# 6D Pose Estimation Using Realsense D435i Depth Camera: A Multi-Method Approach

**By**

**Parth Sanghavi**

**Presented to the Research and Development Team of
Ghost Robotics**

**In Partial Fulfillment of the Requirements for the Completion of
Internship Program**

**Under the Supervision of
Dr. Avik De**

**GHOST ROBOTICS
August 2023**

# Abstract

In the realm of robotics, accurately estimating the 6D pose of objects in diverse environments is paramount for tasks such as manipulation, navigation, and interaction. This study delves into the challenge of 6D pose estimation using the Realsense D435i Depth Camera, particularly emphasizing scenarios with varying degrees of environmental clutter. Three distinct methods were developed to address this:

**[Method-1] Centroid-Based Pose Estimation with Ground Plane Normal Vector Analysis:**
- Tailored for environments with minimal clutter, this method focuses on estimating the centroid of a singular object lying on a plane.
- Additionally, it determines the ground plane's normal vector, providing essential insights for manipulator trajectory planning, ensuring minimal surrounding distractions.

**[Method-2] Model-Assisted 6D Pose Estimation via VFH and Point-to-Point ICP:**
- Designed for moderately cluttered environments, this approach extends the principles of Method-1. It incorporates the Viewpoint Feature Histogram (VFH) and Point-to-Point Iterative Closest Point (ICP) algorithms for detailed pose estimation.
- Object segmentation from the environment is achieved using RANSAC-based plane detection, fast Euclidean clustering, and Region Growing methods. This ensures accurate pose estimation even when objects are spaced as close as 3 cm apart.
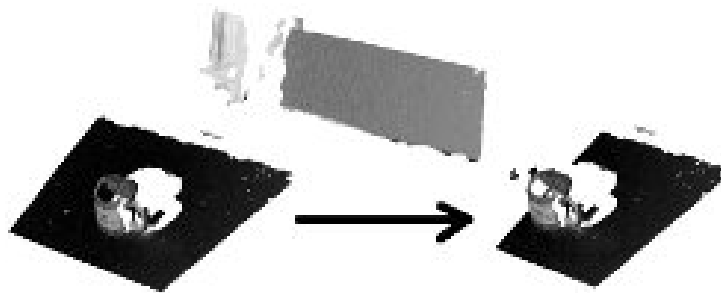
# Background

## 1) Filters for Point Cloud Processing:

The processing of 3D point clouds often necessitates the refinement and segmentation of data to isolate objects of interest and reduce computational overhead. Two commonly employed filters in this domain are the PassThrough filter and the Box Filter (often referred to as the Crop-Box Filter in the PCL library).

The **PassThrough filter** acts as a spatial window in the 3D space, selectively retaining points that lie within specified limits for each axis (x, y, z) of the 3D coordinate system. This filter is particularly effective in scenarios where the region of interest is confined to a specific spatial segment. By focusing on this segment, the filter not only reduces the influence of extraneous data and background noise but also enhances computational efficiency. For instance, in many applications, the filter is applied to the z-axis, representing depth, to create a depth-constrained window. This ensures that objects too far in the background or too close in the foreground are effectively filtered out, leaving behind a refined set of data points that represent the object of interest. Example of Pass through filter is shown in Figure 1.

On the other hand, the **Box Filter** defines a 3D bounding box in the point cloud space, specified by two 3D points: the minimum and maximum corners. All points within this bounding box are retained, while those outside are discarded. This filter is especially useful when the spatial region of interest is more complex or when further refinement post initial processing is required. By defining this 3D bounding box, the Box Filter ensures that only data points within a specific spatial cube are considered for subsequent operations, thereby providing a more focused and refined dataset. In many applications, the Box Filter is employed post other segmentation processes, like after using algorithms to identify planes or surfaces, to further refine and isolate the object of interest.



**Figure 1:** Example of Pass Through Filter applied on the Point Cloud

## 2) K-d Tree Clustering:

A K-dimensional tree, commonly known as a K-d tree, is a specialized data structure designed for organizing and partitioning data points in a k-dimensional space. At its essence, a K-d tree operates as a binary decision tree, facilitating efficient data clustering and search operations. One of its standout applications is in the realm of nearest neighbor searches, where it significantly optimizes the process. For the purposes of this study, we harness the power of the K-d tree implementation from the FLANN library, ensuring rapid and efficient clustering of our data.

## 3) Ransac Plane Segmentation:

RANSAC, renowned for its robust plane-fitting capabilities, is a preferred choice for detecting and segmenting planes in 3D point clouds. In the context of this method, the inliers, or the points that fit the model of the plane, are typically labeled as 'ground'. These inliers encompass all points with a Z value below a certain

threshold. Conversely, the outliers, or the points that deviate from the model, are labeled as 'non-ground' and are subsequently processed in the next segmentation phase.

The essence of the RANSAC method lies in its ability to identify the most substantial set of points that can be fitted to a plane within a specified threshold. The process initiates by randomly selecting three data points to define a plane. The parameters of this plane are then determined based on the equation:

$$ax + by + cz + d = 0$$

where x, y, z represent the 3D coordinates. Given these points, the constants a, b, c, and d are deduced. With the derived plane equation in hand, the inlier 3D points are identified based on a set threshold. After several iterations, the threshold yielding the maximum number of inlier ground points is chosen, and the corresponding plane is segmented from the rest of the point cloud.
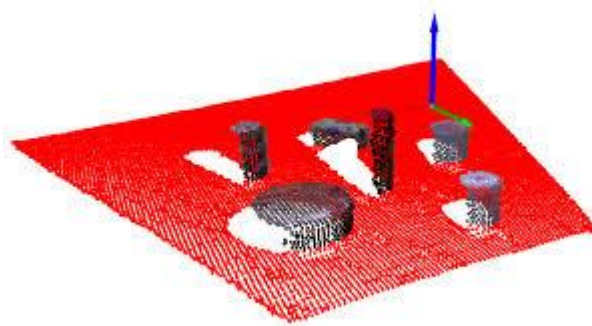
The number of iterations, n, required for the algorithm can be deduced from the probability of success, p using the formula:

$$N = \log(1-p) / \log(1 - (1 - \varepsilon)^s)$$

Where $\varepsilon$ is the outlier ratio and s is the sample size

The segmentation results using RANSAC consistently yielded a point cloud predominantly comprising inliers, or points that lie on the identified ground plane. The efficacy of RANSAC was evident, as it adeptly segmented all planar points within the point cloud that resided on the same plane. The algorithm adeptly treated non-ground points as outliers, which encapsulated above-ground features such as buildings, high-slope terrains, trees, and more.

One of the standout advantages of RANSAC is its robustness in estimating model parameters. It can discern these parameters with remarkable accuracy even when the point cloud is inundated with a significant number of outliers.



**Figure 2:** Example of Ransac Plane segmentation [Plane is shown in Red color]

## 4) Fast Euclidean Clustering (FEC)

Segmentation of point cloud data is pivotal in numerous applications, from remote sensing to autonomous vehicles. However, the inherent sparsity and lack of structure in point clouds, typically captured by 3D range sensors, pose challenges to efficient segmentation. Addressing this, the Fast Euclidean Clustering (FEC) algorithm emerges as a solution designed to expedite point cloud instance segmentation without imposing hefty computational demands.

**Understanding FEC's Core Mechanism:**

At its core, FEC is an evolution of the traditional Euclidean Clustering (EC) approach. While EC measures the proximity of unorganized points using the Euclidean (L2) distance metric and aggregates them into clusters, FEC refines this process. The equation governing this clustering mechanism is:

$$min \; ||P_i - P_i'||_2 \geq d_{th}$$

Here, $C_i = \{P_i \in P\}$ represents a distant cluster, distinct from $C_i' = \{P_i' \in P\}$, and $d_{th}$ is the maximum distance threshold.

**The Distinctiveness of FEC:**

What sets FEC apart is its point-wise scheme, a departure from the cluster-wise scheme prevalent in existing methodologies. This novel approach sidesteps the need to traverse every point repeatedly in nested loops, a process that is both time-intensive and memory-hungry. By focusing on individual points based on their input numbering order, as opposed to entire clusters, FEC achieves a significant reduction in computational overhead.
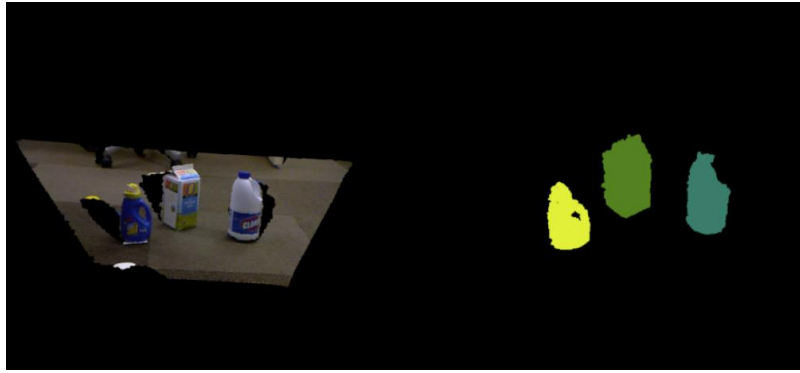
**Complexity and Efficiency: FEC's Winning Edge**

The time complexity of FEC is a testament to its efficiency. While traditional methods like EC might have exponential growth in time complexity with increasing data size, FEC's design ensures a more linear growth. This is primarily due to its point-wise scheme, which minimizes calls to the kd-tree search in the primary loop, especially crucial in scenarios with dense clusters or a high number of clusters.

**Key Advantages of FEC over Traditional Methods:**

- **Computational Efficiency:** FEC's design inherently minimizes the number of operations, particularly in data-rich scenarios. This leads to fewer kd-tree search calls, ensuring rapid processing times.

- **Memory Efficiency:** While methods like EC can be voracious in memory consumption, FEC's point-wise approach is frugal, making it ideal for real-time applications or systems with memory constraints.

- **Robustness:** Whether dealing with varying cluster densities or sizes, FEC's performance remains steadfast and consistent, unlike EC and RG which can see exponential time increases with denser clusters.

- **Detail Handling:** FEC shines in scenarios with intricate details, such as sparse and non-structural data points. It adeptly handles segments where EC and RG might falter, ensuring accurate segmentation.

- **Scalability:** FEC is built for the future. It gracefully handles increasing data sizes, ensuring its efficiency remains uncompromised even with vast point clouds.

In summation, FEC is not just an algorithm; it's a paradigm shift in point cloud segmentation. By re-envisioning the clustering process through a point-wise lens, it presents a streamlined, efficient, and robust solution, outclassing traditional methods in multiple dimensions.

**Figure 3**: Example of Fast Euclidean Clustering segmentation

## 5) Region Growing Method:

Region Growing is an essential image segmentation technique that clusters neighboring pixels with similar attributes to form coherent regions. The primary goal is to partition an image into distinct segments based on local homogeneity, thereby revealing meaningful structures. This approach capitalizes on the concept that adjacent pixels sharing attributes such as intensity or color likely belong to the same object or entity.

**Steps of Region Growing Algorithm:**

- **Seed Selection:** Choose seed pixel(s) from known regions of interest, acting as starting points.
- **Similarity Evaluation:** Examine neighboring pixels' attributes relative to the seed(s).
- **Homogeneity Test:** If neighboring pixels meet a predefined homogeneity criterion, add them to the growing region.
- **Iterative Growth:** Continuously iterate over neighboring pixels, expanding the region by including those meeting the homogeneity conditions.
- **Termination Condition:** Halt the process using a stopping criterion, like reaching a desired region size or attribute variation.

**Description:**

Region Growing excels in scenarios where objects possess consistent attributes or well-defined boundaries. It commences with seed pixel(s) that guide the expansion process. Neighboring pixels' attributes are compared to seeds, and if they match within predefined limits, they're incorporated. This iterative growth continues until the stopping criterion is met.

The method's versatility allows it to handle various data attributes, such as intensity or texture. By employing attribute comparison, Region Growing effectively segments images into regions, with each pixel assigned to the closest-matching region.

Despite its effectiveness, Region Growing faces challenges with noise and gradual attribute changes. Thus, its successful application depends on careful parameter tuning.

**Figure 4:** Example Result of Performing the seeded Region Growing Algorithm

In the context of our application, the Region Growing method could potentially be utilized for segmenting objects or structures in the scene based on their intensity. It can serve as a valuable pre-processing step to isolate relevant regions of interest, which can subsequently be used for more advanced tasks such as object recognition or orientation analysis.

## 6) Object Descriptors

To ascertain the pose of an object, it's imperative first to recognize the object. This recognition is achieved by quantifying the object's geometric attributes, such as surface curvature, surface normal directions, point distribution, and more. This quantification process births what is known as a descriptor. Descriptors fall into two primary categories: local and global, differentiated by their approach to geometric quantification. The choice between them hinges on various factors, including the operating environment, desired accuracy, computational speed, or even individual preference.
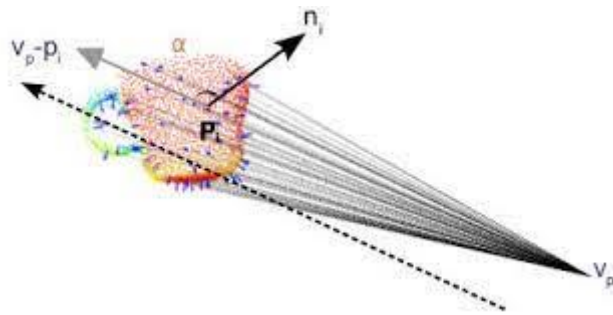
**Local Descriptors:** These zero in on localized areas within the point cloud, capturing the geometry of individual points. Given their focus on singular points, they necessitate calculations across multiple points to aptly describe an object. This multiplicity can be computationally intensive and memory-consuming. To mitigate this, descriptors are typically computed only at keypoints, specific to the descriptor in use. A standout advantage of local descriptors is their resilience to occlusions. They can recognize objects even with partial views, making them particularly adept in cluttered environments.

**Global Descriptors:** In contrast, global descriptors provide a holistic view, encapsulating the overall shape of the point cloud by leveraging all available points. This comprehensive approach means that typically, a single descriptor suffices for an apt description, ensuring quicker computations and lesser memory usage compared to local descriptors. However, they falter when faced with occlusions, as their descriptions rely on complete views. On the upside, they exhibit a higher tolerance to noise. A prerequisite for global descriptors is the segmentation process, which isolates the object's point cluster from the larger point cloud.

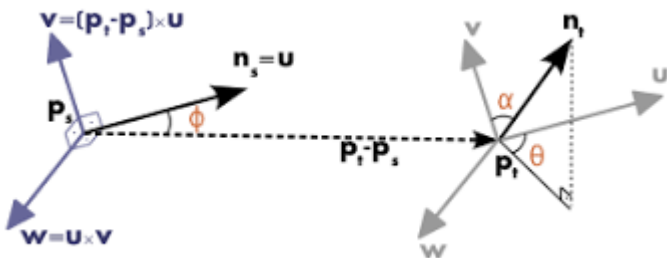## 7) Viewpoint Feature Histogram (VFH)

Given the project's operational conditions, the View Point Feature Histogram (VFH), a global descriptor, emerges as the chosen tool. Renowned for its high object recognition accuracy, noise resilience, and swift computation, the VFH descriptor, originally comprised two segments: the viewpoint direction component and the extended Fast Point Feature Histogram (FPFH) component. To bolster its robustness, the Shape Distribution Component (SDC) was integrated into its PCL implementation.

**Viewpoint Component:** This segment captures the surface normal direction in relation to the camera's central viewpoint. This is achieved by measuring the angle between a vector along the camera's z-axis and the computed normal at each point cluster. These angles populate a 128-bin histogram.



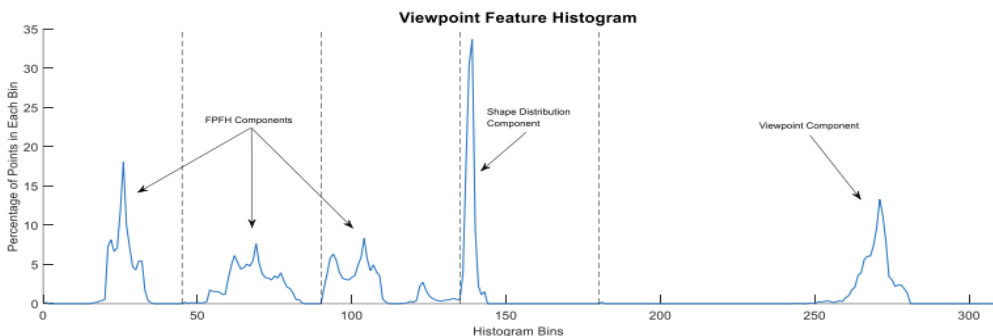**Figure 4:** Figure showing how Viewpoint bins are calculated

**Extended FPFH:** This component elucidates the relationship between surface normals and the point cluster's centroid. The centroid's computation involves averaging all points in the cluster. Given the absence of a natural vector for the centroid, one is drawn from the camera's viewpoint to the centroid. Using this vector as a normal, a Darboux coordinate frame, is constructed. The normals at each cluster point are then measured against this frame, and the resulting angles are stored in three 45-bin histograms.



**Figure 5:** Figure showing how extended FPFH bins are calculated

**Shape Distribution Component (SDC):** Initially absent from the VFH descriptor, the SDC, sourced from the Cluster View Point Histogram (CVFH) descriptor, enhances the VFH's capabilities. It examines the point distribution around the centroid, measuring distances from the centroid. This allows the descriptor to distinguish between objects with similar sizes and point normal distributions but varying point distributions. These distances are cataloged in a 45-bin histogram.

When combined, the histograms from the extended FPFH, SDC, and viewpoint component culminate in the comprehensive VFH, spanning 308 bins. A representative VFH is depicted in Figure 6.



**Figure 6:** Sample Resulting Viewpoint Feature Histogram

## 8) Iterative Closest Point Algorithm:

ICP is fundamentally an iterative method of registration that aligns two point clouds: a source and a target. The primary objective is to find the best possible alignment (in terms of rotation and translation) such that the source point cloud matches the target as closely as possible.

The ICP algorithm can be broken down into a series of steps:

- **Selection**: This involves sampling points from the point clouds to expedite the registration process. Given the potential vastness of point clouds, depending on the sensor's resolution, processing can be computationally intensive. By sampling, or reducing the number of points, the algorithm becomes more efficient. Methods like uniform sampling, which maintains the object's geometry, are often employed.

- **Matching**: This step is about estimating correspondences between points in the source and target point clouds. The most straightforward approach is to find the closest point in the target cloud for each point in the source cloud. Advanced data structures, like kd-trees, can speed up this search.

- **Rejection**: Not all correspondence estimated in the matching phase are valid. Some might be due to noise or other anomalies. The rejection phase filters out these invalid correspondences, improving the accuracy of the registration.

- **Alignment**: Once valid correspondences are identified, the algorithm estimates the transformation (rotation and translation) that best aligns the source to the target. This is often achieved by minimizing an error metric, such as the point-to-point or point-to-plane distance between corresponding points.

- **Iteration**: The above steps are repeated until the algorithm converges, i.e., until the difference between consecutive iterations is below a certain threshold or a set number of iterations is reached.
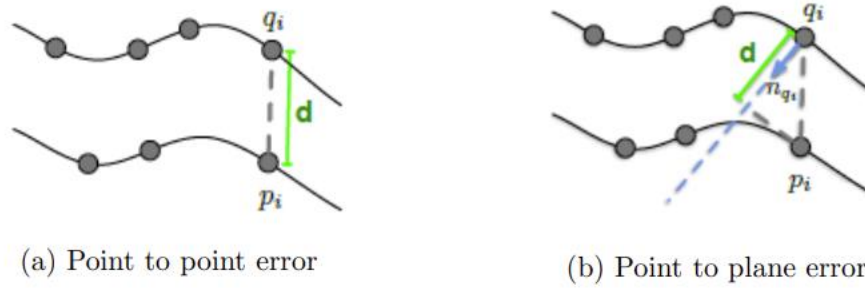
**Error Metrics:**

Two primary error metrics guide the alignment process:

- **Point-to-Point:** This metric measures the Euclidean distance between corresponding points in the source and target clouds. It's the most commonly used metric in ICP.

$$E_{point-to-point}(T) = \sum_{k=1}^{N} ||Tp_k - q_k||^2$$

$$E_{point-to-plane}(T) = \sum_{k=1}^{N} ((Tp_k - q_k).n_{q_k})^2$$

- **Point-to-Plane:** Introduced as an enhancement to the point-to-point metric, the point-to-plane metric measures the distance between a point in the source cloud and a plane defined by a point and its normal in the target cloud. This metric often results in faster and more stable convergence.

(a) Point to point error  (b) Point to plane error

**Figure 8:** Error Metrics

**Termination Criteria:**

For the iterative process to end, certain termination criteria are set:

- **Maximum Iterations:** The algorithm stops after a predefined number of iterations.

- **Transformation Epsilon:** The algorithm stops if the change in transformation between consecutive iterations is below a threshold.

While ICP is renowned for its simplicity and efficacy, it is not without its challenges. The algorithm's sensitivity to the initial alignment can sometimes lead to local minima, resulting in subpar alignments. Moreover, ICP might not perform optimally with point clouds that have significant noise or missing data sections.

# Methodology

## [Method 1] Centroid-Based Pose Estimation with Ground Plane Normal Vector Analysis:

### Objective:
To estimate the 6D pose of an object in a minimally cluttered environment by determining the centroid of the object and the direction of the ground plane with respect to the object.

### Environment Applicability:
This method is best suited for environments with minimal clutter, where there is a single object of interest lying on a plane with minimal surrounding distractions.
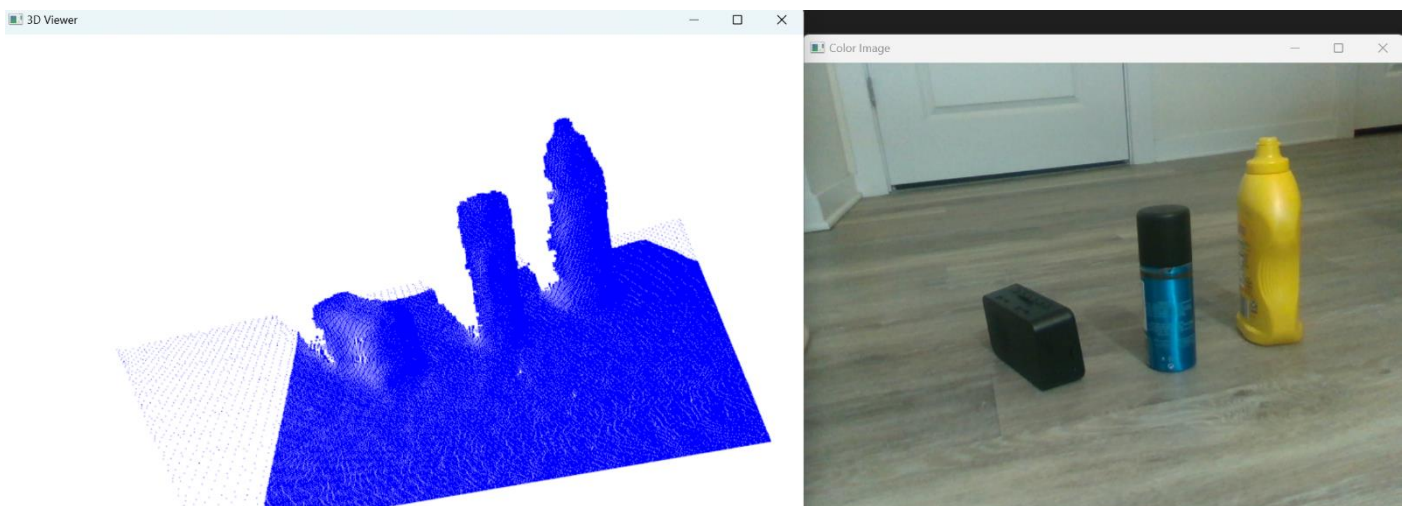
### Procedure:

**1) Device Setup, Data Collection, and User Interaction**

The procedure begins with the initialization and configuration of the Realsense device. Upon activation, the system checks for the device's connectivity, setting up depth and color streams to capture the necessary data. A critical aspect of this phase is obtaining the intrinsic parameters of the depth stream, such as focal lengths and principal points, which are vital for interpreting depth data accurately.

Following device setup, the system delves into data acquisition and pre-processing. Frames are captured, with particular attention given to aligning depth frames with color frames, ensuring data consistency. Given the natural noise associated with depth data, it undergoes a series of filtering processes, including decimation, disparity transformation, and spatial and temporal filters. Once refined, the depth frame is transformed into a point cloud, offering a 3D representation of the scene.

User interaction is facilitated by displaying the color image, allowing users to pinpoint a specific area of interest on the object. This interaction is pivotal, serving as a reference for subsequent processing stages. The 2D point selected on the color image is then mapped to its 3D counterpart in world coordinates, leveraging depth data and camera intrinsic.
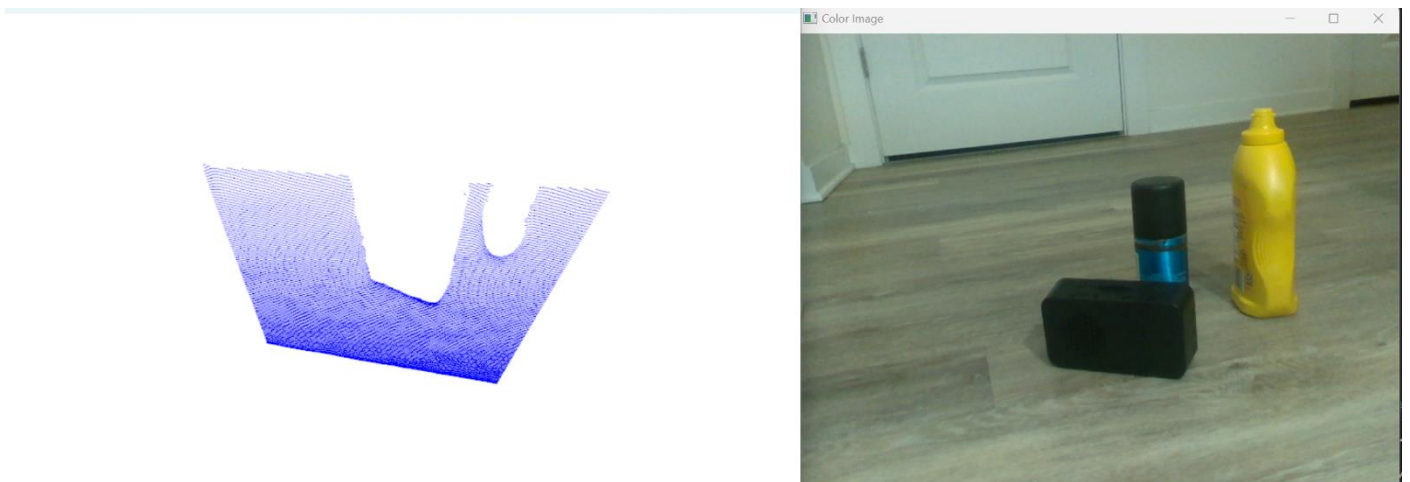


**Figure 9:** Raw Point Cloud of the Frame on which the User Clicks

## 2) Point Cloud Processing and Plane Segmentation:

The point cloud, a 3D representation of the scene, is then subjected to further processing. Initially, the data undergoes refinement using the PassThrough filter, which acts as a spatial window in the 3D space. This filter focuses on a specific region, reducing extraneous data and noise, and enhancing computational efficiency. In the given context, this filter is applied around the user-selected point, particularly emphasizing the z-axis. By creating a window around the depth value of the selected point, the object of interest remains central to the analysis, while unrelated foreground or background objects are filtered out.

RANSAC, a robust algorithm, is then employed for plane segmentation. Despite its random sampling nature, RANSAC's strength lies in its ability to fit models even in noisy environments. In the described scenario, Ransac-plane segmentation algorithm identifies the dominant plane, possibly representing a ground or table surface, in a multi-threaded manner for enhanced efficiency.

Post RANSAC segmentation, the Box Filter is applied to further refine the point cloud. This filter, by defining a 3D bounding box in the point cloud space, retains points within the box and discards those outside. In the provided context, after removing the dominant plane's inliers, the Box Filter refines the point cloud around a user-selected point, ensuring the focus remains on the object of interest.



**Figure 10:** Plane Segmentation using Ransac [inliers point cloud]

## 3) Object segmentation using Clustering algorithm:

After the Box Filter's refinement, the Fast Euclidean Clustering (FEC) algorithm is employed to segment the remaining points into distinct clusters, each potentially representing different objects or parts of the scene. FEC operates by selecting an initial point and forming a cluster by aggregating nearby points based on a distance threshold. This process is iteratively applied to all unprocessed points, resulting in multiple clusters that are distinct yet computationally efficient to identify.
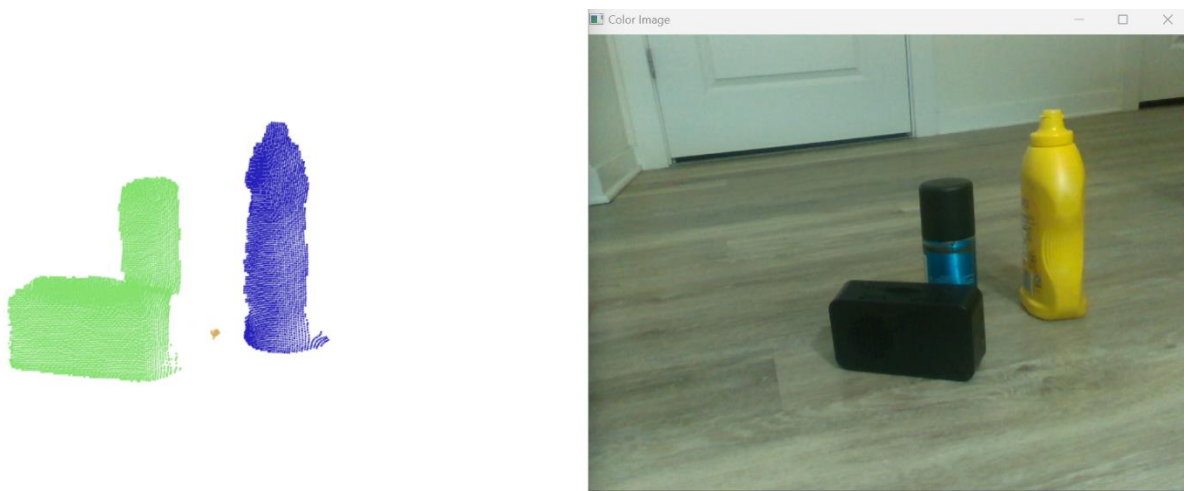
After the application of Fast Euclidean Clustering (FEC), multiple distinct clusters emerge, each potentially representing different objects or parts of the scene. To identify the cluster corresponding to the object of interest, a method centered on the concept of centroids is employed.

Firstly, the centroid of each cluster is computed. A centroid, in this context, is the average position of all the points within a cluster, effectively representing its geometric center. With these centroids in hand, the next step is to determine which cluster is most relevant to the user's point of interest. This is achieved by identifying

the centroid closest to the point clicked by the user. The cluster associated with this nearest centroid is then deemed the object of interest.

Subsequent to this selection, the RANSAC algorithm is applied specifically to this chosen cluster. The goal here is to segment and identify the primary plane or surface of the object. However, a challenge arises: the calculated centroid, being an average of all visible surface points, might not accurately represent the true center of the object, especially if the object is partially occluded or viewed from a particular angle.

To address this, once the primary visible surface of the object is identified using RANSAC, the centroid is adjusted. It's pushed 2 cm in the direction of the plane's normal vector. This adjustment ensures that the centroid is not just representative of the average surface position but is nudged closer to a more accurate representation of the object's center, considering the visible surface. This refined centroid provides a more accurate and meaningful representation of the object's position within the 3D space.



**Figure 11:** Result of Fast Euclidean Clustering Algorithm

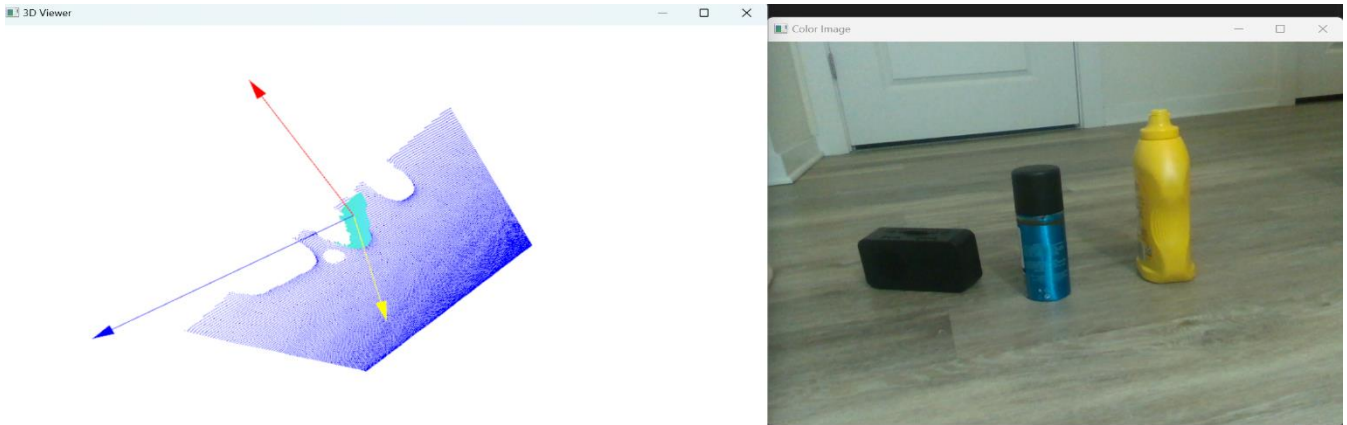**4) Estimation of Object Orientation and Establishing a Coordinate Frame :**

The culmination of the procedure zeroes in on discerning the orientation of the ground plane, a pivotal step that informs subsequent interactions with the object, especially for tasks like robotic gripping.

To achieve this, the procedure leverages the plane's normal vector, which was previously identified using the RANSAC plane segmentation. This normal vector is of paramount importance because it essentially corresponds to the gravity vector. In the context of robotic operations, understanding this gravity vector is crucial. It indicates the direction from which the gripper should not approach, as there would be no space for movement. Instead, the gripper should approach from the opposite direction, ensuring a smooth and unobstructed grip.
However, merely knowing the direction of approach isn't enough for precise robotic operations. It's essential to understand the object's orientation in a more holistic manner. This is where Principal Component Analysis (PCA) comes into play. PCA is employed to identify the longest side of the cluster, which essentially represents the object's primary axis of variance. A vector aligned with this longest side is then determined.

Subsequently, a second vector is computed, which is perpendicular to the first vector and lies within the plane of the object. This vector provides additional orientation information, especially about the object's width.
Finally, a third vector is derived, which is orthogonal to both the previously determined vectors. This trio of vectors, combined, forms a robust coordinate frame. The beauty of this frame is that it's centred at the adjusted position of the centroid, providing a comprehensive spatial understanding of the object.

In essence, this multi-faceted approach not only isolates the object but also meticulously contextualizes it within its environment. By establishing a clear coordinate frame centred on the object, the procedure offers invaluable orientation insights, paving the way for precise robotic interactions or further analytical endeavours.



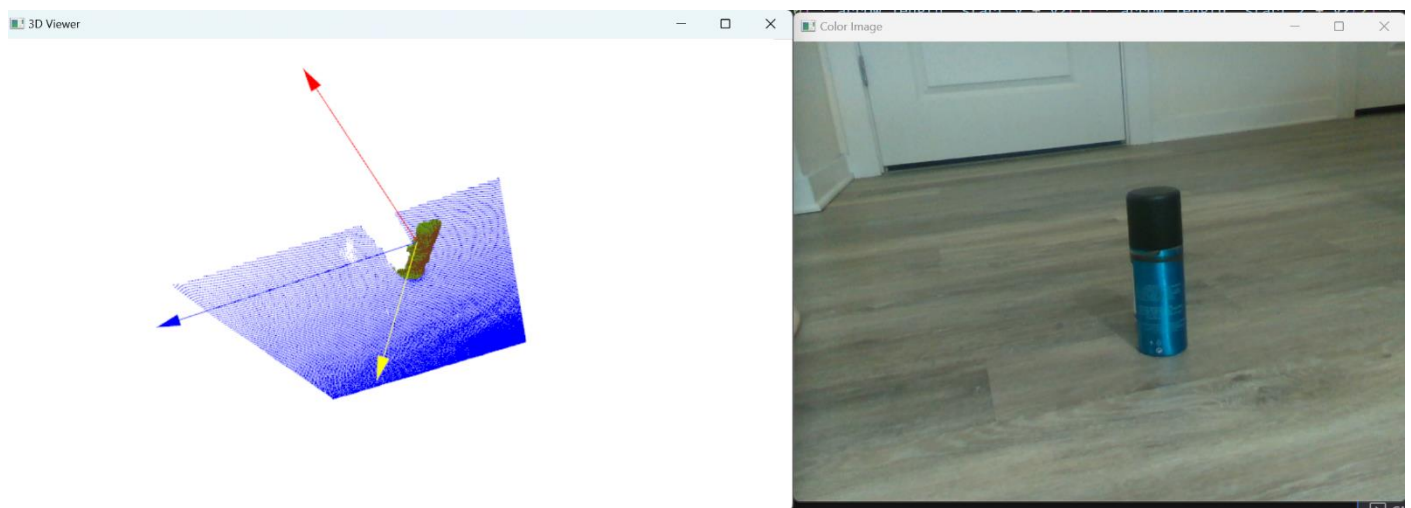**Figure 12:** 6D Pose of the Object Point Cloud in Camera Coordinate frame
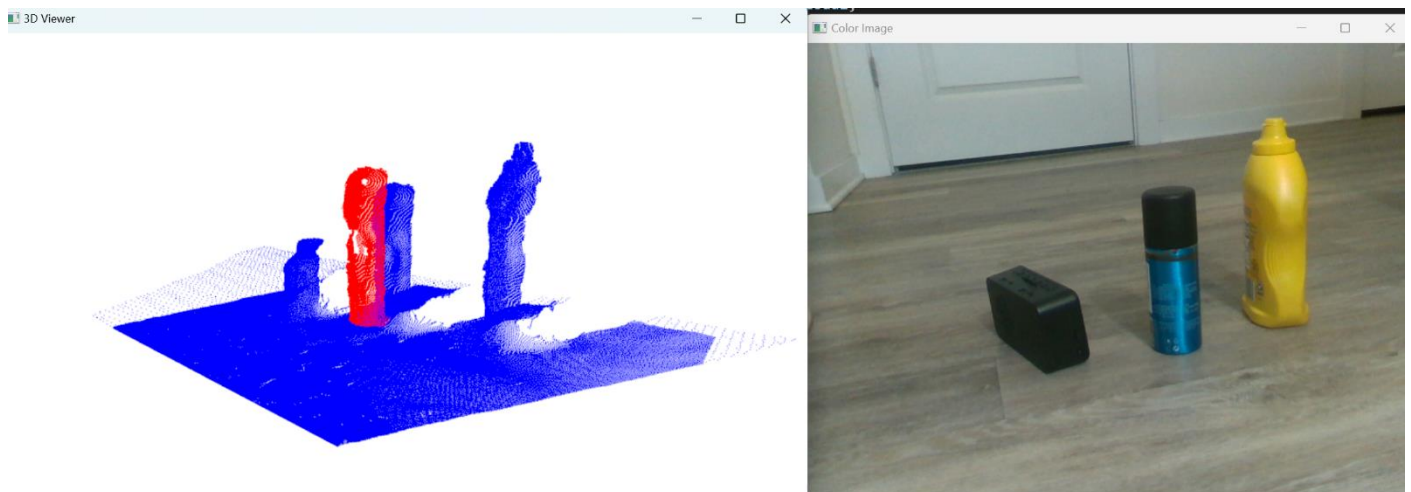
# Evaluation and Testing:

**Test Environment:**

The evaluation was carried out in a controlled setting, using a variety of point objects that are commonly found in kitchen and household environments. These objects were chosen to simulate real-world scenarios where the algorithm would be most applicable. The tests were designed to assess the algorithm's robustness and efficiency, not only when objects are isolated but also when they are in close proximity or amidst clutter. This is crucial as real-world environments are seldom perfectly organized, and the algorithm needs to perform reliably even in less-than-ideal conditions.

**Figure 1: Singular Object Scenario** - This test focuses on the algorithm's ability to detect and process a single object. Here, a lone object is positioned on a flat ground plane. The simplicity of this scenario serves as a baseline, demonstrating the algorithm's capability in an uncluttered environment.



**Figure 13:** Singular Object Scenario - A lone object positioned on a flat ground plane.

**Figure 2: Multiple Object Arrangement** - To challenge the algorithm further, multiple objects were placed on the ground plane, spaced at intervals of 5-6 cm. This test evaluates the algorithm's proficiency in distinguishing between closely spaced objects. It showcases the algorithm's resilience and accuracy in a more complex, cluttered environment, highlighting its potential for real-world applications where objects might be closely packed.

## Performance Metrics:

To gauge the efficiency of the procedure, we measured the time taken for each significant step. The timings, presented in milliseconds, offer insights into the computational demands of each phase and the overall efficiency of the system.

**Timings Breakdown:**

| Steps | Procedure | Time (in milliseconds ) |
|-------|-----------|-------------------------|
| 1 | Decimation Filter of Realsense SDK | 0.436 |
| 2 | Pass Through Filter of PCL | 2 |
| 3 | Box Filter + Voxel Grid Filter | 3.2 |
| 4 | Ransac Plane Segmentation | 60 |
| 5 | Fast Euclidean Clustering | 120 |
| 6 | Cluster Centroid and Basis Vectors calculation | 0.2 |

## Analysis:

From the timings breakdown, it's evident that the most time-intensive steps are the RANSAC Plane Segmentation and the Fast Euclidean Clustering. These steps, while computationally demanding, are crucial for accurate object segmentation and orientation estimation. On the other hand, filtering processes, despite their importance in refining the point cloud data, are relatively swift.
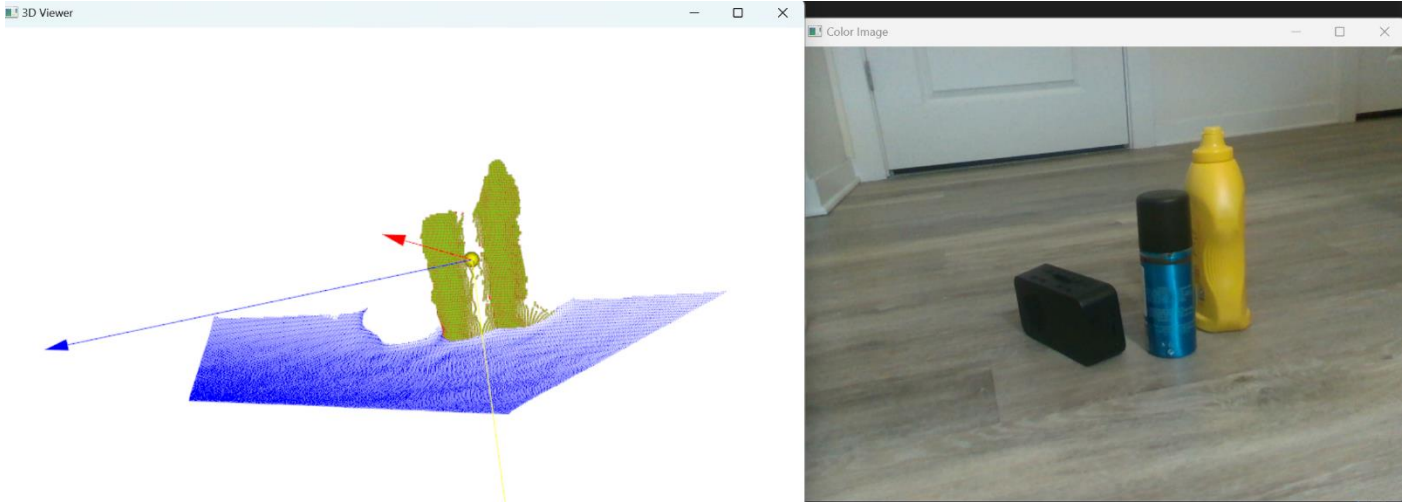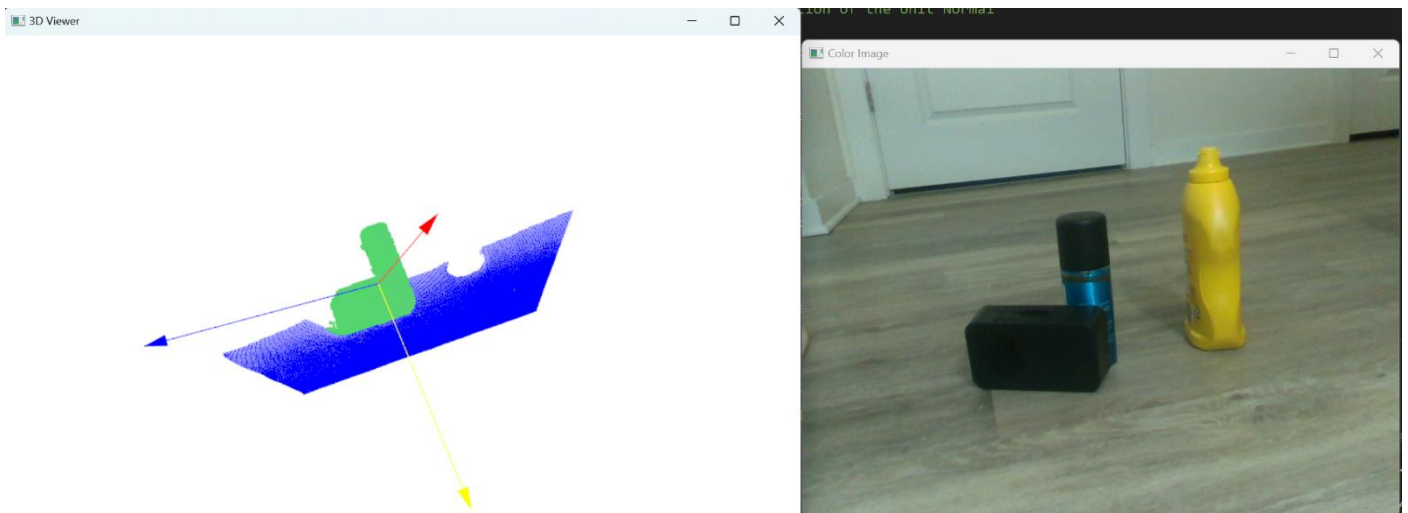
# Challenges and Limitations:

The development and testing of the algorithm revealed several challenges and limitations that need to be addressed for broader applicability and improved accuracy. Here are the primary challenges encountered:

**Object Proximity and Clutter:**
One of the significant challenges faced by the algorithm is the accurate detection and segmentation of objects that are in close proximity or touching each other. When objects are too close or cluttered, the algorithm struggles to distinguish between them, leading to potential inaccuracies in the segmentation process.
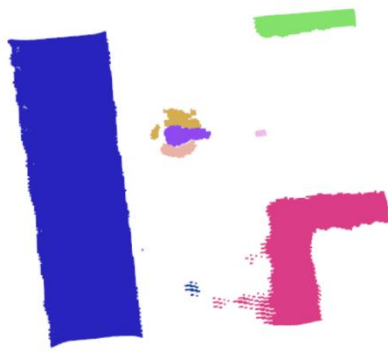


**Figure 15:** Cluttered Environment 1: Objects in Close Proximity, leading the algorithm to fail
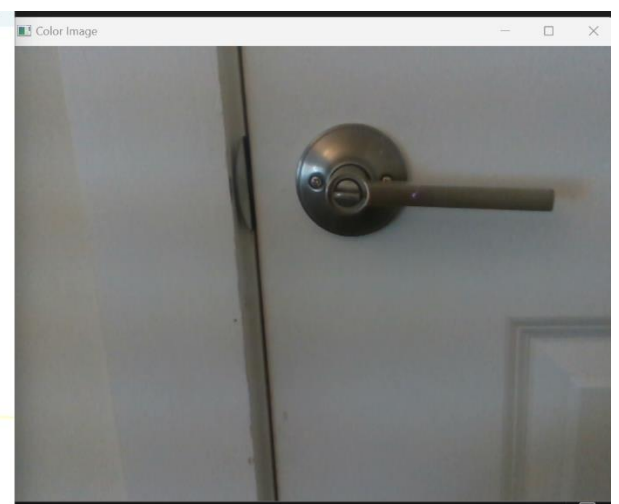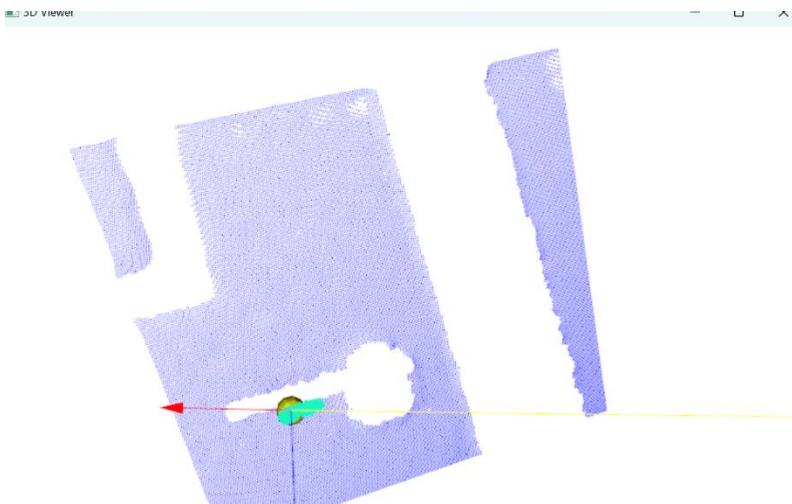


**Figure 16:** Cluttered Environment 2: Objects in Close Proximity, leading the algorithm to fail

**Non-Lambertian Surfaces and Door Handles:**
Objects with shiny surfaces, such as door handles, pose a unique challenge. These surfaces, termed non-Lambertian, reflect light in a manner that doesn't adhere to the typical diffuse reflection model. As a result, the depth sensor struggles to capture a consistent point cloud for such objects. This inconsistency manifests as fragmented or missing sections in the point cloud, making it challenging for the algorithm to detect and process the object accurately.
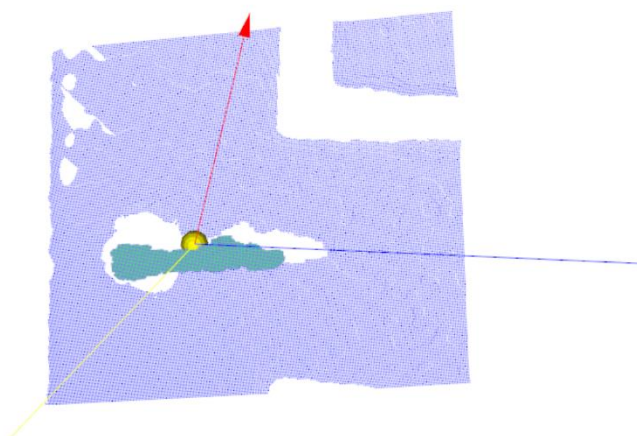
**Figure 17:** Illustration of Point Cloud Degradation due to Non-Lambertian Surface Reflections



**Figure 18:** The algorithm succeeds only when the use selects a point close to the middle of the door handle

**Validation with Modified Surfaces:**

To validate the challenge posed by shiny surfaces, an experiment was conducted where a door handle was wrapped in paper, a Lambertian surface. This modification aimed to provide a consistent reflection model for the depth sensor. The results showed a marked improvement in the algorithm's ability to detect and process the handle, underscoring the challenges posed by non-Lambertian surfaces.



**Figure 19:** Mitigating Reflection Challenges: Door Handle Wrapped in Paper for Consistent Point Cloud Capture

In conclusion, while the algorithm demonstrates significant promise in various scenarios, these challenges highlight areas for further refinement and research. Addressing these limitations will be crucial for enhancing the algorithm's robustness and applicability in real-world, diverse environments.

# [Method-2] Model-Assisted 6D Pose Estimation via Viewpoint Feature Histogram and Point-to-Point ICP:

## Objective:

Expanding on the groundwork established by the initial methodology, this approach endeavors to enhance the precision of object recognition and alignment by leveraging pre-existing models of the objects of interest. By integrating these models, the system is better equipped to discern and align objects even in intricate scenarios.

## Environment Suitability:

Designed with meticulous attention to detail, this methodology excels in environments where objects maintain a separation of roughly 2-3 cm. Such environments, often characterized by their cluttered nature, pose significant challenges due to the close proximity of objects. The inherent risk of misalignment or misrecognition in these settings amplifies the need for a robust and refined approach, which this methodology promises to deliver.

## Procedure:

### 1) Pre-processing and Refinement:

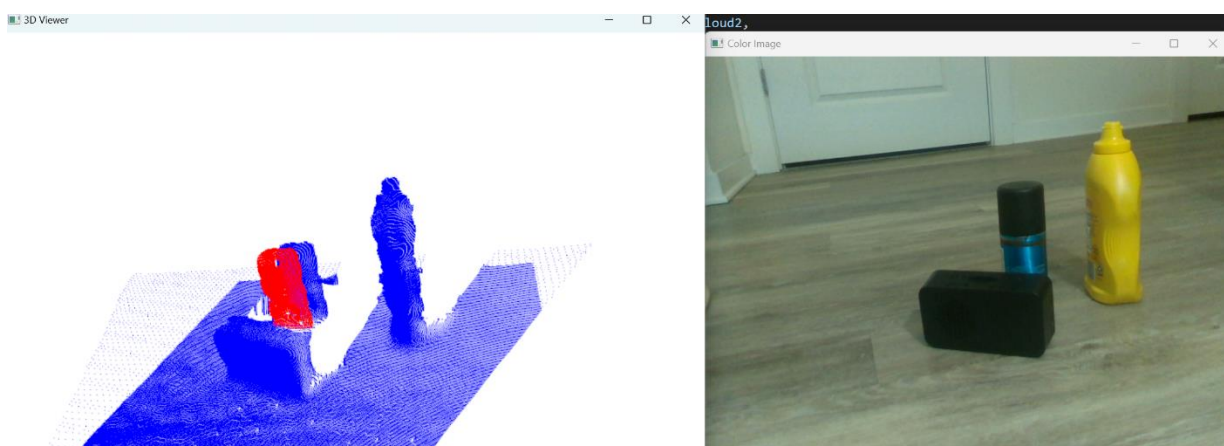**Statistical Outlier Removal:**
Once the cluster closest to the user's selected point is identified via Fast Euclidean Clustering (FEC), the point cloud undergoes a cleansing process. The statistical outlier removal filter is applied to eliminate anomalous points, ensuring a cleaner dataset for further operations.

**Voxel Grid Filter**:
To address the challenges posed by objects in close proximity, especially in densely populated environments, the voxel grid filter is introduced. This filter downsamples the point cloud, enhancing computational efficiency and accentuating the distinction between closely spaced objects.

**Region Growing Segmentation**:
In cluttered environments, FEC might occasionally amalgamate multiple objects into a singular cluster. To counteract this, Region Growing (RW) segmentation is employed. This method, rooted in the principles of surface normals and curvature estimations, is adept at differentiating between objects that share spatial intimacy. The efficacy of RW can hinge on the judicious selection of the seed point, underscoring the need for meticulous parameter tuning.

## 2) Feature Extraction and Initial Alignment:

### Normal Estimation:

The surface normals for both the segmented object point cloud and the model point cloud are computed. These normals are pivotal for the subsequent feature-based alignment process. The code employs the *pcl::NormalEstimation* class, which uses a KdTree to efficiently compute the normals for each point in the point cloud based on its local neighborhood.

### Viewpoint Feature Histogram (VFH):

VFH is a high-dimensional feature descriptor that captures the object's global characteristics. The code computes the VFH for both the segmented object (source) and the model point cloud. This is achieved using the *pcl::VFHEstimation* class. By encapsulating the object's global characteristics, the VFH facilitates an initial coarse alignment, setting the stage for the finer adjustments that follow. The matching of these VFH descriptors is then performed using a KdTree-based search, ensuring an efficient and accurate matching process.
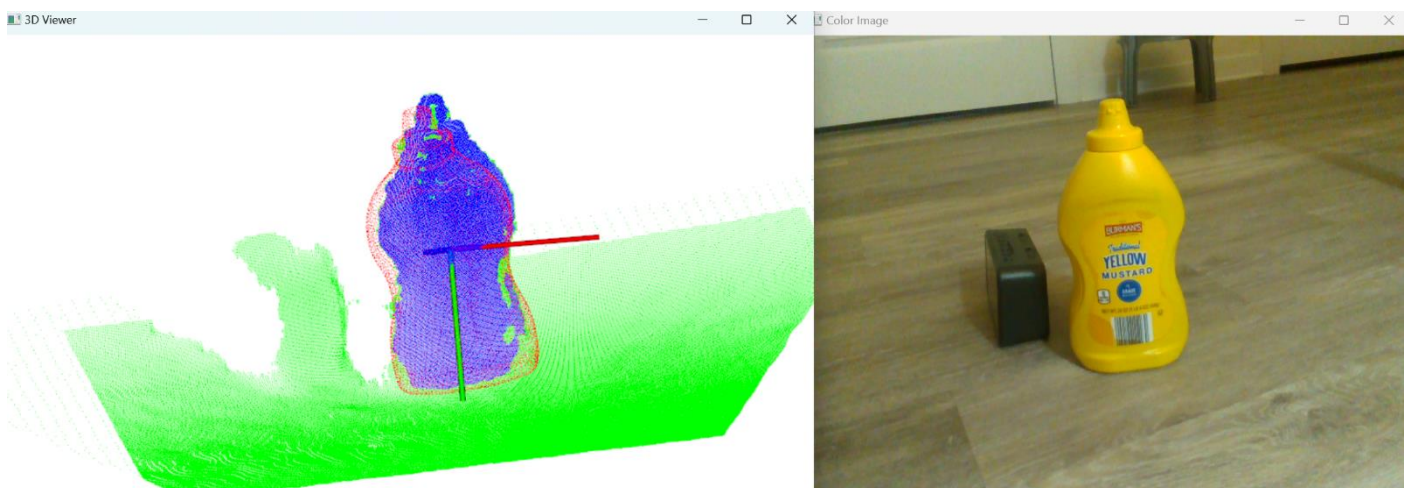
## 3) Fine Alignment using ICP:

### Centroid Calculation and Alignment:

The geometric centers or centroids of both the source and model point clouds are discerned. By translating both point clouds to the origin, a common reference point is established, priming them for the nuanced alignment via ICP.

### ICP with Point-to-Point Error Metric:

The Iterative Closest Point algorithm, renowned for its precision, is then invoked to refine the alignment between the source and model point clouds. The point-to-point error metric ensures a snug alignment between corresponding points in the two clouds, guaranteeing an alignment that's both tight and accurate. The code employs the *pcl::IterativeClosestPoint* class, which iteratively refines the alignment between the source and target point clouds until convergence.



**Figure 21:** Model Point Cloud (Red color) Projection onto Real-World Object with ICP-Driven Transformation
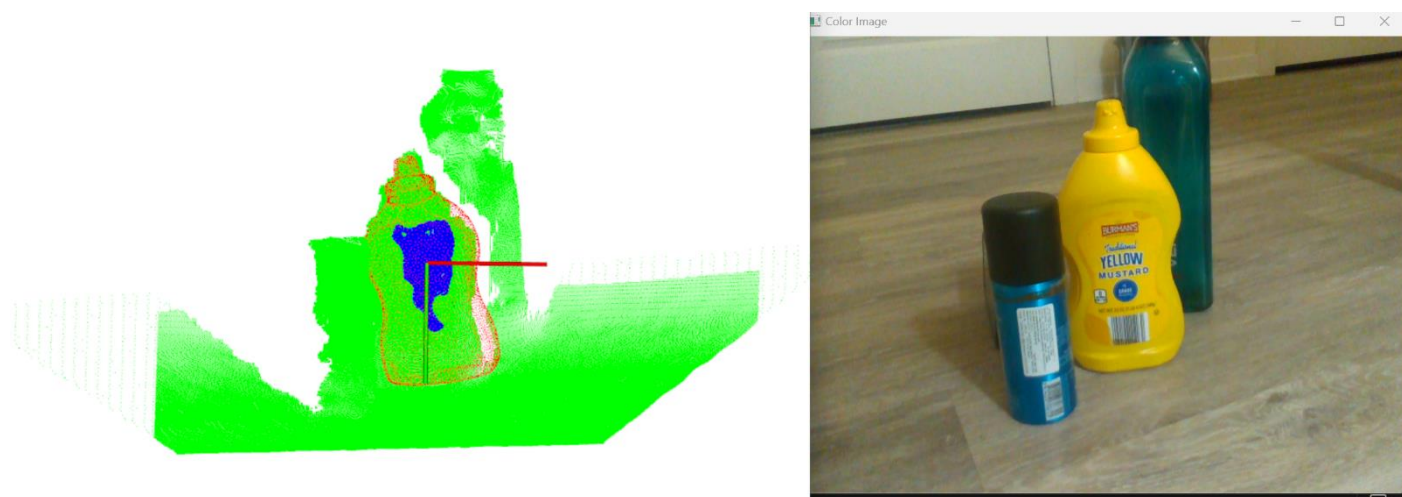
## 4) Visualization and Orientation Estimation:

With the alignment in place, the model point cloud is gracefully transformed back to the camera's coordinate frame. This step bridges the virtual and real, enabling a vivid visualization of the object's orientation within the real-world milieu.

**Coordinate Frame Assignment:**
At the origin, a coordinate frame, symbolic of the object's orientation, is affixed to the model point cloud. When visualized in the camera's coordinate frame, this offers a lucid depiction of the object's orientation shifts relative to its pristine position. Furthermore, the centroid of the transformed model point cloud is pinpointed, and the coordinate frame is anchored to this centroid, offering a holistic and clear orientation reference within the scene.
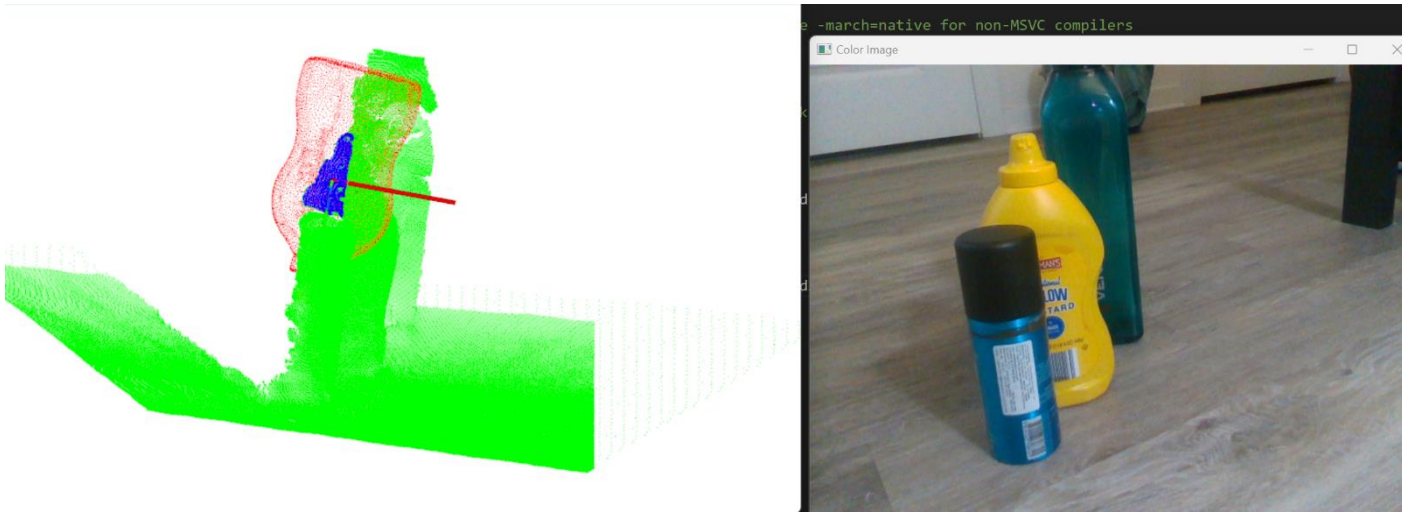
In summation, this methodology, with its intricate blend of robustness and precision, not only elevates object recognition in cluttered terrains but also furnishes a detailed orientation blueprint, setting the stage for nuanced interactions or in-depth analyses.



**Figure 22:** Establishing Object Orientation Reference through Coordinate Frame Assignment

## Limitations and Challenges:

In the case of the second method utilizing VFH + ICP, many of the limitations and challenges overlap with those of the first method. Additionally, a new limitation arises concerning VFH's capability to achieve precise global alignment. One notable constraint is that for VFH to effectively and globally align two point clouds—namely, the source and model point clouds—the source cloud should remain visible and not occlude more than 40 percent of the original object. Otherwise, the alignment process may falter due to VFH's sensitivity to occlusions. This limitation can impact the method's robustness and applicability in scenarios with substantial occlusions.



**Figure 23:** Limitation Due to Occlusion: Object Orientation Incorrect (Occluded > 40%)

## Future Work:

For future endeavors, there are promising avenues to address these limitations. The implementation of alternative versions of ICP holds potential, especially ones designed to handle occlusions effectively. Techniques like trimmed ICP or Go-ICP could prove beneficial in mitigating the impact of occluded regions during alignment, enhancing the method's adaptability to complex scenes.

Furthermore, exploring alternative feature descriptors could enhance alignment performance. Our-CVFH, for instance, offers an improved estimation of both global alignment and feature characteristics. Substituting VFH with Our-CVFH might bolster the accuracy of the alignment process, particularly when occlusions are a concern.

By incorporating these advancements, we can not only expand the method's applicability to more challenging scenarios but also contribute to the refinement of object recognition and orientation analysis in cluttered environments.