

COMP 9517 Computer Vision

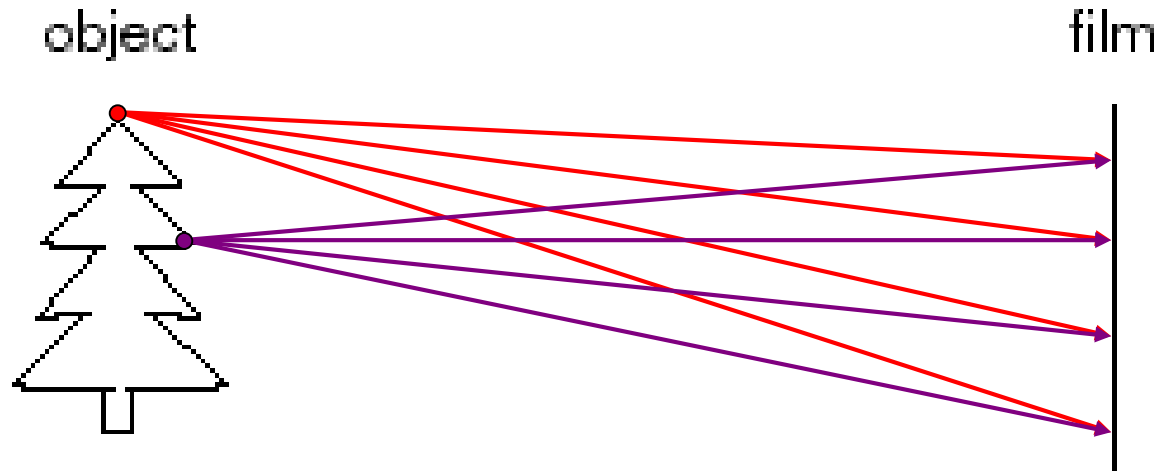
Image Formation

Geometry of Image Formation

Mapping between image and world coordinates

- Pinhole camera model
- Projective geometry
 - Vanishing points and lines
- Projection matrix

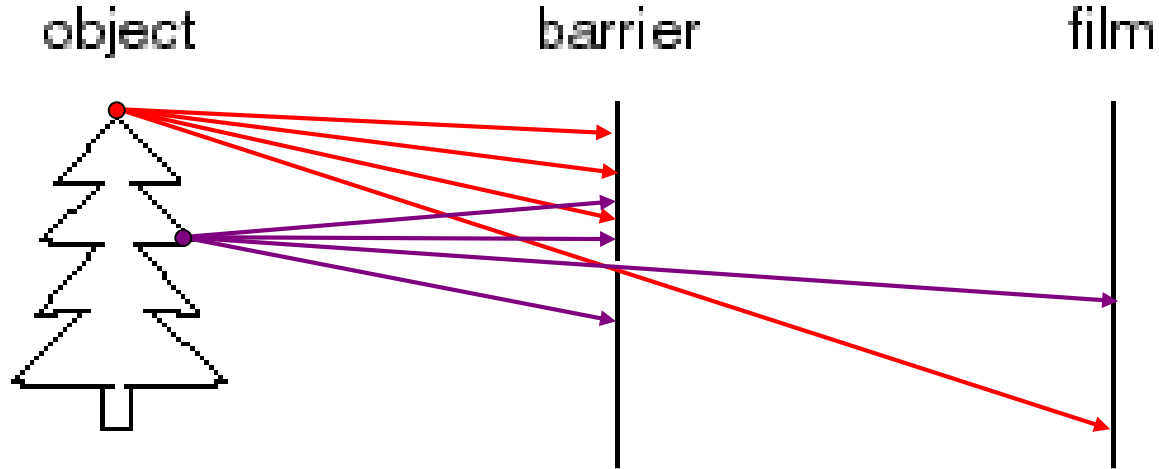
Image formation



Let us design a camera

- Idea 1: put a piece of film in front of an object
- Do we get a reasonable image?

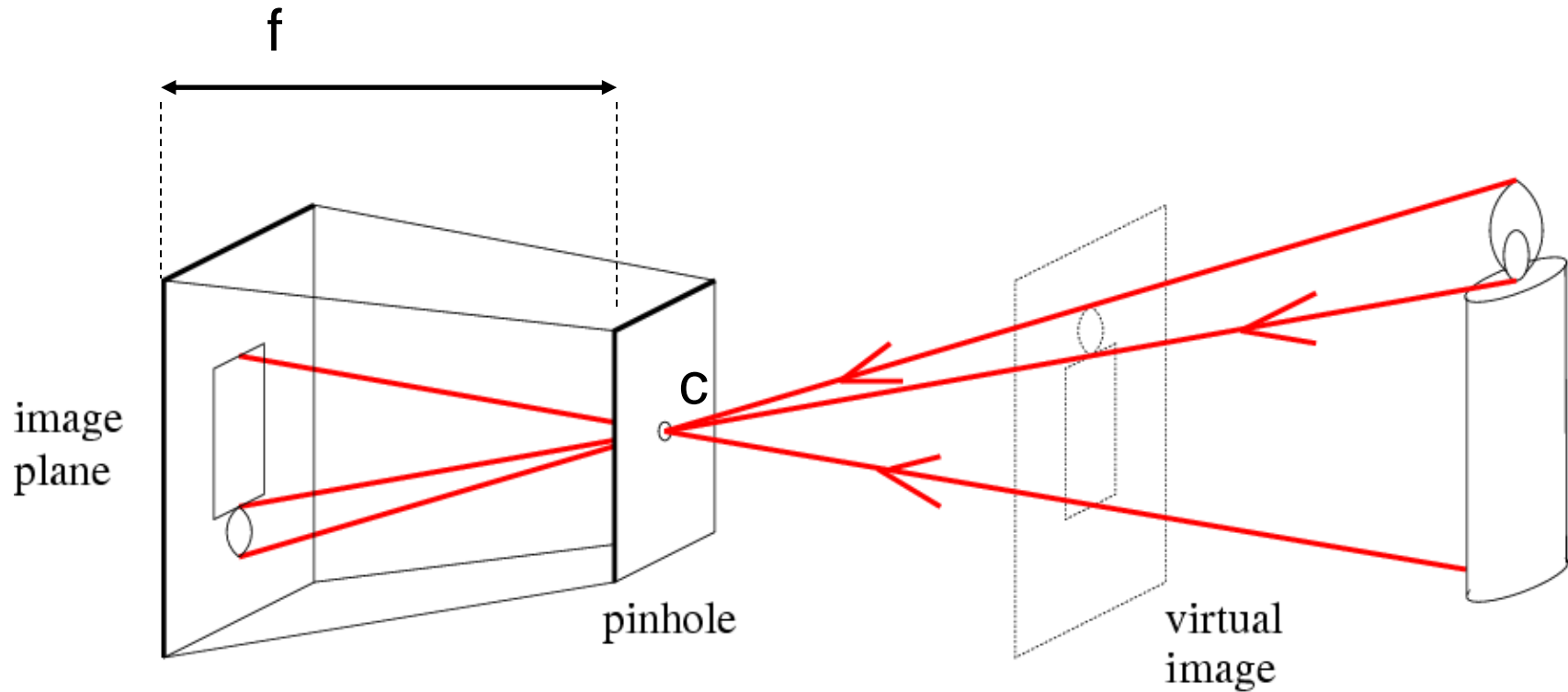
Pinhole camera



Idea 2: add a barrier to block off most of the rays

- This reduces blurring
- The opening known as the **aperture**

Pinhole camera

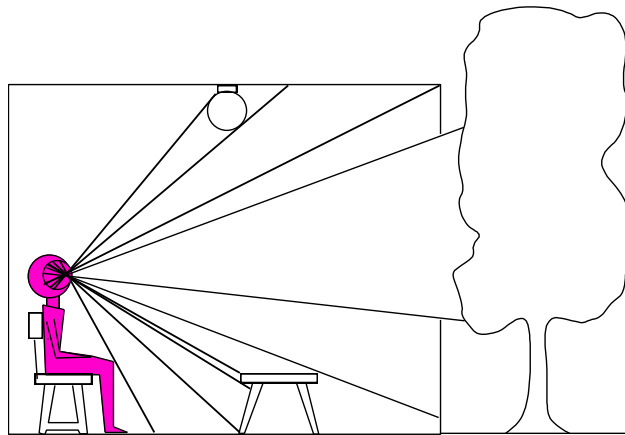


f = focal length

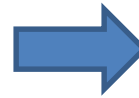
c = centre of the camera

Dimensionality Reduction Machine (3D to 2D)

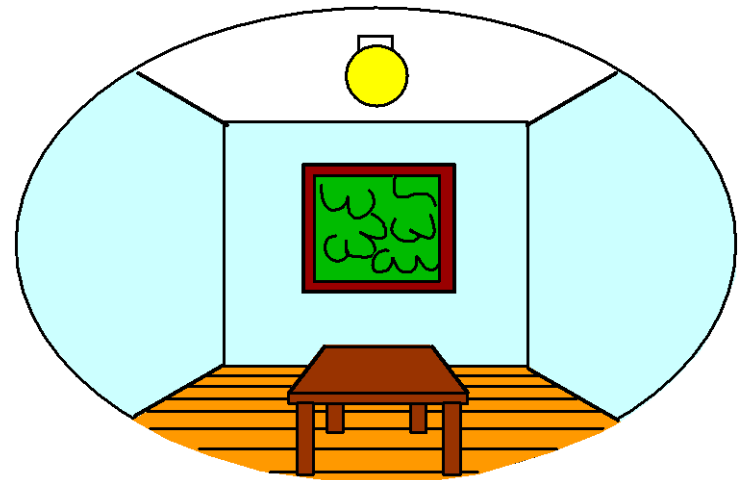
3D world



Point of observation



2D image



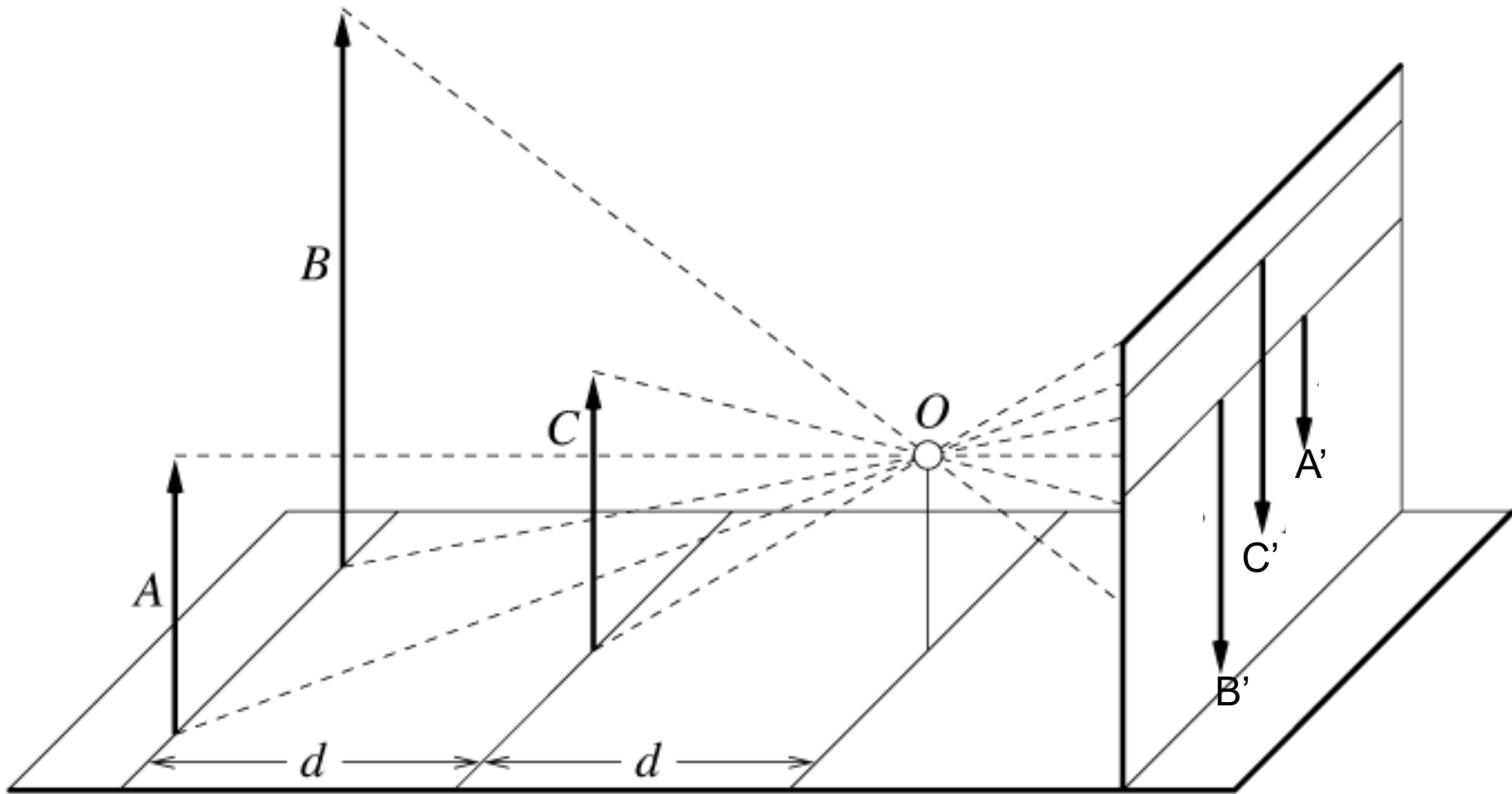
Projection can be tricky...



Projection can be tricky...



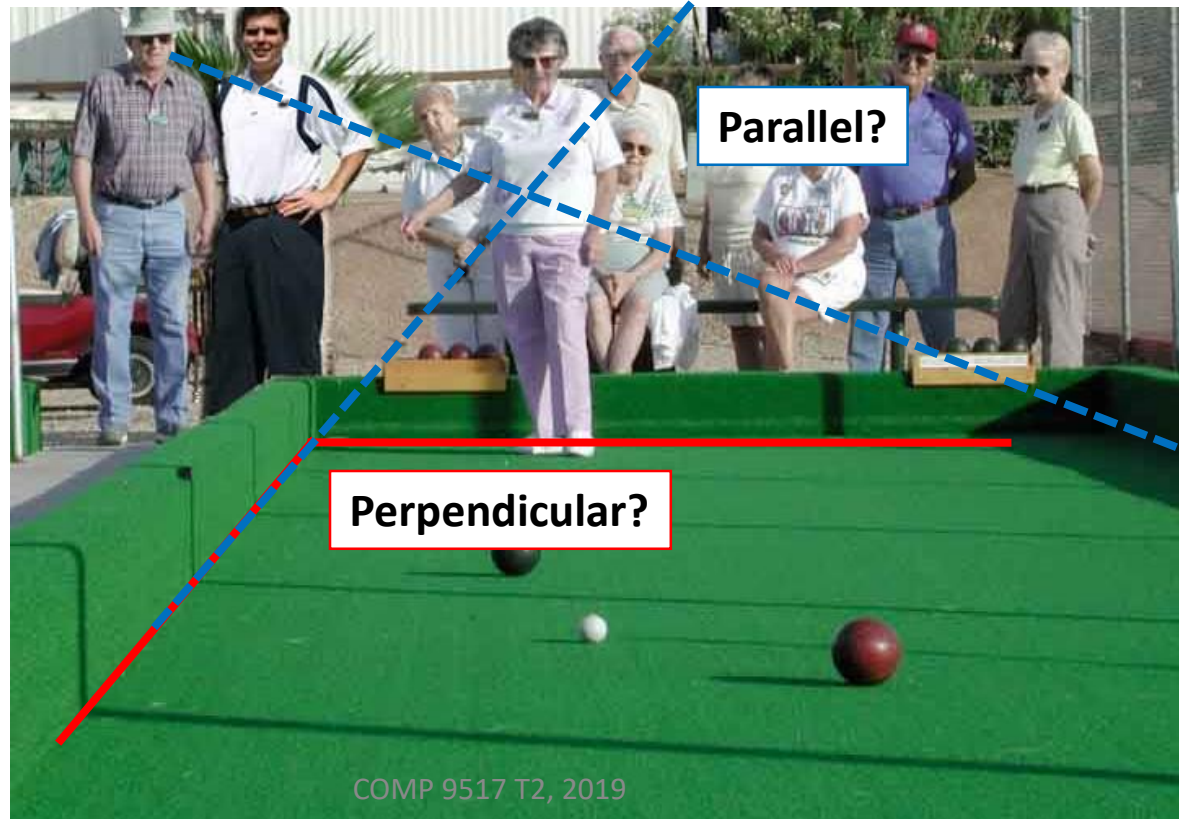
Length and area are not preserved



Projective Geometry

What is lost?

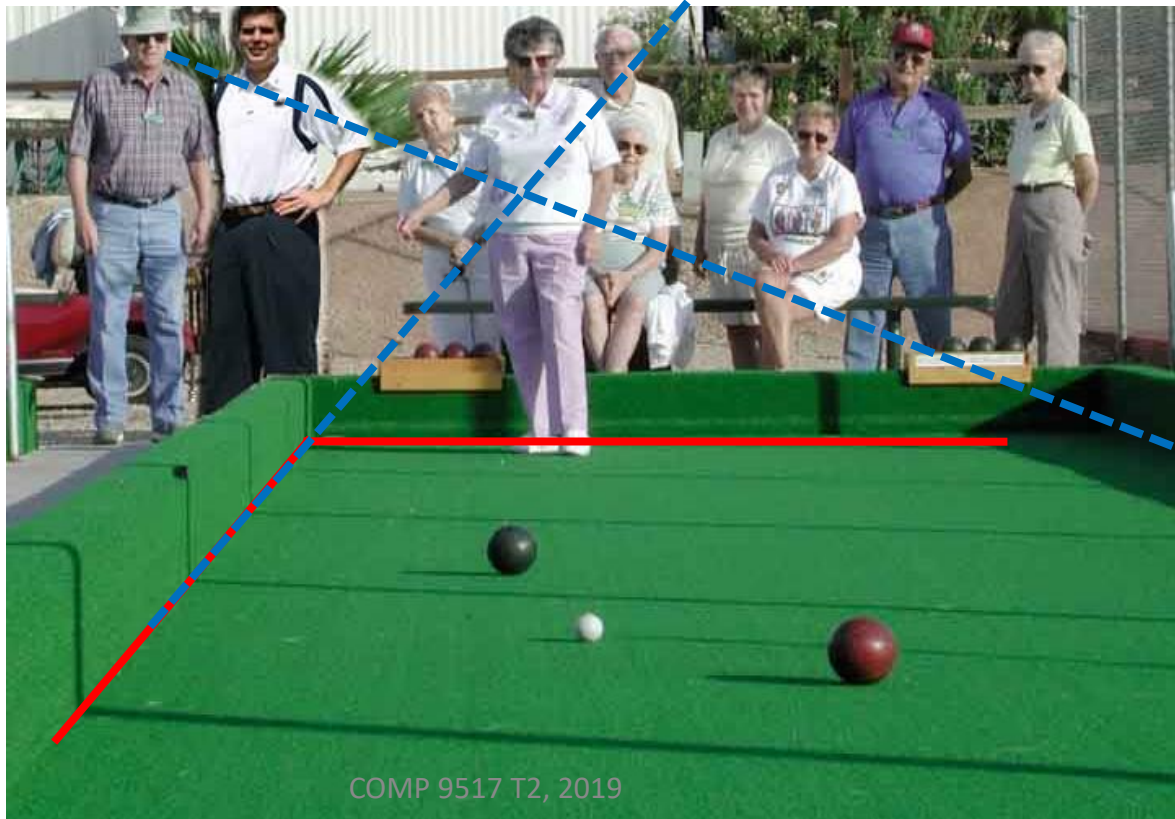
- Length
- Angles



Projective Geometry

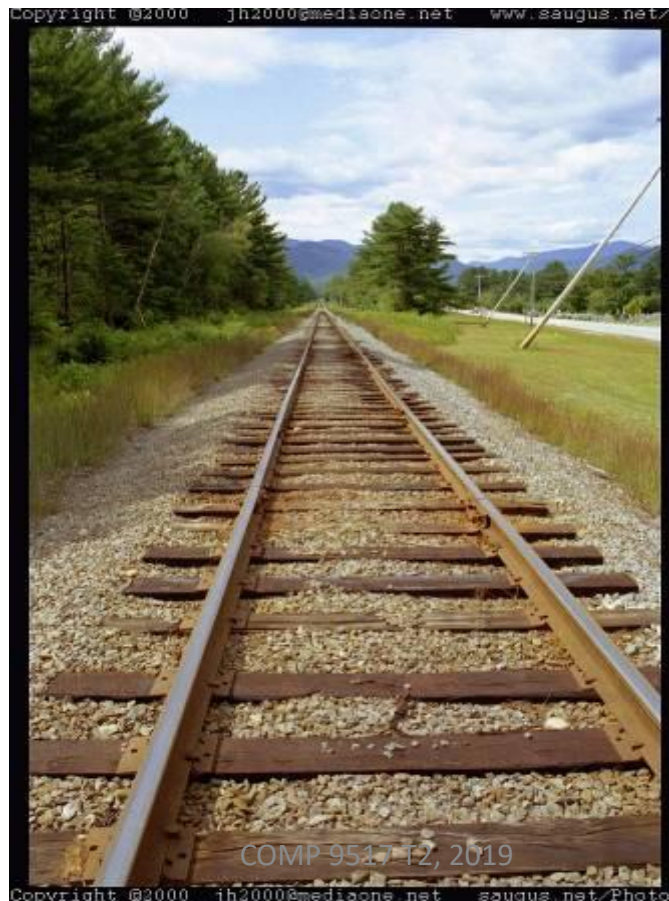
What is preserved?

- Straight lines are still straight

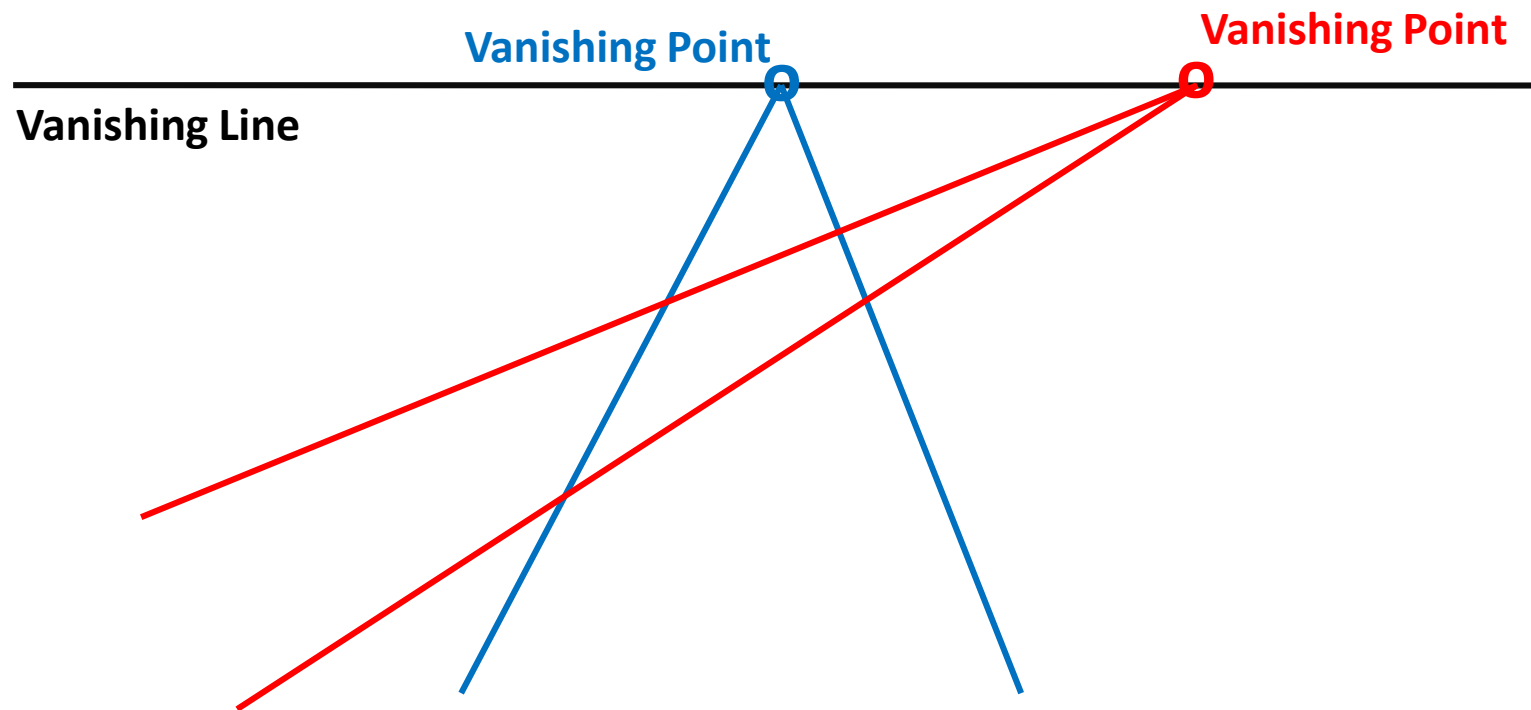


Vanishing points and lines

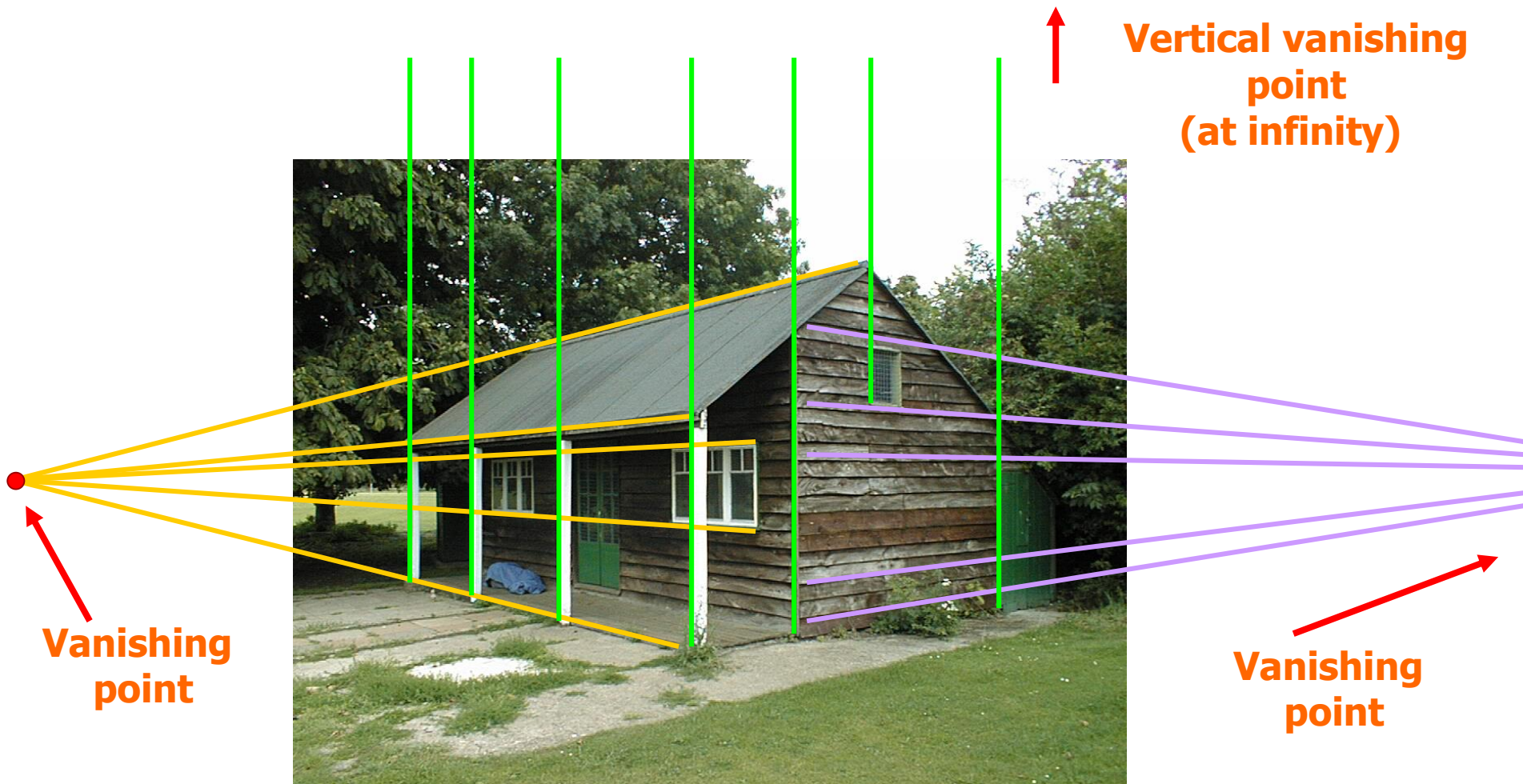
Parallel lines in the world intersect in the image at a “vanishing point”



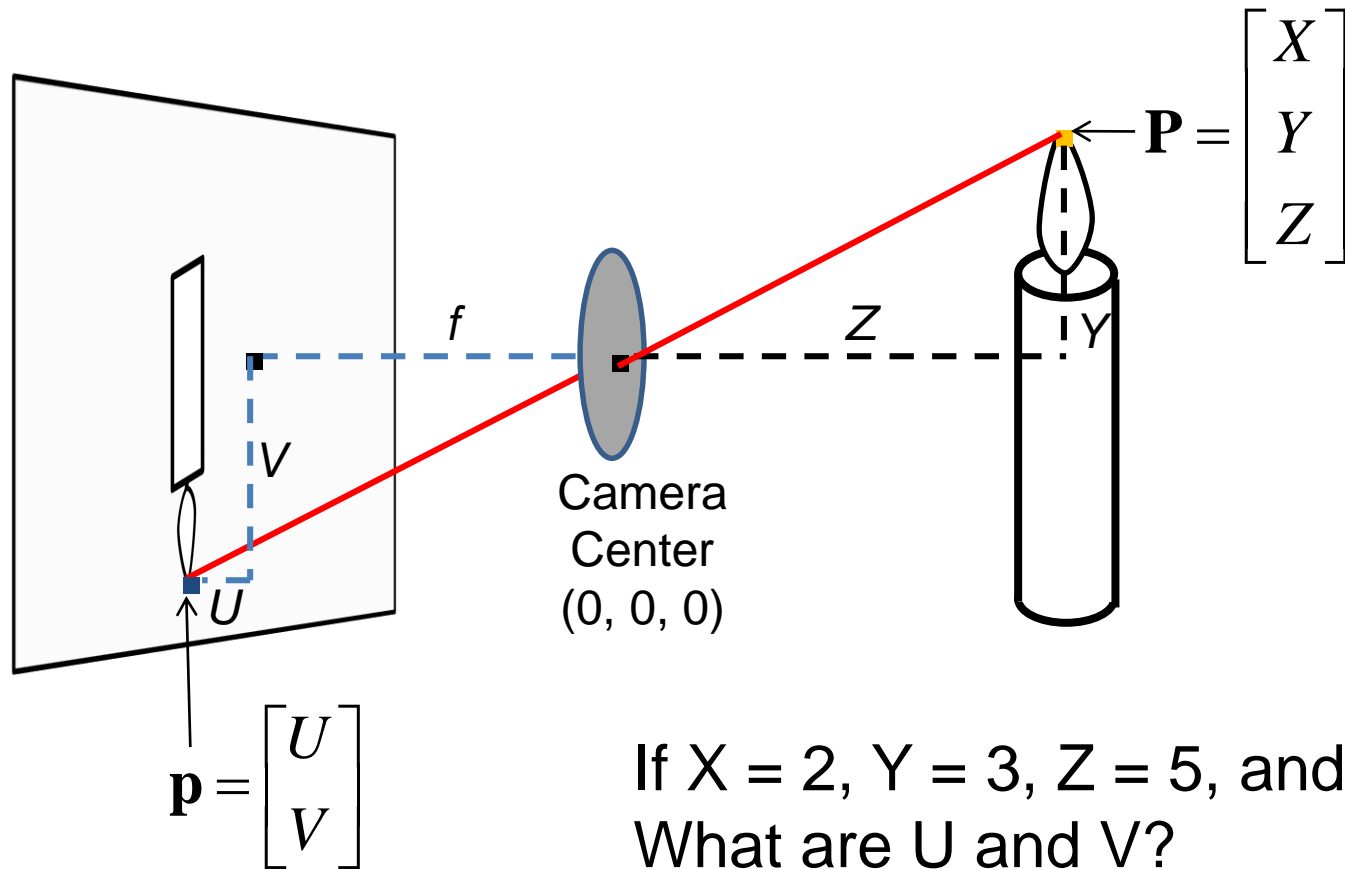
Vanishing points and lines



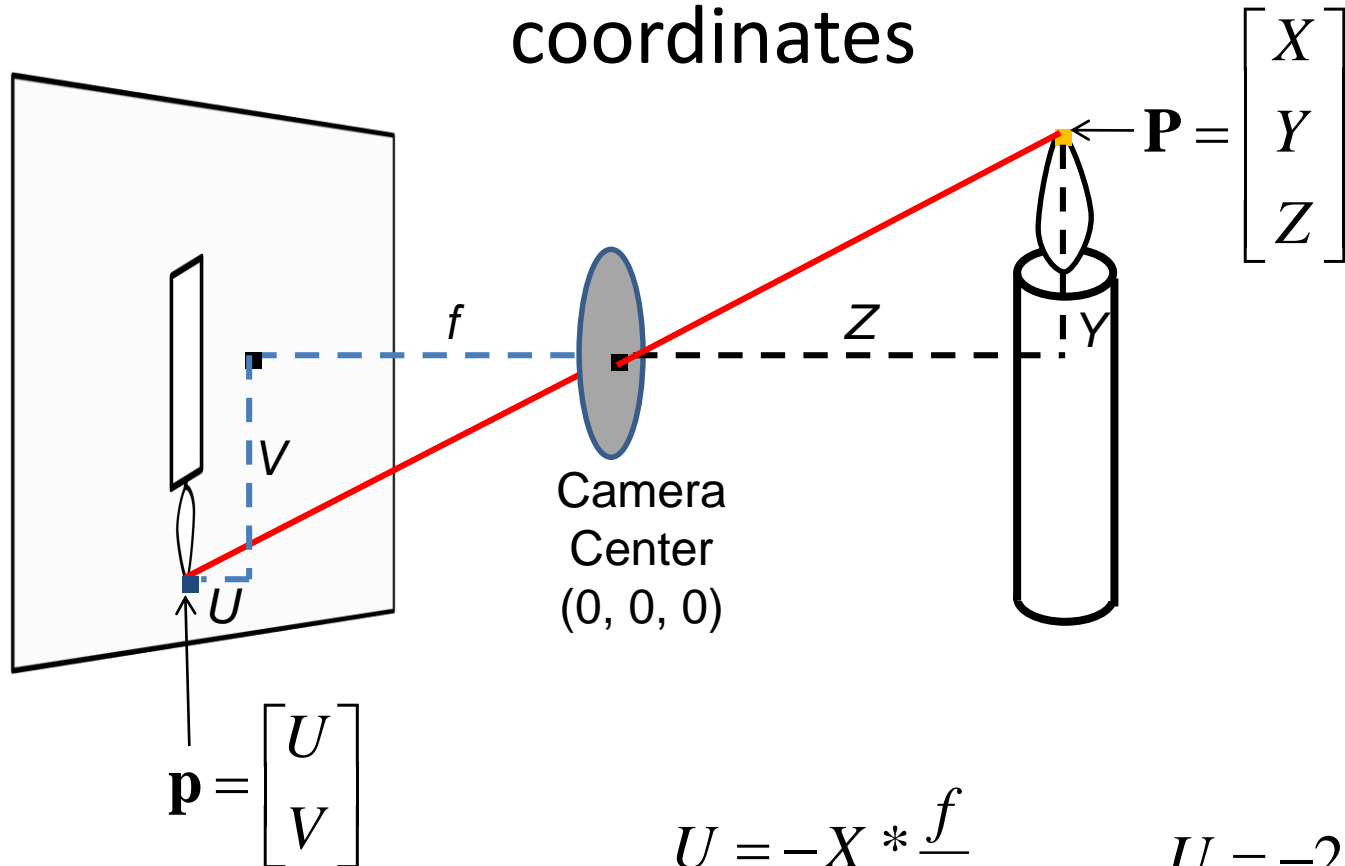
Vanishing points and lines



Projection: world coordinates \rightarrow image coordinates



Projection: world coordinates \rightarrow image coordinates



$$U = -X * \frac{f}{Z}$$

$$U = -2 * \frac{2}{5}$$

$$V = -Y * \frac{f}{Z}$$

$$V = -3 * \frac{2}{5}$$

Perspective Projection

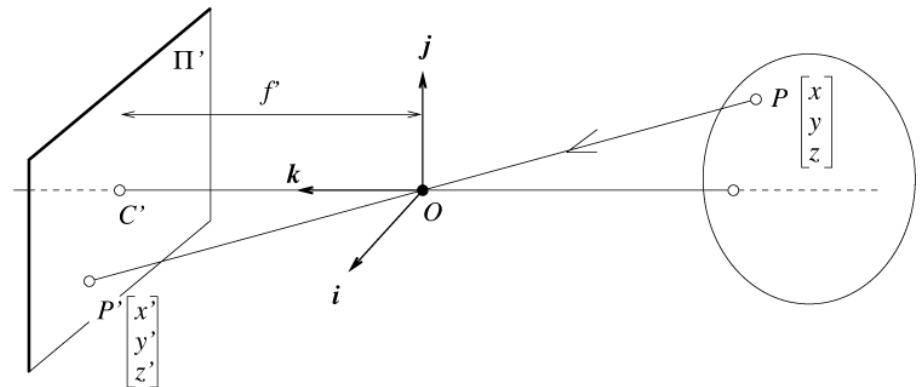
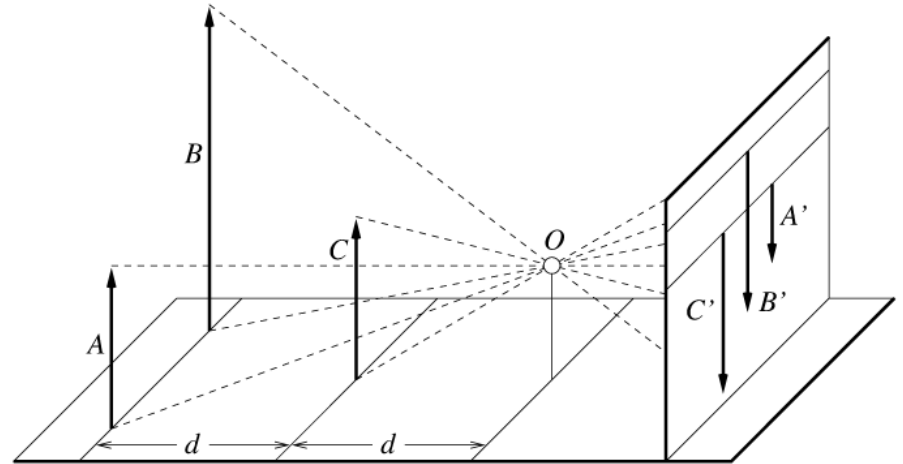
- Apparent size of object depends on its distance: far objects appear smaller

- By similar triangles

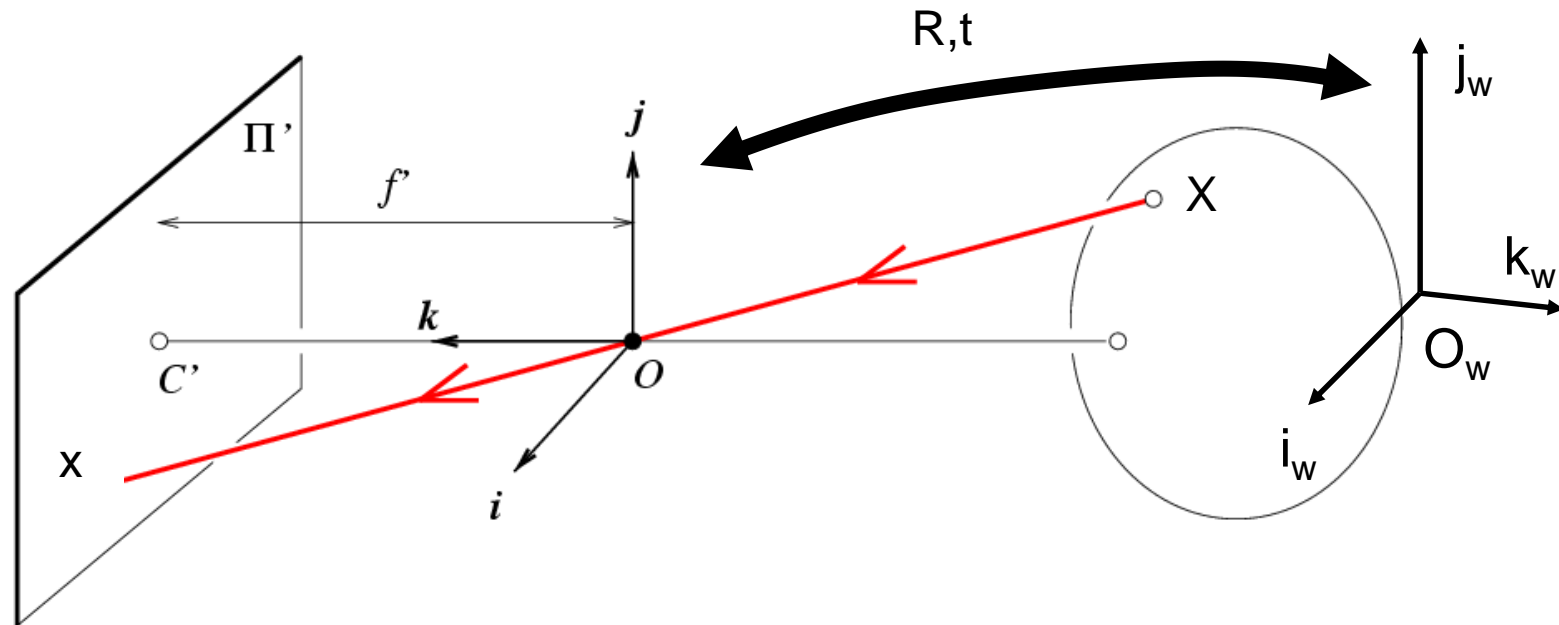
$$(x', y', z') \rightarrow (f \frac{x}{z}, f \frac{y}{z}, -f)$$

- Ignore the third coordinate, and get

$$(x', y') \rightarrow (f \frac{x}{z}, f \frac{y}{z})$$



Projection matrix



$$\mathbf{x} = \mathbf{K}[\mathbf{R} \quad \mathbf{t}] \mathbf{X}$$

\mathbf{x} : Image Coordinates: $(u, v, 1)$

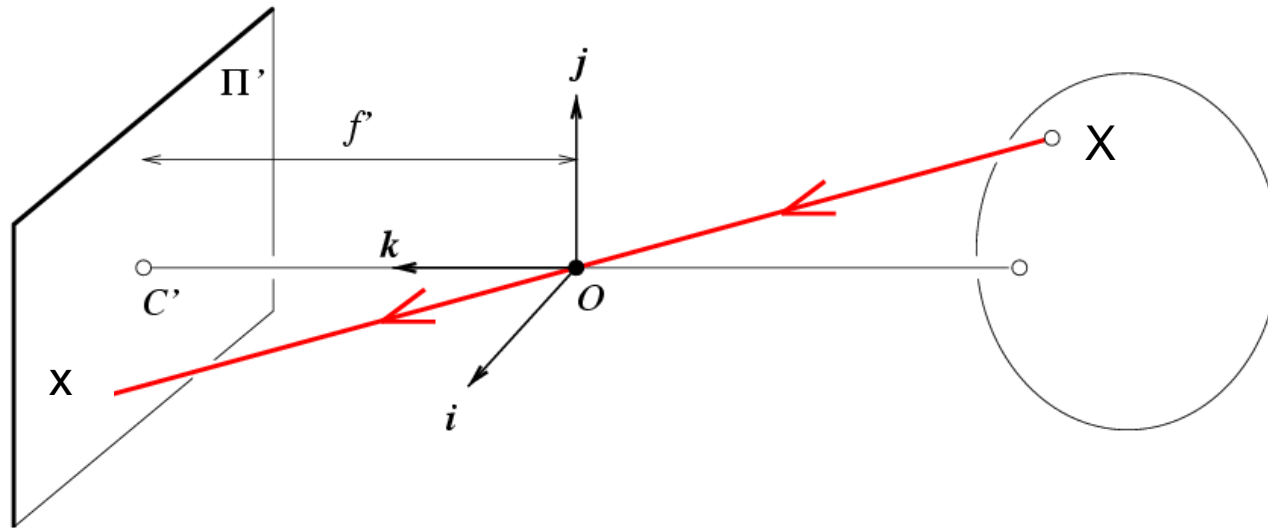
\mathbf{K} : Intrinsic Matrix (3×3)

\mathbf{R} : Rotation (3×3)

\mathbf{t} : Translation (3×1)

\mathbf{X} : World Coordinates: $(X, Y, Z, 1)$

Pinhole Camera Model



Intrinsic Assumptions

- Unit aspect ratio
- Optical center at $(0,0)$
- No skew

Extrinsic Assumptions

- No rotation
- Camera at $(0,0,0)$

$$\mathbf{x} = \mathbf{K} \begin{bmatrix} \mathbf{I} & \mathbf{0} \end{bmatrix} \mathbf{X} \Rightarrow {}^w \begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = \begin{bmatrix} f & 0 & 0 & 0 \\ 0 & f & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} x \\ y \\ z \\ 1 \end{bmatrix}$$

The matrix \mathbf{K} is indicated by a red dashed arrow pointing to the first three rows of the transformation matrix.

Homogeneous coordinates

- Line equation: $ax + by + c = 0$

$$line_i = \begin{bmatrix} a_i \\ b_i \\ c_i \end{bmatrix}$$

- Append 1 to pixel coordinate to get homogeneous coordinate

$$p_i = \begin{bmatrix} u_i \\ v_i \\ 1 \end{bmatrix}$$

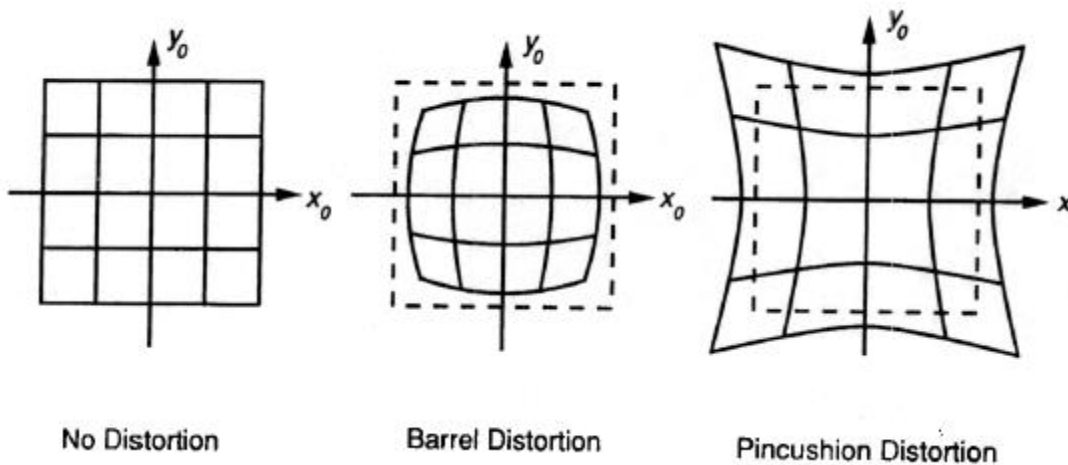
- Line given by cross product of two points

$$line_{ij} = p_i \times p_j$$

- Intersection of two lines given by cross product of the lines

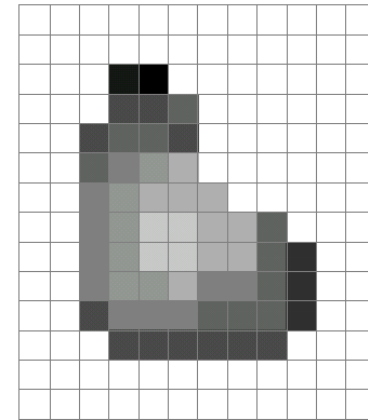
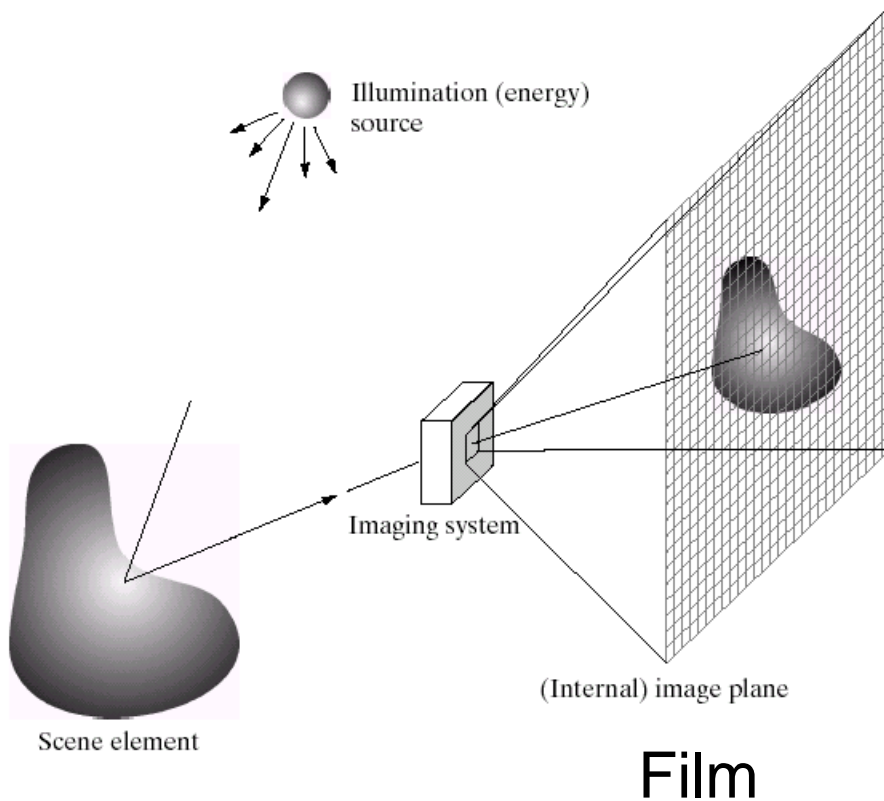
$$q_{ij} = line_i \times line_j$$

Beyond Pinholes: Radial Distortion

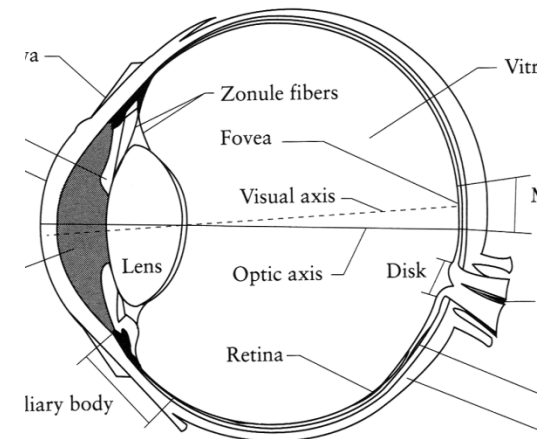


Corrected Barrel Distortion

Image Formation



Digital Camera

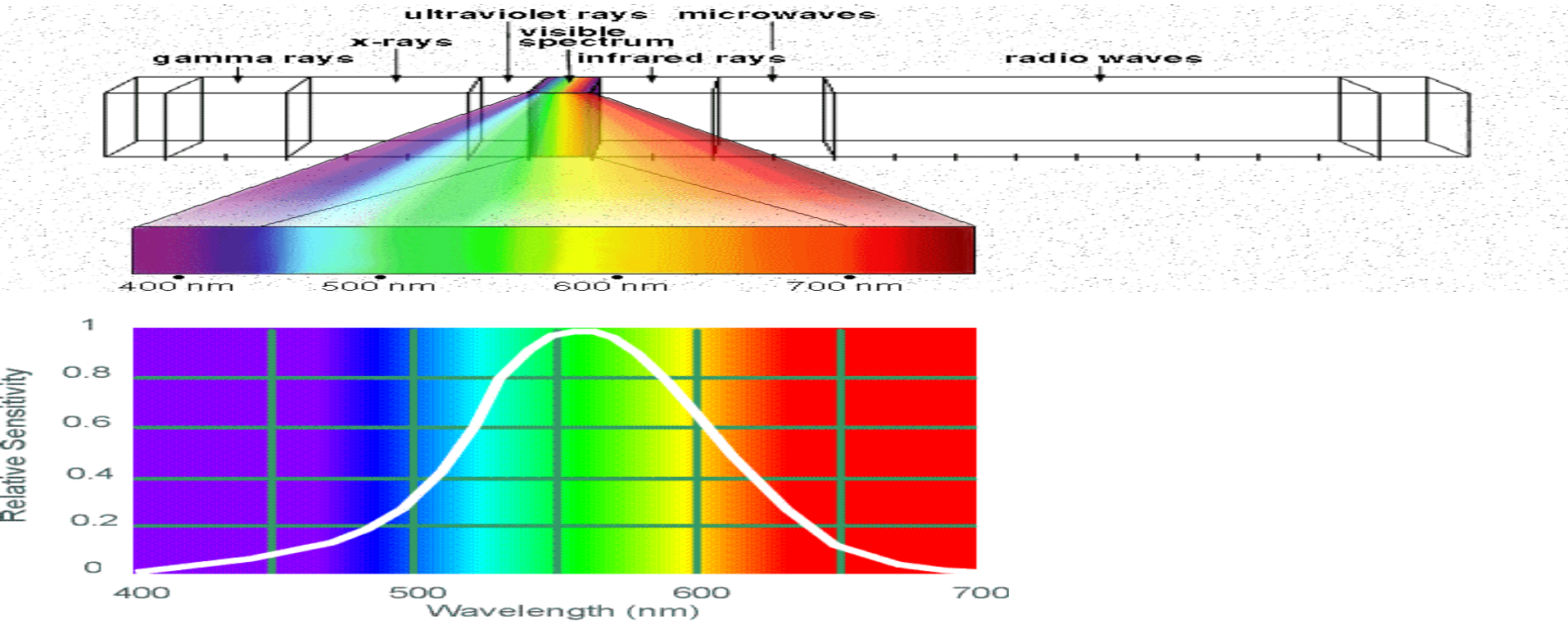


The Eye

Why do we care about human vision?

- Cameras necessarily imitate the frequency response of the human eye, so we should know that much.
- Computer vision probably would not get as much attention if biological vision (especially human vision) had not proved that it was possible to make important judgements from 2D images.

Electromagnetic Spectrum



Human Luminance Sensitivity Function

Colour Image

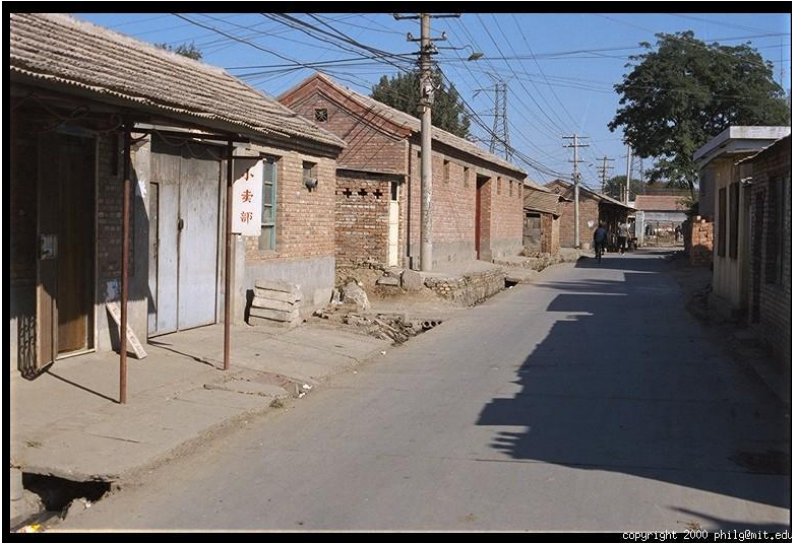
R



G



B



Images represented as a matrix

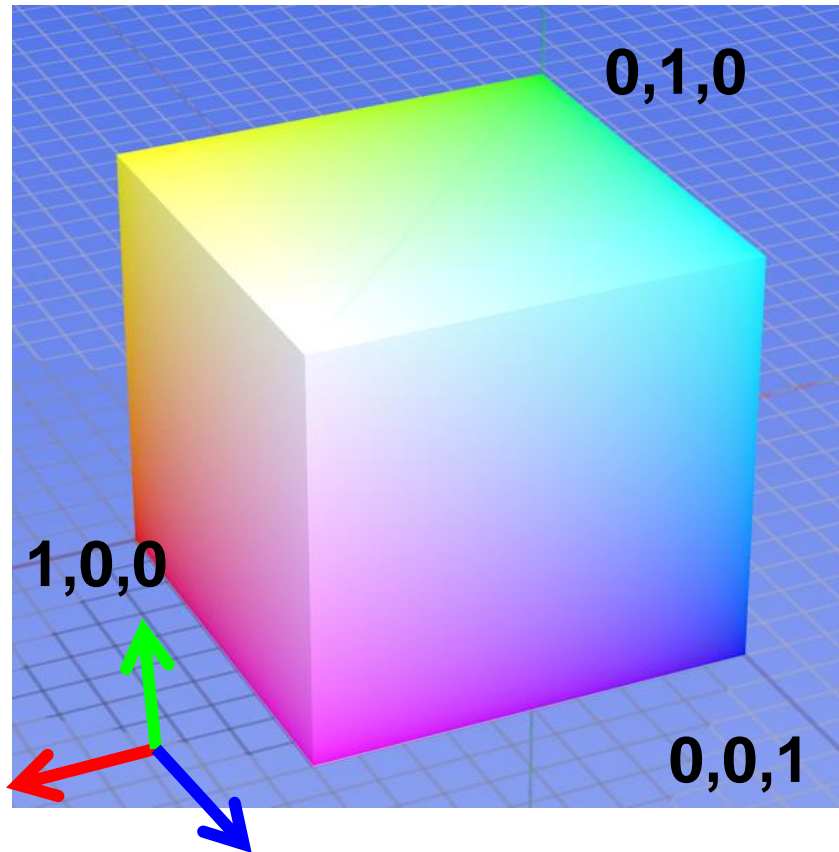
Diagram illustrating an image represented as a matrix, showing the Red (R), Green (G), and Blue (B) color channels. The main matrix is labeled with 'row' (vertical arrow) and 'column' (horizontal arrow).

0.92	0.93	0.94	0.97	0.62	0.37	0.85	0.97	0.93	0.92	0.99	0.92	0.99	0.92	0.99
0.95	0.89	0.82	0.89	0.56	0.31	0.75	0.92	0.81	0.95	0.91	0.95	0.91	0.95	0.91
0.89	0.72	0.51	0.55	0.51	0.42	0.57	0.41	0.49	0.91	0.92	0.91	0.92	0.91	0.92
0.96	0.95	0.88	0.94	0.56	0.46	0.91	0.87	0.90	0.97	0.95	0.97	0.95	0.97	0.95
0.71	0.81	0.81	0.87	0.57	0.37	0.80	0.88	0.89	0.79	0.85	0.79	0.85	0.79	0.85
0.49	0.62	0.60	0.58	0.50	0.60	0.58	0.50	0.61	0.45	0.33	0.45	0.33	0.45	0.33
0.86	0.84	0.74	0.58	0.51	0.39	0.73	0.92	0.91	0.49	0.74	0.49	0.74	0.49	0.74
0.96	0.67	0.54	0.85	0.48	0.37	0.88	0.90	0.94	0.82	0.93	0.82	0.93	0.82	0.93
0.69	0.49	0.56	0.66	0.43	0.42	0.77	0.73	0.71	0.90	0.99	0.90	0.99	0.90	0.99
0.79	0.73	0.90	0.67	0.33	0.61	0.69	0.79	0.73	0.93	0.97	0.93	0.97	0.93	0.97
0.91	0.94	0.89	0.49	0.41	0.78	0.78	0.77	0.89	0.99	0.93	0.99	0.93	0.99	0.93

The diagram shows the Red (R), Green (G), and Blue (B) color channels. The main matrix is labeled with 'row' (vertical arrow) and 'column' (horizontal arrow). The R channel is the first 10 columns, the G channel is the next 2 columns, and the B channel is the final 2 columns. The values are repeated for the G and B channels.

Colour spaces: RGB

Default colour space



Some drawbacks

- Strongly correlated channels
- Non-perceptual



R
(G=0,B=0)



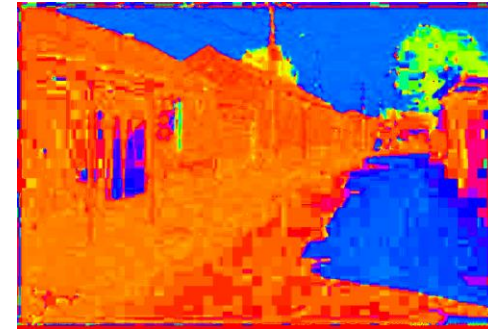
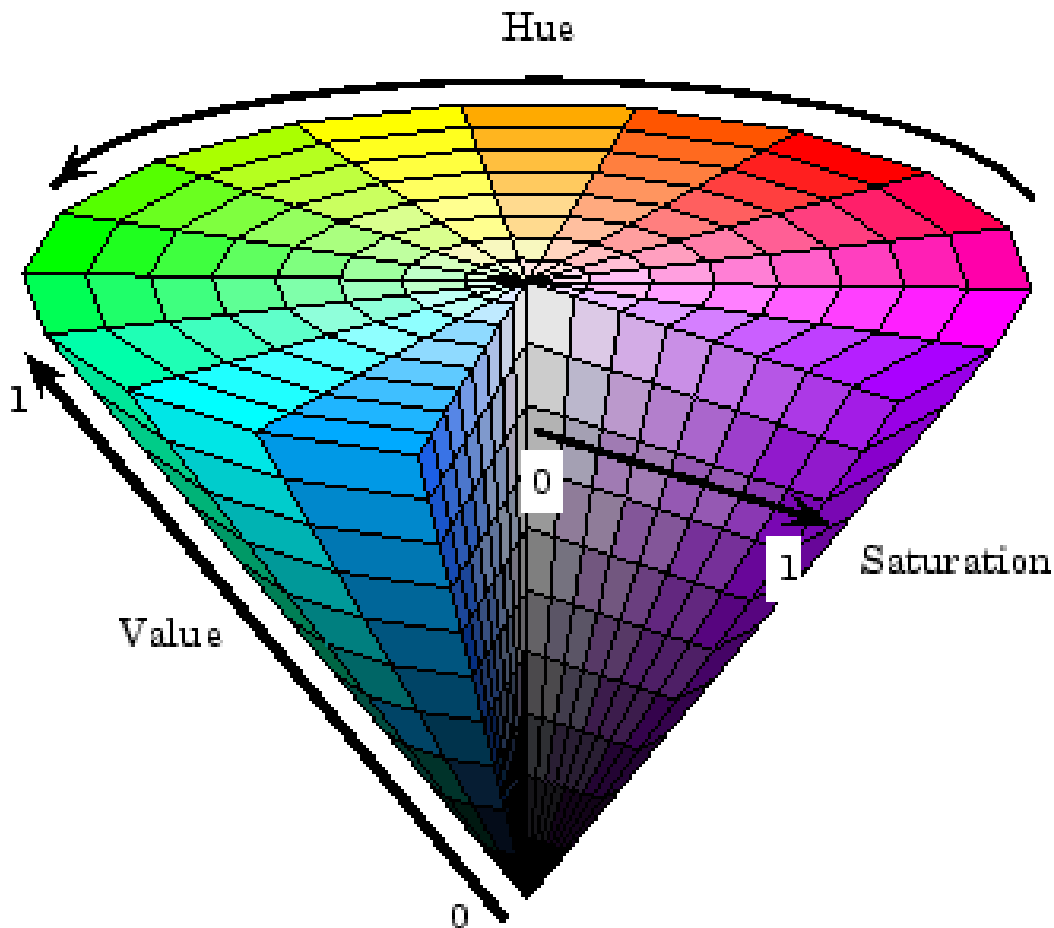
G
(R=0,B=0)



B
(R=0,G=0)

Colour spaces: HSV

Intuitive colour space



H
(S=1,V=1)



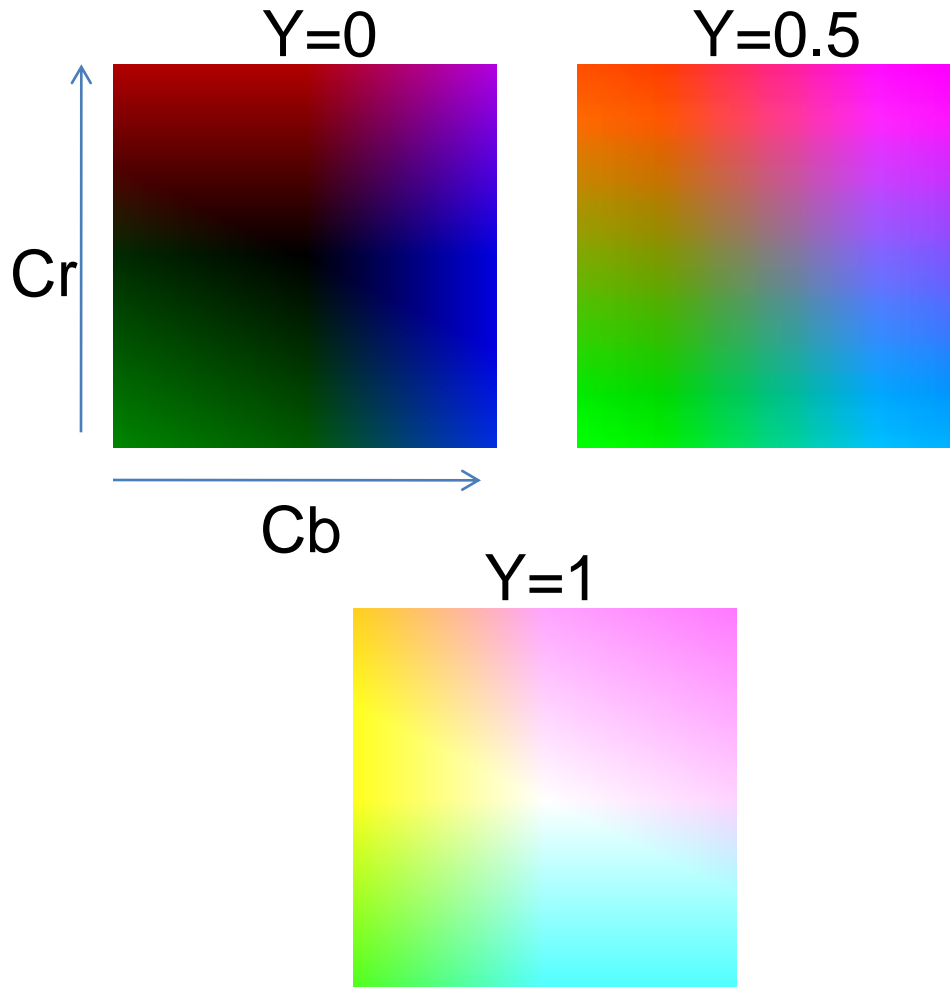
S
(H=1,V=1)



V
(H=1,S=0)

Colour spaces: YCbCr

Fast to compute, good for compression, used by TV



Y
(Cb=0.5,Cr=0.5)



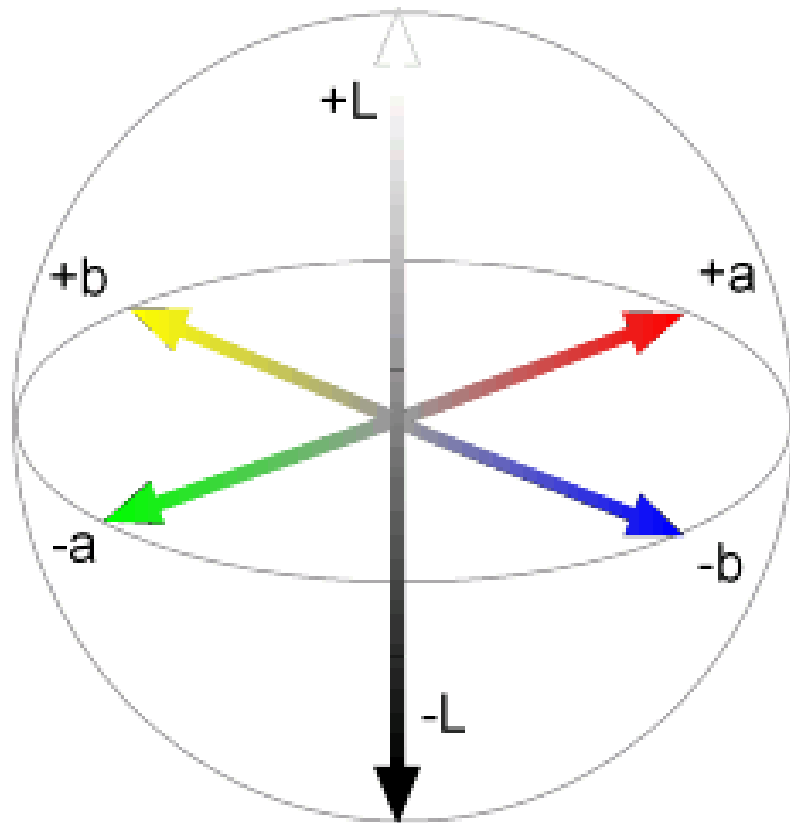
Cb
(Y=0.5,Cr=0.5)



Cr
(Y=0.5,Cb=0.5)

Colour spaces: $L^*a^*b^*$

“Perceptually uniform”^{*} colour space



L
($a=0, b=0$)



a
($L=65, b=0$)



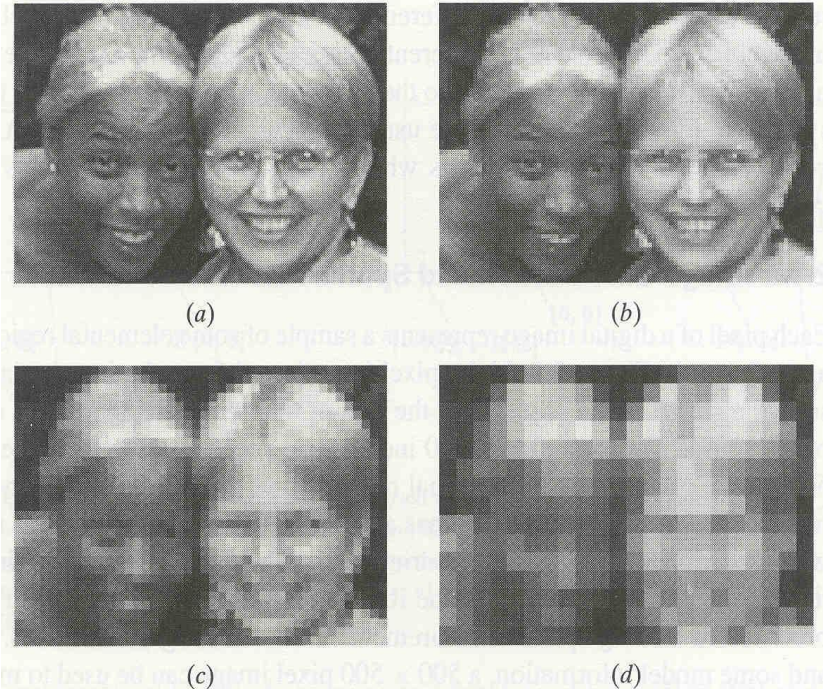
b
($L=65, a=0$)

Digitisation and Sampling

- Digitisation: converts analog image to digital image
- **Sampling** digitises the coordinates x and y :
 - *spatial discretisation* of picture function $F(x,y)$
 - use a grid of sampling points, normally rectangular: image sampled at points $x = j \Delta x$, $y = k \Delta y$, $j = 1 \dots M$, $k = 1 \dots N$.
 - Δx , Δy called the **sampling intervals**.

Spatial Resolution

- Spatial Resolution: number of pixels per unit of length
- Resolution decreases by one half- see right
- Human faces can be recognized at 64 x 64 pixels per face



- Appropriate resolution is essential:
 - too little resolution, poor recognition
 - too much resolution, slow and wastes memory

Quantisation

- **Quantisation** digitises the intensity or amplitude values, ie $F(x, y)$
 - called intensity or gray level quantisation
 - Gray-level resolution:
 - usually has 16, 32, 64, ..., 128, 256 levels
 - number of levels should be high enough for human perception of shading details - human visual system requires about 100 levels for a realistic image.

For Reading

- Szeliski, Chapter 2
- Shapiro and Stockman, Chapter 2

Acknowledgement

- Slides from Derek Hoiem, Alexei Efros, Steve Seitz, and David Forsyth
- Image sources credited where possible
- Some material, including images and tables, were drawn from the referenced textbooks and associated online resources.