

Early History: Initial Forays

Outline

- 1 Early History: Initial Forays
- 2 Towards Algorithms and Practice: Low-level Understanding
- 3 Towards Algorithms and Practice: Next Level of Understanding
- 4 The Deep Learning Era

Vincent N.B. (MIT-H)

5.1.2 History



Let us move on to the topic of this lecture. As most of you may know images are formed when a light source hits the surface of an object and light is reflected and some of that light is reflected onto an image plane which is then captured through optics on to a sensor plane. So, that is the overall information and the factors that affect the image formation are the light source strength and direction, the surface geometry, material of the surface such as its texture as well as other nearby surfaces that, whose light could get reflected onto the surface, the sensor capture properties we will talk more about that as we go and the image representation and color space itself. We will talk about some of these as we go. So, to study all of these one would probably need to study this from geometrical perspective, where you study 2D transformations, 3D transformations, camera calibration, distortion. So, we would not cover all of these but cover a few relevant topics from these in this particular lecture. If you are interested in a more detailed coverage of these topics please read chapters 1 to 5 of the book by Forsyth and Ponce. Starting with how light gets reflected off a surface the more typical morals of reflection state that when light hits a surface there are 3 simple reactions possible, there are more than 3 but 3 simple reactions to start with. Firstly, some light is absorbed and that depends on a factor called albedo ρ and typically when you have a surface with low albedo more light gets absorbed. Some light is reflected diffusively.

Early History¹

1959 1963 1966 1971'73 1979-82

- David Hubel and Torsten Wiesel publish their work "Receptive fields of single neurons in the cat's striate cortex"
- Placed electrodes into primary visual cortex area of an anesthetized cat's brain
- Showed that simple and complex neurons exist, and that visual processing starts with simple structures such as oriented edges

¹Credit: Rostyslav Demush, medium.com

Vincent N.B. (MIT-H)


5.1.2 History



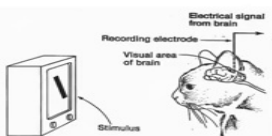
It scatters in multiple directions, so that happens independent of the viewing angle.

Early History: Initial Forays




Early History¹



- David Hubel and Torsten Wiesel publish their work "Receptive fields of single neurons in the cat's striate cortex"
- Placed electrodes into primary visual cortex area of an anesthetized cat's brain
- Showed that simple and complex neurons exist, and that visual processing starts with simple structures such as oriented edges



¹Credit: Rostyslav Demush, medium.com
Vincent N B. (B.T.H) 11.2 History

Example of surfaces where lights scatters diffusively is brick, cloth, rough wood or any other texture material and in this scenario Lambert's cosine law states that the amount of reflected light is proportional to cosine of angle from which you are viewing the reflection. And finally, there are also phenomena such as fluorescence, where the output wavelength could be different from the input wavelength or other phenomena such as phosphorescence. And there are models to evaluate how bright the surface appears. So from a view point of colour itself we all know that visible light is 1 portion of the vast electromagnetic spectrum, so visible light is one small portion of the vast electromagnetic spectrum, so we know that infrared falls on one side, ultraviolet falls on the other side and there are many other forms of light across the electromagnetic spectrum. So, coloured light which arrives at a sensor typically involves two factors, colour of the light source and colour of the surface itself. So, an important development in sensing of colour in cameras is what is known as the Bayer Grid or the Bayer Filter.

Early History: Initial Forays

Early History⁴



- Seymour Papert (with Gerald Sussman) from MIT launched the *Summer Vision Project*
- Aimed to develop a platform to automatically segment background/foreground and extract non-overlapping objects from real-world images



⁴Credit: Rostyslav Demush, medium.com
Vincent N B. (B.T.H) 11.2 History

For accessing this content for free (no charge), visit : npTEL.ac.in





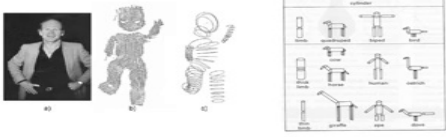
So, not every sensing element in a camera captures all three components of light you may be aware that typically we represent light as RGB at least coloured light as RGB; Red Green and Blue. We will talk a little bit more about other ways of representing coloured light a little later, but this is the typical way of representing coloured light and not every sensing element on the camera captures all three colours instead a person called Bayer proposed this method in a grid manner where you have 50 percent green sensors, 25 percent red sensors and 25 percent blue sensors which is inspired by human visual receptors. And this is how these sensors are checkered, so in a real camera device you would have a sensor array and there is a set of sensors that captures only red light, there is set of sensors that captures the green light, there is set of sensors that captures the blue light and to obtain the full colour image demosaicing algorithms are used where surrounding pixels are used to contribute the value of the exact colour at a given pixel. These are known as demosaicing algorithms. This is not the only kind of colour filter.

Early History: Initial Forays




Early History

1959 1963 1966 1971-73 1979-82

- Object recognition through shape understanding
 - Binford 1971, Generalized Cylinders
 - Marr and Nishihara 1978, Skeletons and Cylinders
- MIT's Artificial Intelligence Lab offers a "Machine Vision" course



Vineeth N.B. (BT-H) 11.2 History


So, this is the general pipeline of image capture. So, let us try to revisit, visit some of these components over the next few minutes. So, first thing is the camera sensor itself so you all must have heard of CCD and CMOS. This is often common decision to be made when you buy a camera these days a lesser issue but earlier days it used to be even more.

Early History: Initial Forays

Early History⁶




1959 1963 1966 1971 1979-82

- David Marr, *Vision: A computational investigation into the human representation and processing of visual information*, 1982
- Established that vision is hierarchical
- Introduced a framework where low-level algorithms that detect edges, curves, corners, etc., are used to get high-level understanding of visual data



⁶Credit: Rostyslav Demush, medium.com

Vineeth N.B. (BT-H) 11.2 History

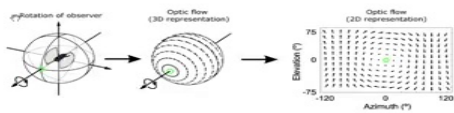
What is the difference? So, the main difference between CCD and CMOS is that in CCD it stands for Charged Coupled Device. So, the many properties that you may actually see when you look at, when you take a picture on a camera. Sampling pitch, which defines a spacing between the sensor cells on the imaging chip. Fill factor or also known as active sensing area size, sorry, which is the ratio of the active sensing area size with respect to the theoretically available sensing area on the sensing element.

Towards Algorithms and Practice: Low-level Understanding




Towards Algorithms and Practice: Low-level Understanding

1981 1986 '87 '88 '89

- Optical Flow:** Horn and Schunck develop method to estimate the direction and speed of a moving object across two images
- Flow is formulated as a global energy functional which is minimized



Vineeth N.B. (BT-H) 11.2 History


Fill factor or also known as active sensing area size, sorry, which is the ratio of the active sensing area size with respect to the theoretically available sensing area on the sensing element. Chip size which is the entire size of area of the chip itself. Analog gain which is the amplification of the sense signal using automatic gain control logic we would not going to the details of each of this once again if you are interested you can read the references provided at the end of this lecture to get more details of all of them. Typically, analog gain is what you control using your ISO setting on your camera, you can also have sensor noise that comes from various sources in the sensing process. Your resolution tells you how many bits is specified for each pixel which is also decided by an analog to digital conversion module in CCD or in case of CMOS in the sensing, in the sensing elements.

Towards Algorithms and Practice: Low-level Understanding




1981 1986 '87 '88 '89

↓

- **Canny Edge Detector:** Multi-stage edge detection operator, with a computational theory of edge detection
- Used calculus of variations to find the function that optimizes a given functional
- Well-defined method, simple to implement, became very popular for edge detection



Vineeth N.B. (BT-H) 11.2 History






Towards Algorithms and Practice: Low-level Understanding




1981 1986 '87 '88 '89

↓

- **Snakes or active contour models** delineate an object outline from a possibly noisy 2D image
- Widely used in applications like object tracking, shape recognition, segmentation, edge detection and stereo matching



Vineeth N.B. (BT-H) 11.2 History

Mirrored cameras such as DSLR also give you a physical shutter mechanism variable focal length and aperture so on and so forth. That is the reason there is value for DSLR cameras despite the advancement in smartphones cameras. So, the other factors that you need to understand when you talk about image formation is the concept of sampling and Aliasing, we will talk about this in more details bit later but a brief review now is Shannon Sampling Theorem states that if the maximum frequency of your data on your image is f_{max} you should at least sample at twice that frequency. Why so, we will see a bit later but for the moment that frequency that you captured it is also called the Nyquist frequency and if you have frequencies about the Nyquist frequency in your image then the phenomenon called Aliasing happens. So, why is this bad and what impact can it have on image formation? This can often create issues when you up sample or down sample an image.



Towards Algorithms and Practice: Next Level Understanding

1991 1997 '98 '99 2001 2005 '06 '07 2009

↓

- **Eigenfaces for face recognition** (Turk & Pentland, 1991)
- **Computational theories of object recognition** (Edelman, 1997)
- **Perceptual grouping, Normalized cuts** (Shi & Malik, 1997)
- **Particle filters, Mean shift for tracking** (Liu & Chen, 1998)(Cheng, 1998)
- **SIFT** (Lowe, 1999) (Lowe, 2004)
- **Viola-Jones face detection** (Viola & Jones, 2001)
- **Conditional Random Fields** (Lafferty et al, 2001)

Vineeth N.B. (BT-H) 11.2 History


Other colour spaces that are used in practice are XYZ, YUV, Lab, YCbCr, HSV so on and so forth. There is actually an organization call the CIE which establishes standards for colour spaces because this is an important, this is actually important for the printing and scanning industry, I think this is extremely important people working in that space. So, that is the reason there are standards establish for these kinds of spaces, we would not get into more details here once again if you are interested please go through these links below to know more about colour spaces what do you mean by additive, subtractive, so on and so forth, please look at these links. Finally, the last stage in image formation is image compression, because you have to store the image that you captured, so

typically you convert the signal into a form called YCbCr where Y is luminance CbCr talks about chrominance what is known as colour factor or the chrominance and the reason for this is that you typically try to compress luminance with a higher fidelity than chrominance.

The Deep Learning Era

2010 2012 '13 '14 '15 '16 '17 '18 '19

• ImageNet arrives

Vineth N B. (BT-H)
5.2 History

Because of the way humans or the human visual system perceives light, luminance is a bit more important than chrominance, so you ensure that luminance is actually compressed with a higher fidelity which means your reconstruction is better for luminance than for chrominance, so that is one reason why YCbCr is used as a popular colour space before storage, once again if you do not understand YCbCr, go back to the previous slide look at all of these links to understand YCbCr is one of the colour space representations that are available in practice. And as I just mentioned so the most common compression technique that used to store an image is called the Discrete Cosine Transform which is popularly used in standard such as MPEG and JPEG Discrete Cosine Transform is actually a variant of Discrete Fourier Transform and it is a you can call it as a reasonable approximation of an eigen decomposition of image patches.

The Deep Learning Era

2010 2012 '13 '14 '15 '16 '17 '18 '19

• ImageNet arrives
• AlexNet wins the ImageNet challenge

Vineth N B. (BT-H)
5.2 History

So, we would not get into in for the time now, videos this is how images are compressed using method call DCT, videos also use what is known as block level motion compensation, so you also divide images into frames and set of frames into block and then you store certain frames based on concepts from motion compensation, this is typically used in the MPEG standard which uses, which divides all frames into what are known as i frames, p frames and b frames and then uses strategies to decide how each frame should be coded, that is how videos are compressed. And compression quality finally is measured through a metric called PSNR, apologies for the typo, it will be fixed before the slides are uploaded, which stands for Peak Signal to Noise Ratio, sorry for these typos. PSNR is defined as $10 \log_{10} \frac{i_{\max}^2}{\text{MSE}}$, where i_{\max} is the maximum intensity and MSE is simply talks about the mean squared error between the original image and the compressed image, how much is the mean squared error pixel wise between these two images. And the numerator talks about the maximum intensity that you can have in an image, so this is typically called as PSNR which is used to measure the quality of image compression, there are other kinds of matrix which are based on human perception but this is the most popular statistical metric that is used.