



Intelligence Academy

**Mathematical Foundations of Image Processing and
Computer Vision : Chapter 02: What Is Computer
Vision**

Mejbah Ahammad
Intelligence Academy

What Is Computer Vision?

Introduction to Computer Vision

Computer Vision is a rapidly advancing field within artificial intelligence (AI) that enables machines to interpret and make decisions based on visual data. It integrates mathematical models, algorithms, and data-driven techniques to simulate the human visual system. This section introduces the foundational motivation and historical context of the discipline.

What is Vision in Humans vs. Machines?

Human vision is a highly evolved biological system capable of understanding spatial and temporal patterns, recognizing objects, and inferring meaning from complex visual environments. It involves photoreceptor cells (rods and cones) converting light into neural signals that are processed by the brain for perception.

In contrast, machine vision—or computer vision—emulates this process us-

Introduction to Computer Vision

ing digital sensors, computational algorithms, and artificial intelligence. While humans effortlessly recognize faces or navigate through crowded scenes, machines must rely on mathematical models and pattern recognition to replicate such behavior.

- **Human Vision:** Biological, intuitive, context-aware, and adaptive.
- **Machine Vision:** Algorithmic, data-driven, and dependent on computational rules and training data.

Modern computer vision systems strive to bridge this gap by incorporating deep learning, geometric modeling, and statistical inference to mimic cognitive capabilities.

Historical Evolution of Computer Vision

Introduction to Computer Vision

The field of computer vision has undergone significant evolution since its inception:

- **1960s–1970s:** Early work focused on symbolic reasoning and manual feature extraction. Image understanding systems aimed to recognize simple geometric shapes.
- **1980s:** The rise of edge detection (e.g., Canny), region growing, and Hough transforms marked a shift toward mathematical modeling of visual data.
- **1990s:** Advances in machine learning, optical flow, and 3D reconstruction techniques brought increased realism and applicability in robotics and industrial systems.
- **2000s:** The development of local feature descriptors like SIFT and SURF,

Introduction to Computer Vision

combined with bag-of-words models, empowered object recognition tasks.

- **2010s–present:** The deep learning revolution introduced Convolutional Neural Networks (CNNs), enabling end-to-end learning from raw pixels. Architectures like AlexNet, ResNet, YOLO, and Vision Transformers have dramatically increased performance in vision benchmarks.

Formal Definition

Computer Vision is a subfield of artificial intelligence that focuses on enabling computers to interpret and process visual data similarly to how humans perceive images. It combines techniques from image processing, machine learning, and pattern recognition to extract high-level understanding from digital images or videos.

While often used interchangeably with image processing, Computer Vision has broader objectives that extend beyond basic image enhancement or transformation.

Computer Vision vs. Image Processing

Though interconnected, image processing and computer vision have distinct purposes:

- **Image Processing** focuses on transforming or enhancing images. It deals

Formal Definition

primarily with low-level operations such as noise removal, contrast adjustment, filtering, and compression.

- **Computer Vision** aims at interpreting and understanding the content of an image. It performs high-level tasks such as object detection, scene understanding, face recognition, and activity tracking.

In summary:

Formal Definition

Aspect	Image Processing	Computer Vision
Objective	Enhancement or transformation of images	Interpretation and understanding of visual data
Level of Abstraction	Low-level operations	High-level reasoning
Examples	Denoising, filtering, contrast adjustment	Object detection, facial recognition, motion tracking
Output	Processed image or extracted features	Semantic labels, actions, or decisions

Definitions from Academia and Industry

Formal Definition

Several definitions help formalize the scope of Computer Vision:

- **Academic Definition:** “Computer Vision is the science of enabling machines to perceive, interpret, and make decisions based on visual inputs such as images and videos.” – *Szeliski, Computer Vision: Algorithms and Applications*
- **Industry Definition (Microsoft):** “Computer Vision is a field of AI that trains computers to interpret and understand the visual world.”
- **IEEE Definition:** “A discipline that studies how computers can gain high-level understanding from digital images or videos, and automate tasks that the human visual system can do.”

These definitions emphasize the goal of computer vision: to build intelligent

Formal Definition

systems that can interpret visual environments with minimal human intervention.

Components of a Computer Vision System

A complete computer vision system consists of several interconnected components that work in a pipeline to capture, analyze, and interpret visual data. Each component plays a specific role in transforming raw images into structured, actionable outputs.

Image Acquisition

Image acquisition is the initial step of any computer vision pipeline. It involves capturing visual data using various sensing technologies:

- **Sensors:** Cameras (RGB, IR), LiDAR, depth sensors, stereo rigs.
- **Formats:** Still images, video streams, or real-time feeds.
- **Challenges:** Illumination variation, motion blur, occlusions.

Components of a Computer Vision System

The quality of image acquisition directly affects the performance of downstream vision algorithms.

Feature Detection

Feature detection refers to identifying key structures or elements within an image that are useful for analysis. These features serve as input to higher-level recognition and learning algorithms.

- **Types:** Edges, corners, blobs, textures, contours.
- **Techniques:** Harris corner detector, SIFT, SURF, ORB.
- **Applications:** Matching, tracking, 3D reconstruction.

Robust feature detection is crucial for object tracking, pose estimation, and image matching tasks.

Components of a Computer Vision System

Pattern Recognition

Pattern recognition enables the classification or identification of visual patterns based on extracted features. It may be rule-based or learned from data.

- **Classical Methods:** K-NN, SVM, decision trees.
- **Deep Learning:** Convolutional Neural Networks (CNNs), ResNet, EfficientNet.
- **Tasks:** Object recognition, digit classification, face detection.

Effective pattern recognition bridges the gap between low-level features and high-level understanding.

Decision Making and Output

Components of a Computer Vision System

The final component involves interpreting the recognized patterns and making decisions based on the application's goals.

- **Output Types:** Bounding boxes, class labels, trajectories, segmentation masks.
- **Systems:** Human-computer interfaces, autonomous vehicles, medical diagnostics.
- **Integration:** Often combined with control systems, robotics, or user feedback loops.

Fundamental Tasks in Computer Vision

Computer vision encompasses a wide range of tasks aimed at extracting various forms of information from visual data. These tasks range from simple classification to complex spatial and temporal understanding. This section highlights the most prominent and widely used tasks in modern computer vision systems.

Image Classification

Image classification involves assigning a single label to an entire image based on the dominant object or scene it contains.

- **Input:** An image.
- **Output:** A class label (e.g., “cat”, “car”, “tree”).

Fundamental Tasks in Computer Vision

- **Techniques:** Convolutional Neural Networks (CNNs), Vision Transformers (ViTs).
- **Use Cases:** Medical image diagnosis, document classification, content moderation.

Object Detection

Object detection extends classification by identifying and localizing multiple objects within an image.

- **Input:** An image.
- **Output:** Class labels with bounding box coordinates.
- **Popular Models:** YOLO, SSD, Faster R-CNN.

Fundamental Tasks in Computer Vision

- **Applications:** Autonomous driving, surveillance, retail analytics.

Semantic and Instance Segmentation

Segmentation refers to assigning a label to each pixel in an image.

- **Semantic Segmentation:** Classifies each pixel into a predefined category (e.g., road, sky, person), without distinguishing between object instances.
- **Instance Segmentation:** Differentiates between individual objects of the same class (e.g., multiple persons).
- **Models:** U-Net, DeepLab, Mask R-CNN.
- **Use Cases:** Medical imaging, autonomous navigation, robotics.

Pose Estimation

Fundamental Tasks in Computer Vision

Pose estimation refers to detecting the spatial configuration of an object or human body in an image.

- **Types:** 2D pose (joint positions in image plane), 3D pose (coordinates in real-world space).
- **Applications:** Human-computer interaction, sports analytics, animation.
- **Popular Libraries:** OpenPose, MediaPipe.

Optical Flow and Motion Tracking

These tasks deal with the temporal aspect of vision, analyzing changes between frames.

- **Optical Flow:** Computes the apparent motion of pixels between consecutive frames.

Fundamental Tasks in Computer Vision

- **Motion Tracking:** Follows objects over time to build trajectories.
- **Applications:** Video stabilization, action recognition, drone navigation.

3D Vision and Depth Estimation

3D vision allows systems to understand spatial relationships and depth in the environment.

- **Goals:** Estimate distance, reconstruct 3D surfaces, infer scene geometry.
- **Techniques:** Stereo vision, LiDAR fusion, monocular depth prediction.
- **Applications:** AR/VR, robotics, 3D modeling.

Mathematical and Algorithmic Foundations

The power of computer vision stems from its strong mathematical and algorithmic foundations. These mathematical tools allow computers to process and understand visual data with precision and efficiency. This section outlines the most essential mathematical techniques and algorithms used in modern vision systems.

Linear Algebra and Geometry in Vision

Linear algebra and geometry are the backbone of image transformations and 3D modeling.

- **Matrices and Vectors:** Images are represented as matrices; pixel operations are performed using matrix algebra.
- **Transformations:** Translation, rotation, scaling, and affine transformations are modeled using transformation matrices.

Mathematical and Algorithmic Foundations

- **Camera Geometry:** Concepts such as projection matrices, epipolar geometry, and homographies enable 3D reconstruction from 2D images.
- **Applications:** Stereo vision, structure-from-motion (SfM), visual odometry.

Convolution and Filtering

Convolution is the fundamental operation in both classical image processing and deep learning-based computer vision.

- **Spatial Filtering:** Applies kernels to enhance or detect image features such as edges or textures (e.g., Sobel, Gaussian, Laplacian).
- **Frequency Domain Filtering:** Uses Fourier transforms for denoising and compression.

Mathematical and Algorithmic Foundations

- **Convolutional Layers:** In CNNs, convolutional kernels learn patterns during training to extract hierarchical features.

Filtering enables both handcrafted feature extraction and automatic learning from visual data.

Feature Descriptors (SIFT, HOG, SURF)

Feature descriptors are compact representations of key visual patterns, enabling matching, recognition, and tracking.

- **SIFT (Scale-Invariant Feature Transform):** Detects and describes local features invariant to scale and rotation.
- **HOG (Histogram of Oriented Gradients):** Captures edge orientations, widely used for pedestrian detection.

Mathematical and Algorithmic Foundations

- **SURF (Speeded-Up Robust Features):** Faster alternative to SIFT with robust performance in real-time settings.
- **Applications:** Object matching, registration, panorama stitching, motion estimation.

Machine Learning & Deep Learning Integration

Modern computer vision systems increasingly rely on data-driven models for recognition and prediction.

- **Classical ML:** Algorithms such as KNN, SVM, and decision trees are used for structured prediction and classification.
- **Deep Learning:** Models like CNNs, RNNs, and Vision Transformers automatically learn features from large datasets.

Mathematical and Algorithmic Foundations

- **End-to-End Training:** Input-to-output mappings are learned directly from annotated data, reducing dependence on manual feature engineering.
- **Frameworks:** TensorFlow, PyTorch, OpenCV, Keras.

