# Exploring Depth Estimation Methods with Mono, Stereo, and RGB-D Cameras

Yousif Elraey

September 16, 2023

**Abstract**

In the realm of computer vision and robotics, understanding the three-dimensional structure of the world is a fundamental task. Estimating depth from cameras and obtaining a 3D view of the environment has a multitude of applications, ranging from autonomous navigation for robots to augmented reality experiences for humans. To achieve this, various methods have been developed, each leveraging different camera setups and principles. In this article, we will explore different approaches to estimating depth using mono cameras, stereo cameras, and RGB-D cameras.

## 1 Mono Camera: Monocular Vision

A monocular or mono camera is a single-lens camera, similar to what you find in most smartphones. Depth estimation with a mono camera is a challenging task as it relies on deriving depth information from a single 2D image. Despite the limitations, several techniques have been developed:

### 1.1 Structure from Motion (SfM)

Structure from Motion (SfM) is a classical method used to estimate 3D structure from a sequence of 2D images taken from different viewpoints. By tracking the movement of distinct features across frames and considering camera calibration parameters, SfM can estimate the depth of these features in the scene. However, this method is typically suited for sparse point clouds.

One key advantage of SfM is its ability to work with historical imagery. Archaeologists, for example, have used SfM to reconstruct 3D models of ancient ruins based on old photographs.

### 1.2 Monocular Depth Estimation Neural Networks

Advancements in deep learning have revolutionized monocular depth estimation. Convolutional Neural Networks (CNNs), particularly architectures like

Monodepth and SfMLearner, have demonstrated impressive results in predicting depth maps from single images. These networks are trained on large datasets with corresponding depth information to learn the relationship between image features and depth.

The advantage of neural networks is their ability to capture complex relationships between image features and depth cues. They can handle a wide range of scenes and objects, making them suitable for applications such as autonomous driving, where real-time depth estimation is critical for obstacle detection.

# 2  Stereo Camera: Binocular Vision

Stereo cameras use two synchronized cameras placed at a known baseline distance to capture images from slightly different perspectives. This configuration enables depth estimation through triangulation:

## 2.1  Stereo Correspondence Matching

In stereo vision, the process involves matching corresponding points in the left and right images. By calculating the disparity (horizontal shift) between matched points and knowing the camera baseline, you can compute depth using the simple formula:

$$depth = \frac{baseline \times focal\_length}{disparity}$$

. Stereo vision is highly accurate and suitable for various applications, including obstacle detection and 3D reconstruction.

One of the significant advantages of stereo vision is its ability to handle occlusions. When an object is partially obstructed in one of the stereo images, the matching algorithm can still find corresponding points in the unobstructed image, allowing for accurate depth estimation.

## 2.2  Stereo Neural Networks

Similar to monocular depth estimation, deep learning has made its mark in stereo vision. Stereo neural networks, such as PSMNet and GC-Net, utilize convolutional neural networks to learn stereo correspondences directly from stereo image pairs. These networks can handle complex scenes and occlusions more effectively than traditional algorithms.

Stereo neural networks are particularly useful in applications where real-time performance is essential, such as robotics and augmented reality. They can provide dense and accurate depth maps in milliseconds, enabling fast decision-making.

# 3 RGB-D Camera: Depth as a Sensor Modality

RGB-D cameras, like the popular Microsoft Kinect, combine a standard RGB camera with a depth sensor (usually based on Time-of-Flight or structured light technology). These cameras provide direct depth information for each pixel in the image:

## 3.1 Depth Sensing Technology

RGB-D cameras emit infrared light patterns or use time-of-flight techniques to measure the distance to objects in the scene. The depth information is often represented as a depth map, where each pixel encodes the distance from the camera. This approach is incredibly accurate and reliable for applications like gesture recognition, 3D scanning, and robotics.

The high accuracy and reliability of depth sensing technology make RGB-D cameras the preferred choice for applications where precision is paramount. In fields like medical imaging, where submillimeter accuracy is required, RGB-D cameras play a crucial role in diagnosis and treatment planning.

## 3.2 Simultaneous Localization and Mapping (SLAM)

RGB-D cameras are frequently employed in Simultaneous Localization and Mapping (SLAM) algorithms, where the camera's pose (position and orientation) is estimated simultaneously with the 3D map of the environment. Systems like Google's Project Tango and Apple's ARKit leverage RGB-D sensors to create augmented reality experiences and indoor navigation.

SLAM with RGB-D cameras is used in a wide range of applications, from indoor drone navigation to interactive gaming. The combination of RGB data and depth information allows for more robust tracking of the camera's position and the creation of detailed 3D maps.

# 4 Choosing the Right Method

The choice between mono, stereo, or RGB-D cameras for depth estimation depends on the specific application and the trade-offs between accuracy, cost, and complexity. Monocular vision is more affordable but less accurate, while stereo and RGB-D cameras offer higher accuracy but come with increased hardware requirements.

In recent years, a hybrid approach combining the strengths of these camera types has gained popularity. For example, integrating a mono camera with a depth sensor can provide accurate depth information even in challenging environments.

When choosing the right method, consider factors such as:

- Accuracy requirements: Is high precision necessary, or is a rough estimation sufficient?

- Real-time processing: Does the application require fast depth estimation for timely decision-making?

- Cost constraints: What is the budget for the camera setup and associated hardware?

- Environmental conditions: Will the camera be used in well-controlled conditions, or will it encounter challenging lighting or weather?

In conclusion, estimating depth and obtaining a 3D view of the world is a critical task in computer vision and robotics. The choice of camera setup and depth estimation method depends on the desired application and specific constraints. With the advent of deep learning and advancements in sensor technologies, we are witnessing rapid progress in the field, enabling more sophisticated and versatile 3D understanding of our surroundings.