scitas-doc.epfl.ch/user-guide/software/python/python-venv/> and activate it.

Step 6: Install all the necessary libraries for Gen6D to run, listed in the `requirements.txt` file, except for `openmpi`, `py-torch`, `py-torchvision`, and `cuda`.

Step 7: Compose your `execute.sh` file, in which you put the commands you want to run (use the provided template, do not modify the modules load). Notice that you need to replace the `<username>` section with yours.

```
1 #!/bin/bash
2 #SBATCH --chdir /scratch/izar/<username>
3 #SBATCH --partition=gpu
4 #SBATCH --qos=gpu_free
5 #SBATCH --gres=gpu:1
6 #SBATCH --nodes=1
7 #SBATCH --ntasks-per-node=1
8 #SBATCH --cpus-per-task=1
9 #SBATCH --mem 16G
10
11 echo STARTING AT `date`
```

```
12
13 echo "Loading modules"
14 module load gcc openmpi py-torch py-torchvision cuda
15
16 echo "Launching the virtual environment"
17 source ~/opt/izar1/venv-gcc/bin/activate
18
19 echo "Navigating to the directory and executing the task"
20 cd ~/Gen6D
21 python eval.py --cfg configs/gen6d_pretrain.yaml --object_name spacecraft/hubble
    --symmetric
22 python eval.py --cfg configs/gen6d_pretrain.yaml --object_name spacecraft/jwst
    --symmetric
23 python eval.py --cfg configs/gen6d_pretrain.yaml --object_name spacecraft/cosmos
    --symmetric
24 python eval.py --cfg configs/gen6d_pretrain.yaml --object_name spacecraft/rocket
    --symmetric
25
26 echo FINISHED AT `date`
```

Listing B.2: Bash script `execute.sh` to run a machine learning model on Scitas Izar EPFL. While the overall structure remains consistent, this script is specific to Gen6D's architecture.

Step 8: Use the following command to submit the job:

```
1 sbatch execute.sh
```

Step 9: Find the pose estimation results in folder `/Gen6D/data`. You can get them back on your computer using the commands:

```
1 rsync -azP <username>@izar.epfl.ch:~/Gen6D/data/performance.log <your/folder>
2 rsync -azP <username>@izar.epfl.ch:~/Gen6D/data/      <your/folder>
3 rsync -azP <username>@izar.epfl.ch:~/Gen6D/data/vis_inter <your/folder>
4 rsync -azP <username>@izar.epfl.ch:~/Gen6D/data/vis_final <your/folder>
```

**Pro Tip** Once the job submitted, you'll have access to the console log in the folder `scratch/izar/<username>`. Once on Scitas Izar, to go there, do the following command:

```
1 cd && cd /scratch/izar/<username>
```

You can find a lot of useful command to control your job at this link: <https://scitas-doc.epfl.ch/user-guide/using-clusters/running-jobs/>

I will keep the `README.md` file in the repository <https://github.com/JCHAVEROT/Gen6D/> updated to provide more information if necessary. Additionally, you'll find a compressed 'zip' file that includes well-organized SPACECRAFT images, allowing you to promptly test your setup and providing a template folder for your convenience.

## Bibliography

- [1] Y. Liu, Y. Wen, S. Peng, C. Lin, X. Long, T. Komura, and W. Wang. *Gen6D: Generalizable Model-Free 6-DoF Object Pose Estimation from RGB Images*. 2023. arXiv: 2204.10776.
- [2] Y.-B. Jia. “Quaternions.” In: (2022).
- [3] W. R. Hamilton and W. E. Hamilton. *Elements of Quaternions*. London: Longmans, Green, & Company, 1866.
- [4] S. Hinterstoisser, V. Lepetit, S. Ilic, S. Holzer, G. Bradski, K. Konolige, and N. Navab. “Model Based Training, Detection and Pose Estimation of Texture-Less 3D Objects in Heavily Cluttered Scenes.” In: *Computer Vision – ACCV 2012*. Springer Berlin Heidelberg, 2013, pp. 548–562.
- [5] T.-Y. Lin, P. Dollár, R. Girshick, K. He, B. Hariharan, and S. Belongie. *Feature Pyramid Networks for Object Detection*. 2017. arXiv: 1612.03144.