

Disclaimer

This summary is part of the lecture “ETH Image Analysis & Computer Vision” by Prof. Van Gool, Prof. Konukoglu and Prof. Goksel (HS19). It is based on the lecture slides and script.

Please report errors to doerm@student.ethz.ch such that others can benefit as well.

The upstream repository can be found at <https://github.com/mrrebot/Summaries>

Image Analysis & Computer Vision

Marco Dober

1st March 2020

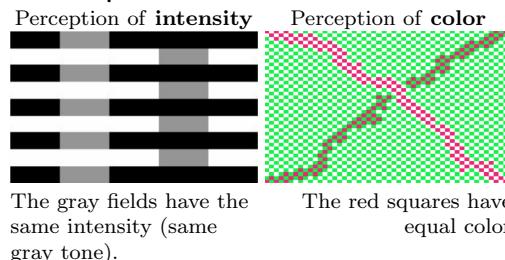
1 Introduction

Vision is important:

- Half our brain is devoted to it
- Developed many times during evolution
- It is non-contact
- It can be implemented with high-resolution
- Works with ambient EM-waves
- yields color, texture, depth, motion, shape

Take home message:
For people vision is their most crucial sense, for good reason

1.1 Perception of vision



The horizontal lines are equally long.

The lines do not have any curvature.

All lines are parallel.

There is no spiral.

Perception of parallelism

Perception of curvatures

Perception of length

Lines being straight

Perception of motion



The pole rotates about the vertical, it does not translate vertically.

The role of context



All encircled patterns are identical!

Augmented Reality, e.g. sports



Take home message:
It is feasible now to let most things see and interpret their environment.

Computer-assisted surgery



The visible range differs from humans to animals and also cameras may have different spectral sensitivities. There are also cameras for non-visible light such as infrared. The following picture shows the three color cones humans have and their sensitivity range: nm 350 400 450 500 550 600 650



1.3.2 Interactions with matter

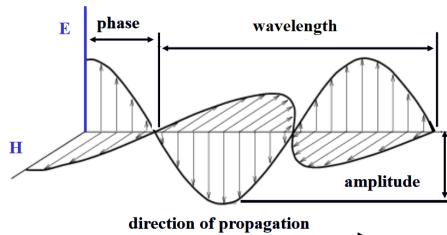
We look at the following types of interaction with matter:

1. **Absorption**
→ blue water
2. **Scattering**
→ blue sky
→ red sunset
3. **Reflection**
→ colored ink
4. **Refraction**
→ dispersion by a prism
5. **Diffraction**

We look at few of those in more detail:

1. Absorption

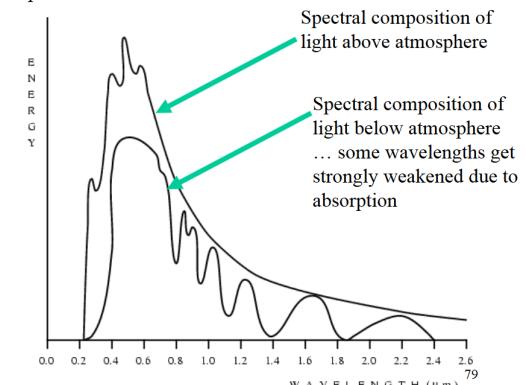
A nice example of absorption is earth's atmosphere which absorbs certain wavelengths of the incoming light. The absorbed frequencies correspond to resonance frequencies of molecules in earth's atmosphere.



- **Wavelength**
- **Direction of propagation**
- **Amplitude of E**
- **Phase**
- **Direction of polarization**

The spectrum:

Normal ambient light is a mixture of wavelengths, polarization directions and phases. The visible range for humans is only a small fraction of the EM-waves-spectrum.



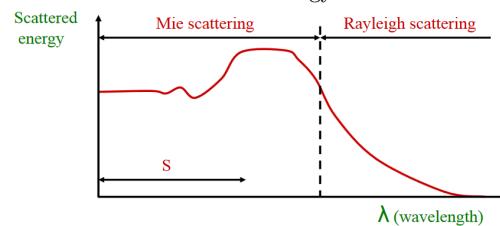
Wavelength [nm]	Color
380 - 450	→ violet
450 - 490	→ blue
490 - 560	→ green
560 - 590	→ yellow
590 - 630	→ orange
630 - 760	→ red

2. Scattering

There are three types of scattering depending on the relative sizes of particles and wavelengths:

- (a) Small particles: **Rayleigh** (strong wavelength dependent)
- (b) Comparable size: **Mie** (weakly wavelength dependent)
- (c) Large particles: **Non-selective** (wavelength independent)

If we look at the scattered energy it looks as follows:



Let's see some examples of these different scatter types in our atmosphere:



Rayleigh: Tyndall effect (blue sky, red setting sun)
Non-selective: Grey clouds



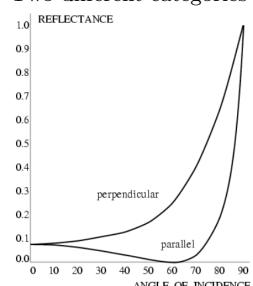
Mie: Colored cloud from volcanic eruption

3. Reflection:

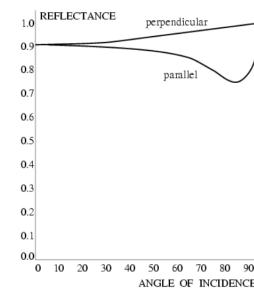
In mirror reflection we have:

angle of reflection = angle of incident.

Two different categories of reflective materials:



Dielectric:
For parallel polarization there exists the Brewster angle where $r = 0$.

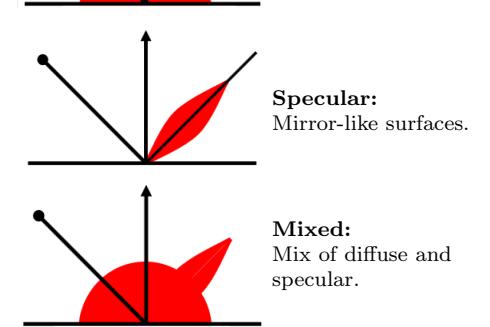


Conductor:
Strong reflectors under all angles, more or less preserve polarization.

We differentiate three types of reflection which depend on the surface structure:

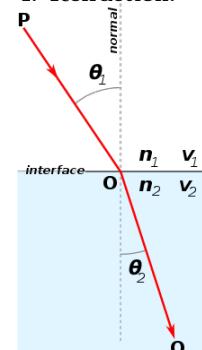


Specular:
Mirror-like surfaces.



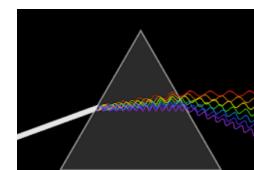
Mixed:
Mix of diffuse and specular.

4. Refraction:



Effect of the bending of light if it hits an interface of two materials with different refraction index $n = \sqrt{\epsilon\mu}$. The bending is described through **Snell's law**:

$$n_1 \sin \theta_1 = n_2 \sin \theta_2$$



Dispersion:
The bending is dependent of the frequency (wavelength) of the light.

2 Image Acquisition

2.1 Illumination

Well designed illumination often is key in visual inspection and can extremely simplify the image processing. Here is an overview of different illumination techniques:

2.1.1 Back-lighting



2.1.3 Diffuse-lighting

Left: Direct lighting produced large changes in brightness due to specular reflection.

right: Diffuse lighting reduces bright spots.

How:

Do not directly shine with light source on object, but rather indirectly with the help of a diffusive surface. It does not reduce the specular reflection, but increases the diffuse reflection component, yielding in less variations.

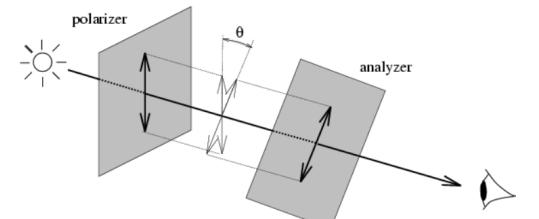
Why:

Prevents sharp shadows and large intensity variations over glossy surface.

2.1.4 Polarized-lighting

The polarization direction is the one of the E-Wave. Normally, light is composed of many waves with different polarizations.

Following picture shows a polarizer/analyzer configuration:



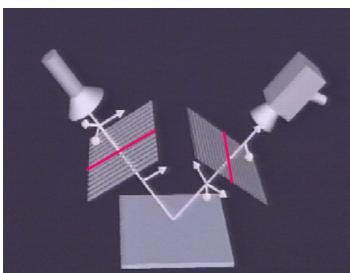
The intensity seen from the observer depends on the angle θ and is described by the **law of Malus**:

$$I(\theta) = I(0) \cos^2 \theta$$

There are 2 uses for polarized lighting:

1. Improve contrast between Lambertian and specular reflection.
2. Improve contrast between dielectrics and metals.

1. Specular vs. Lambertian



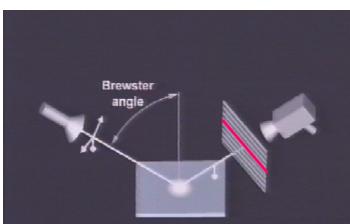
How:

Polarizer and analyzer in crossed arrangement. Specular reflection keeps polarization, Lambertian reflection depolarizes, because of this the arrangement reduces the large dynamic range caused by glare.

Why:

Increases contrast between Lambertian and specular reflection (specular reflection gets blocked).

2. Dielectric vs. Metal



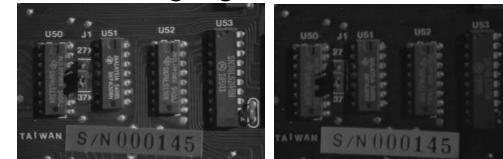
How:

Shine non-polarized light at Brewster angle on object. The dielectric will not reflect the parallel polarized component (Brewster angle...) and the perpendicular component is filtered out by the analyzer, hence the dielectric parts of the object will be really dark.

Why:

Increases contrast between dielectrics and metals (dielectric reflection gets blocked).

2.1.5 Colored-lighting



The contrast between the red label and the green PCB is increased when red light is shined on it (right) in comparison to white light (left).

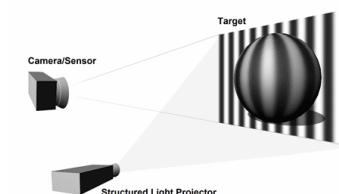
How:

Shine colored light at an object and maybe also use an bandpass filter. You need to have in mind the spectral sensitivity of your sensor!

Why:

- Highlight regions of similar color
- Differentiate between specular and diffuse reflection
- Comparing color

2.1.6 Structures-lighting



How:

Spatially or temporally modulated light patterns are shined on a 3D object.

Why:

Obtain 3D info of object.

More on this later...

2.1.7 Stroboscopic-lighting



How:

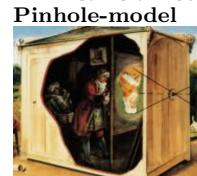
High intensity light flashes. The light flash artificially shortens the sensors integration time. Mostly this is much cheaper than fast cameras.

Why:

Eliminate motion blur

2.2 Cameras

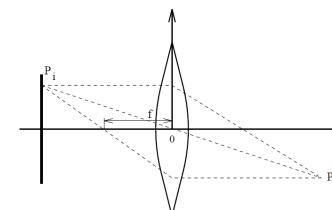
2.2.1 Camera models



Light enters the box only through a small hole and an image is formed on a plane inside the box. If the hole is too small not enough light enters and additional diffraction occurs, if the hole is too big, the image gets blurry. The solution to this problem is a lens. From similar triangles the following formula follows:

$$\frac{X_i}{X_o} = \frac{Y_i}{Y_o} = \frac{f}{Z_o} = -m = \text{linear magnification}$$

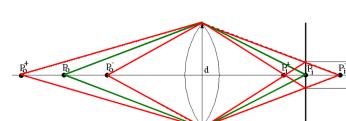
The thin-lens model



A lens captures more light and focuses it which gets rid off the problems of the pinhole. The price we pay is that only points at certain plane will be sharply imaged. Similar to the pinhole model the **thin lens equation** reads as:

$$\frac{1}{Z_o} - \frac{1}{Z_i} = \frac{1}{f}$$

The depth-of-field



As already mentioned only at a specific plane the image will be sharp, we can define an interval ΔZ_o^- in which we have a reasonable sharpness of the image:

$$\Delta Z_o^- = Z_o - Z_o^- = \frac{Z_o(Z_o - f)}{Z_o + fd/b - f}$$

The depth-of-field decreases with d and increases with Z_o . The more we focus with the lens, the smaller the depth-of-field is, because de rays diverge stronger outside of the focal point. There is also a trade off between collecting a lot light (big d) and a large depth-of-field (usable depth range).

Take home message:

In summary, introducing a lens to some extent solves the problem of insufficient light reaching the light detecting area of the camera. The price we pay is the loss of overall sharpness, i.e. points at a distance outside some range are no longer in focus.

2.2.2 Aberrations

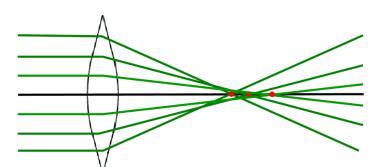
In the above lens-model we made three assumptions:

1. All rays from a point are focused onto 1 image point
2. All image points in a single plane
3. Magnification is constant

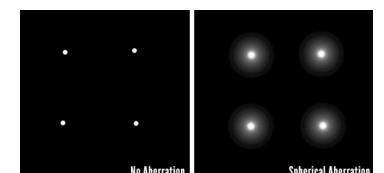
Deviations from this ideal are called **aberrations**. We differentiate between two types of aberrations:

1. **Geometrical:** visible as image distortions or degradation like blurring
 - Spherical aberration
 - astigmatism
 - **radial distortion** (most important)
 - coma
2. **Chromatic:** visible as different behavior for different wavelengths (color)

Spherical Aberration:



Rays parallel to the optical axis do not converge, because outer parts of the length yield smaller focal lengths. This results in blurry edges on the image.



Radial distortion:

Different magnification for different angles of incident. This results in:

- Lines become curves

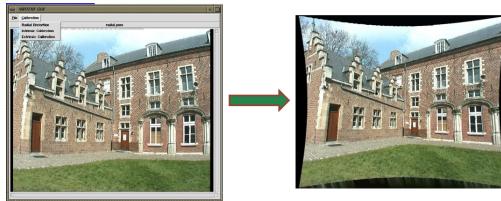
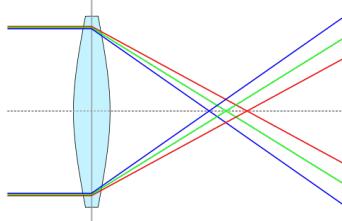
→ Curvature increases as you move away from the center of distortion

→ Model assume this is the image center and there is a multiplicative factor on the pixel depending on the distance r to the center:

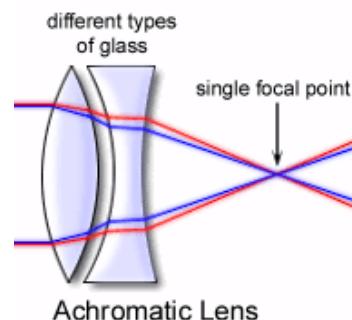
$$d = 1 + \kappa_1 r^2 + \kappa_2 r^4 + \dots$$

Only even factors because effect is symmetric.

This aberration can be corrected by software if the parameters $\kappa_1, \kappa_2, \dots$ are known.

**Chromatic distortion:**

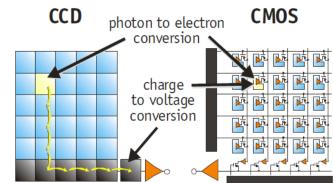
Rays of different wavelength focused on different planes. This can not be removed completely, but **achromatization** can be achieved at some well chosen wavelength pair, by combining lenses made of different glasses:



Sometimes achromatization is achieved for more than two wavelengths.

2.2.3 Device technologies

We consider 2 types:



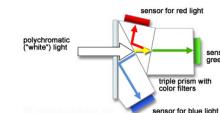
CCD	CMOS
Niche applications	Consumer cameras
Specific technology	Standard IC technology
Expensive product.	Cheap
High power	Low power
Higher fill rate	Less sensitive
Blooming	Per pixel amplif.
Sequential readout	Random pixel access
	Smart pixels
	On chip with other comp.

In 2006 was year of sales cross-over and in 2015 Sony said to stop CCD chip production.

2.2.4 Color cameras

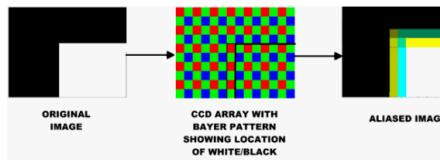
We consider 3 concepts:

1. Prism (with 3 Sensors)
2. Filter mosaic
3. Filter wheel

1. Prism color camera

Separates light in 3 different beams using dichroic prism

- Requires 3 Sensors & precise alignment
- + Good color separation

2. Filter mosaic

Filter is directly coated on sensor. **Microlenses** are used to gain more light on pixels, because effective resolution is reduced by the filter.

- Reduces resolution
- + Cheap and easy

3. Filter wheel

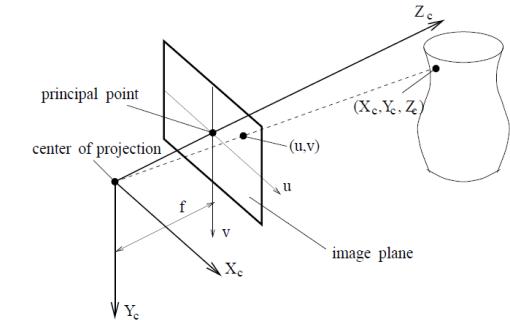
Rotate multiple filters in front of lens

- Only suitable for static scenes
- + Allows more than 3 color bands

	Prism	Mosaic	Wheel
#Sensors	3	1	1
Resolution	High	Average	Good
Cost	High	Low	Average
Framerate	High	High	Low
Artefacts	Low	Aliasing	Motion
Bands	3	3	3 or more
	High-end	Low-end	Scientific

2.2.5 Geometric models**Perspective projection**

Perspective projection is the projection of the object to a **virtual image plane** in front of the lens. We choose the virtual image plane so that we do not need to care about the rotation of the picture. The **center of projection** is the **center of the lens** (pinhole).



- Camera coordinate frame (hence subscript c) lies at center of projection.
- Z_c coincides with the optical axis of the lens/objective.
- X_c is parallel to image rows, Y_c is parallel to image columns.
- u and v axes are parallel to X_c and Y_c axes.
- Principal point = point where optical axis intercepts image plane.
- The image of a point P (X_c, Y_c, Z_c) is the intersection of the line through P and the center of projection with the image plane (plane at f).

The (u, v) -coordinates of this image point are directly found through similar triangles:

$$u = f \frac{X_c}{Z_c} \quad v = f \frac{Y_c}{Z_c}$$

This model is an approximation, because for the image plane to lie a focal distance from the center of projection assumes the object to lie at a large distance compared to the focal length.

Pseudo-orthographic projection

First, if Z is constant for all points of an object:

$$x = kX \quad y = kY, \quad k = f/Z$$

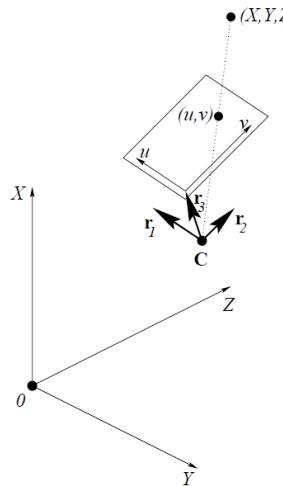
Second, this is also a good approximation if $f/Z \approx \text{constant}$, i.e. objects are small compared to distance from camera.

Projection matrices

The perspective projection model is not very practical and incomplete in at least 2 ways:

1. 3D coordinates are specified in a *world coordinate frame*, not the camera coordinate frame.
2. Image coordinates (u, v) must be translated to row and column numbers (pixels).

No additional refinements such as radial distortions are considered.



The position of the Camera is described by a point **C** (center of projection) and a 3×3 rotation matrix \mathbb{R} where \mathbf{r}_i denotes the i th column of the matrix. \mathbf{r}_1 corresponds to X_c , \mathbf{r}_2 to Y_c and \mathbf{r}_3 to Z_c of the camera coordinate frame. The corresponding image point then has (u, v) -coordinates:

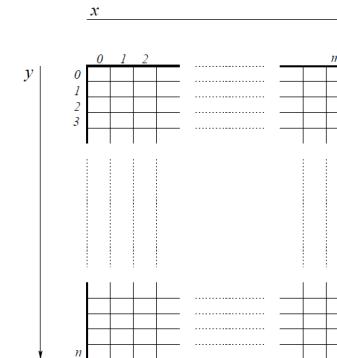
$$u = f \frac{\langle \mathbf{r}_1, \mathbf{P} - \mathbf{C} \rangle}{\langle \mathbf{r}_3, \mathbf{P} - \mathbf{C} \rangle} \quad \text{and} \quad v = f \frac{\langle \mathbf{r}_2, \mathbf{P} - \mathbf{C} \rangle}{\langle \mathbf{r}_3, \mathbf{P} - \mathbf{C} \rangle}$$

If **P** has coordinates (X, Y, Z) and **C** has coordinates (C_1, C_2, C_3) with respect to the *world frame*, then:

$$u = \frac{r_{11}(X_1 - C_1) + r_{12}(Y_1 - C_2) + r_{13}(Z_1 - C_3)}{r_{31}(X_1 - C_1) + r_{32}(Y_1 - C_2) + r_{33}(Z_1 - C_3)}$$

$$v = \frac{r_{21}(X_1 - C_1) + r_{22}(Y_1 - C_2) + r_{23}(Z_1 - C_3)}{r_{31}(X_1 - C_1) + r_{32}(Y_1 - C_2) + r_{33}(Z_1 - C_3)}$$

When working with digital images we want to define the position of an image point in so-called *pixel coordinates*. A digital image thus is a 2-dimensional array of pixels, say m -columns and n -rows:



We want to indicate the position of an image point by its column and row number instead of the (u, v) -coordinates.

$$x = k_x u + sv + x_0 \\ y = k_y v + y_0$$

Here:

- (x_0, y_0) are the pixel coordinates of the principal point.
- k_x is the number of pixels per unit length in horizontal direction (describes 1/width of a pixel)
- k_y is the number of pixels per unit length in vertical direction (describes 1/height of a pixel)
- k_x/k_y aspect ratio, if $\neq 1$, then pixel ware not square
- s indicates the skew, how much pixel deviates from a rectangular.

k_x, k_y, s, x_0, y_0 are called **internal camera parameters**. When they are known the camera is **internally calibrated**.

Vector **C** and matrix **R** are the **external camera parameters**. When they are known the camera is **externally calibrated**.

Fully calibrated means internally and externally calibrated.

Perspective projection being a projective mapping, the image formation process in a pinhole camera (as described above) can be represented **linearly** by using **homogeneous coordinates**. Examples:

$$\text{2D: } \begin{pmatrix} x \\ y \\ z \end{pmatrix} \rightarrow \begin{pmatrix} x/z \\ y/z \end{pmatrix} \quad \text{3D: } \begin{pmatrix} X \\ Y \\ Z \\ W \end{pmatrix} \rightarrow \begin{pmatrix} X/W \\ Y/W \\ Z/W \end{pmatrix}$$

Exploiting homogeneous coordinates we can write (u, v) as:

$$\tau \begin{pmatrix} u \\ v \\ 1 \end{pmatrix} = \begin{pmatrix} fr_{11} & fr_{12} & fr_{13} \\ fr_{21} & fr_{22} & fr_{23} \\ r_{31} & r_{32} & r_{33} \end{pmatrix} \begin{pmatrix} X - C_1 \\ Y - C_2 \\ Z - C_3 \end{pmatrix}$$

and (x, y) as:

$$\tau \begin{pmatrix} x \\ y \\ 1 \end{pmatrix} = \begin{pmatrix} k_x & s & x_0 \\ 0 & k_y & y_0 \\ 0 & 0 & 1 \end{pmatrix} \tau \begin{pmatrix} u \\ v \\ 1 \end{pmatrix}$$

Concatenating the results:

$$\tau \begin{pmatrix} x \\ y \\ 1 \end{pmatrix} = \begin{pmatrix} k_x & s & x_0 \\ 0 & k_y & y_0 \\ 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} f & 0 & 0 \\ 0 & f & 0 \\ 0 & 0 & f \end{pmatrix} \cdot \begin{pmatrix} r_{11} & r_{12} & r_{13} \\ r_{21} & r_{22} & r_{23} \\ r_{31} & r_{32} & r_{33} \end{pmatrix} \begin{pmatrix} X - C_1 \\ Y - C_2 \\ Z - C_3 \end{pmatrix}$$

This yields the **calibration matrix K**:

$$K = \begin{pmatrix} k_x & s & x_0 \\ 0 & k_y & y_0 \\ 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} f & 0 & 0 \\ 0 & f & 0 \\ 0 & 0 & f \end{pmatrix} = \begin{pmatrix} fk_x & fs & x_0 \\ 0 & fk_y & y_0 \\ 0 & 0 & 1 \end{pmatrix}$$

We define: $\mathbf{p} = \begin{pmatrix} x \\ y \\ 1 \end{pmatrix}$; $\mathbf{P} = \begin{pmatrix} X \\ Y \\ Z \end{pmatrix}$; $\tilde{\mathbf{P}} = \begin{pmatrix} X \\ Y \\ Z \\ 1 \end{pmatrix}$

yielding the **projection equation** for a pinhole camera, for some non-zero $\rho \in \mathbb{R}$:

$$\rho \mathbf{p} = K R^t (\mathbf{P} - \mathbf{C})$$

If also homogeneous coordinates are used for the scene points, then the projection equation becomes:

$$\rho \mathbf{p} = K (R^t \mid -R^t \mathbf{C}) \tilde{\mathbf{P}} \quad \mid := \text{concatenate}$$

or,

$$\rho \mathbf{p} = K (M \mid \mathbf{t}) \tilde{\mathbf{P}} \quad \text{with rank } M = 3$$

The last equation follows from linear algebra from which we see that $K R^t$ can be any invertible 3×3 matrix M and that $-K R^t \mathbf{C}$ can be any column vector $\mathbf{t} \in \mathbb{R}^3$.

2.2.6 Photometric camera model

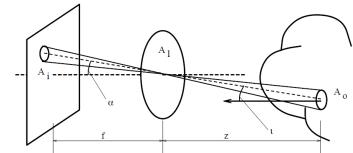
From object radiance to pixel gray level in 2 steps:

1. From object radiance to image irradiance
2. From image irradiance to pixel gray level.

1. From object radiance to image irradiance: We look at the irradiance that an object patch will cause in the image. Assumptions:

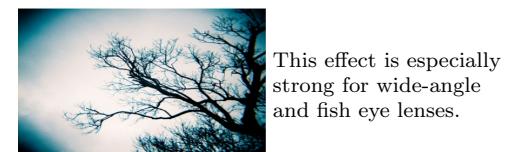
- Radiance R is known
- Object at far distance compared to focal length

Is image irradiance directly related to radiance of the object patch? We look at the following viewing condition:



With a few calculations and a little black magic we obtain the **cos⁴ law**:

$$I = R \frac{A_i}{f^2} \cos^4 \alpha$$



This effect is especially strong for wide-angle and fish eye lenses.

2. From irradiance to gray levels

The relation of irradiance to gray level has the form of:

$$f = g I^\gamma + d$$

f: resulting gray level

I: irradiance

g: so-called *gain* of camera (constant)

γ : non-linearity of the camera

d: *dark reference*

The *gain* g can be set with the diaphragm size. Nowadays γ is pretty close to 1. The *dark reference* d is the signal measured with the lens cap on.

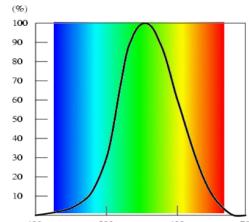
3 Feature Extraction

3.1 Color

3.1.1 Radiometry vs. Photometry

Photometry: subjective impression

Radiometry: objective, physical measurements



Luminous efficiency function $v(\lambda)$ relates radiometry & photometry.
C.I.E standard

at 555nm: $1lm = 1/683W = 1.46mW$

For light with spectral composition $c(\lambda)$ (radian flux):

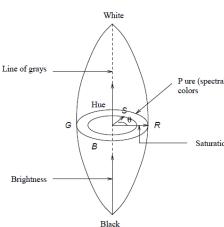
$$l = k \int_{\lambda=0}^{\infty} c(\lambda)v(\lambda)d\lambda, k = 683 \text{ lumens/Watt}$$

3.1.2 Perceptual dimensions of color

The perceptual dimensions of color are:

- Luminance (brightness)
- Hue
- Saturation

This concept can be represented in the football-shaped color space:

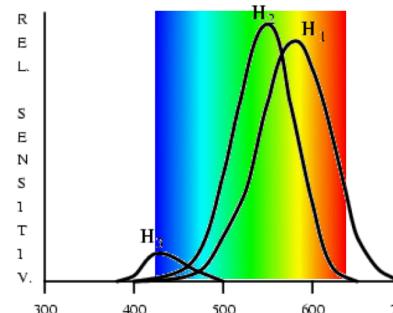


shaped color space:
brightness varies along
vertical axis.
hue θ varies along
circumference
saturation varies along
radius.

3.1.3 Tristimulus model

Young said we can represent any color, by the mixture of **three primary colors**. This does **not** mean, we can actually generate any color with the mixture of those three, but humans can't tell the difference. This model was underpinned as three different kind of cones were found in the human retina.

The following picture shows the three cone sensitivity curves $H_1(\lambda), H_2(\lambda), H_3(\lambda)$:



A source with spectral radiant flux $C(\lambda)$ produces responses R_i :

$$R_i(c) = \int H_i(\lambda)C(\lambda)d\lambda, i = 1, 2, 3$$

An entire distribution over λ is projected to only three numbers R_i !

The three by CIE recommended **primaries** $P_j(\lambda)$, $j = 1, 2, 3$ are:

- $\lambda_1 = 700.0nm$, red
- $\lambda_2 = 546.1nm$, green
- $\lambda_3 = 435.8nm$, blue

Different organizations have defined different primaries for practical applications, e.g. TV.

We now want to match the source $C(\lambda)$ by the primaries with $\sum_{j=1}^3 m_j P_j(\lambda)$ such that R_i will be the same:

$$\begin{aligned} R_i(\lambda) &= \int \sum_{j=1}^3 m_j P_j(\lambda) H_i(\lambda) d\lambda \\ &= \sum_{j=1}^3 m_j \int P_j(\lambda) H_i(\lambda) d\lambda \end{aligned}$$

The integrals can be calculated off-line, because they are fix once the primaries have been chosen, hence we can write: $l_{i,j} = \int P_j(\lambda) H_i(\lambda) d\lambda$
Knowing the human responses R_i to the source, we obtain a linear system of equations to be solved for the m_j :

$$\sum_{j=1}^3 m_j l_{i,j} = R_i$$

To solve this we need to invert the matrix $l_{i,j}$, therefore the primaries have to be independent with respect to human vision (none of them can be produced as a linear combination of the other two). When we use different primaries, we obtain a linear transformation between the two m_j s of the two

system of primaries.

$$L \cdot m = R$$

$$L' \cdot m' = R$$

gives

$$m' = L'^{-1} \cdot L \cdot m$$

3.1.4 Tristimulus values

Since "white" can be considered a natural reference, one will usually specify relative values with respect to the amounts of the primaries needed for some standard white source, written w_j . These **tristimulus values** for the source $C(\lambda)$ are then:

$$T_j = \frac{m_j}{w_j}$$

By definition, the standard white source has **tristimulus values 1** ($T_1 = T_2 = T_3 = 1$). For the CIE primaries, the corresponding tristimulus values are generally called **R,G,B**. The CIE standard white is defined to have a flat energy spectrum ($w_1 = w_2 = w_3$).

The **spectral matching curves** $T_j(\lambda')$ give the tristimulus values for monochromatic sources $C_{\lambda'}$ with wavelength λ' :

$$R_i(C_{\lambda}) = H_i(\lambda) = \sum_{j=1}^3 m_j l_{i,j} = \sum_{j=1}^3 w_j l_{i,j} T_j(\lambda)$$

3.2 Texture