University of Cambridge Engineering Part IB Paper 8 Information Engineering

Introduction to Computer Vision



Roberto Cipolla May 2025

What is computer vision?

Vision is about discovering from images what is present in the scene and where it is. It is our most powerful sense.



In **computer vision** a camera is linked to a computer. The computer automatically processes and interprets the images of a real scene to obtain useful information (**3R's**: recognition, registration and reconstruction) and **representations** for decision making and **action** (e.g. for navigation, manipulation or communication).

Why study computer vision?

- 1. Intellectual curiosity how do we see?
- 2. Replicate human vision to allow a machine to see many industrial, commercial and healthcare applications.

Computer Vision is not:

- **Image processing:** image enhancement, image restoration, image compression. Take an image and process it to produce a new image which is, in some way, more desirable.
- **Pattern recognition:** classifies patterns into one of a finite set of prototypes. There is an Infinite variation in images of objects and scenes due to changes in viewpoint, lighting, occlusion and clutter.

Applications

- Autonomous vehicles and self-driving cars
- Industrial and agricultural automation
 - Visual inspection
 - Object recognition.
 - Robot hand-eye coordination
- Human-computer interaction
 - Face detection and recognition.
 - Gesture-based and touch free interactions
 - Cashierless transactions
 - Image search in video and image databases
- Augmented reality and enhanced interactions
 - AR with mobile phones and we arable computers
- Surveillance and Security
- Medical Imaging
 - Detection, segmentation and classification
- \bullet 3D modelling, measurement and visualisation
 - 3D model building and photogrammetry
 - Human body and motion capture
 - 3D Virtual fitting and e-commerce
 - Avatar creation and talking heads

Applications

Examples of recent computer vision research that has led to new products and services.



- Microsoft Kinect Human pose detection and tracking for game interface using gestures
- Microsoft Hololens Smart glasses for Augmented Reality
- Orcam Wearable camera using text-recognition to help visually-impaired
- Wayve and Waymo autonomous driving using cameras
- Dogtooth Technologies addressing labour shortages in fruit picking with robotics
- Synthesia and Toshiba Europe Photorealistic 3D avatars and Talking Heads
- Metail and Trya Virtual fitting of clothes and shoes by estimating shape from images
- Amazon Prime Air Drone delivery services with visual localisation and navigation
- Boston Dynamics and Unitree Vision for robot navigation and hand-eye co-ordination
- Infrastructure visual inspection

How to study vision? The eye





- Retina measures about 1000 mm^2 and contains about 10^8 sampling elements (rods) (and about 10^6 cones for sampling colour).
- The eye's spatial resolution is about 0.01° over a 150° field of view (not evenly spaced, there is a fovea and a peripheral region).
- Intensity resolution is about 11 bits/element, spectral resolution is about 2 bits/element (400–700 nm).
- Temporal resolution is about 100 ms (10 Hz).
- Two eyes (each about 2cm in diameter), separated by about 6cm.
- A large chunk of our brain is dedicated to processing the signals from our eyes a data rate of about 3 GBytes/s!

Why not copy the biology?

- There is no point copying the eye and brain human vision involves over 60 billion neurons.
- Evolution took its course under a set of constraints that are very different from today's technological barriers.
- The computers we have available cannot perform like the human brain.
- We need to understand the underlying principles rather than the particular implementation.

Compare with flight. Attempts to duplicate the flight of birds failed.



<u>The camera</u>



- A typical digital SLR CCD measures about 24×16 mm and contains about 6×10^6 sampling elements (pixels).
- Intensity resolution is about 8 bits/pixel for each colour channel (RGB).
- Most computer vision applications work with monochrome images.
- Temporal resolution is about 40 ms (25 Hz)
- One camera gives a raw data rate of about 400 MBytes/s.

The CCD camera is an adequate sensor for computer vision.

Image formation



Image formation is a many-to-one mapping. The image encodes nothing about the *depth* of the objects in the scene. It only tells us along *which* ray a feature lies, not *how far* along the ray. The inverse imaging problem (inferring the scene from a single image) has no unique solution.

Ambiguities in the imaging process



Two examples showing that image formation is a many-toone mapping. The Ames room and two images of the same 3D structure.



Vision as information processing

David Marr, one of the pioneers of computer vision, said: "One cannot understand what seeing is and how it works unless one understands the underlying information processing tasks being solved."

From an information processing point of view we must convert the huge amount of unstructured data in images into useful and actionable representations:

images	\rightarrow	generic salient features
100 MBytes/s		100 KBytes/s
(mono CCD)		
salient features	\rightarrow	representations and actions
100 KBytes/s		1-10 bits/s

Vision resolves the ambiguities inherent in the imaging proces by drawing on a set of constraints (AI). But where do the constraints come from? We have the following options:

- 1. Use more than one image of the scene.
- 2. Make assumptions about the world in the scene.
- 3. Learn (supervised and unsupervised) from the real world.

Feature extraction

The first stages of most computer vision algorithms perform feature extraction. The aim is to reduce the data content of the images while preserving the useful information they contain.



The most commonly used features are edges, which are detected as discontinuities in the image. This involves filtering (by convolution) and differentiating the image. Automatic edge detection algorithms produce something resembling a noisy line drawing of the scene.



Corner detection is also common. Corner features are localised in 2D and are particularly useful for finding correspondences in motion analysis using correlation.

Feature descriptors which are invariant to scale, orientation and lighting (e.g. SIFT) facilitate matching over arbitrary viewpoints and in different lighting.

Perspective Projection

To interpret the image (using the features extracted from the image), we have to understand how the image was formed. In other words, we have to develop a **camera model**.



Camera models must account for the position of the camera, perspective projection and CCD imaging. These geometric transformations have been well-understood since the C14th. They are best described within the framework of **projective geometry**.

Projection and Camera models



With a camera model, we can predict how known objects will appear in an image and can also recover their position and orientation (pose) in the scene.



Cluttered scene



Spanner pose recovered

Stereo vision

Having two cameras allows us to triangulate on features in the left and right images to obtain depth. It is even possible to infer useful information about the scene when the cameras are not **calibrated**.



Stereo vision requires that features in the left and right image be matched. This is known as the **correspondence problem**.





Structure from motion

Related to stereo vision is a technique known as **structure from motion**. Instead of collecting two images simultaneously, we allow a single camera to move and collect a sequence of images from different viewpoints.



As the camera moves, the motion of some features (in this case corner features) is **tracked**.

The trajectories allow us to recover the 3D translation and rotation of the camera and the 3D structure of the scene.



Geometrical framework

The first approach is to focus on generic computer vision techniques which make minimal assumptions about the outside world. This means concentrating on the theory of perspective, stereo vision and structure from motion.

We typically use a geometric framework:

- 1. Reduce the information content of the images to a manageable size by extracting salient features, typically **edges** or **blobs**. (These features are generic and substantially *invariant* to a variety of lighting conditions.)
- 2. Model the imaging process, usually as a perspective projection and express using projective transformations.
- 3. Invert the transformation using as many images and constraints as necessary to extract 3D structure and motion.

Statistical framework

Geometry alone is only a part of the solution. We will also need techniques which learn from the visual world. They are part of a statistical framework to understanding vision and for building systems which:

- 1. Have the ability to test hypotheses
- 2. Deal with the ambiguity of the visual world
- 3. Are able to fuse information
- 4. Have the ability to learn

Many of these requirements can be addressed by reasoning with probabilities and are the subject of other advanced courses Machine Learning.

Deep Learning for Computer Vision

Modern computer vision uses Deep Learning architectures based on **Convolutional Neural Networks** (CNNs).

CNNs have multiple layers of feature responses which are obtained by filtering/convolutions and non-linear activation functions. The weights of each filter are learned from training examples and deep networks will typically have millions (and even billions!) of parameters.



CNNs have been shown to be very effective at learning a hierarchy of features and representations for computer vision tasks. In particular they are used in many **recognition** tasks including text and face recognition, object detection and semantic segmentation.

These architectures (and the simple algorithms to train them) were first introduced in the 1980's. It is only in the last-decade that they have achieved state-of-the art performance on computer vision tasks. This is due to the availability of very large amounts of labelled training data; deeper networks and specialised computing hardware (GPUs) that can speed up the training algorithms (based on stochastic gradient descent optimisation) by many orders of magnitude.

Syllabus

1. Introduction

- Computer vision: what is it, why study it and how?
- Vision as an information processing task
- Perspective projection, stereo and structure from motion
- Geometrical, statistical frameworks and machine learning frameworks for vision

2. Detecting and matching image features

- Edge detection, corner detection, blob detection
- Convolution with gaussians and derivatives of gaussians to provide bandpass filters.
- Edge detection using directional filters.
- Scale-space and image pyramids for feature detection
- The SIFT feature descriptor for matching image features.

3. Demonstrations

• State-of-the-art object detection, semantic segmentation and localisation systems

Further reading

Books

D.A. Forsyth and J. Ponce. *Computer Vision - A Modern Approach*. Prentice Hall 2003.

D. Marr. Vision: a computational investigation into the human representation and processing of visual information. Freeman, 1982.

R. Szeliski. Computer Vision: algorithms and applications. Springer, 2011.