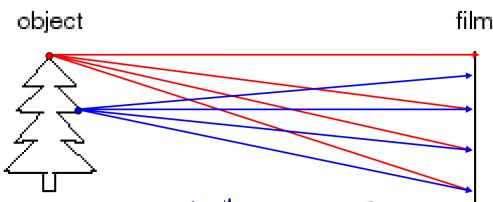


Midterm Review (First half lecture note)

Scope: Beginning → Oct. 9, stereo

Lec 1: Camera

物体不是主动发出光的，而是光照在物体上反射。以下为分析简便，



认为是物体发射 rays of light.

How to record (film / take picture) of them?

If we just put a film -----

⇒ Average of scene! Blurry!

What if we add a barrier?

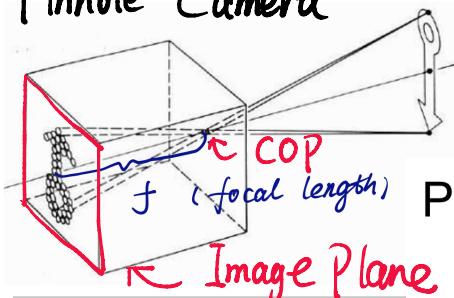
- Reduce blurring

- Opening is known as aperture (光圈)

- Images are upside-down

With this, we can build pinhole camera!

Pinhole Camera

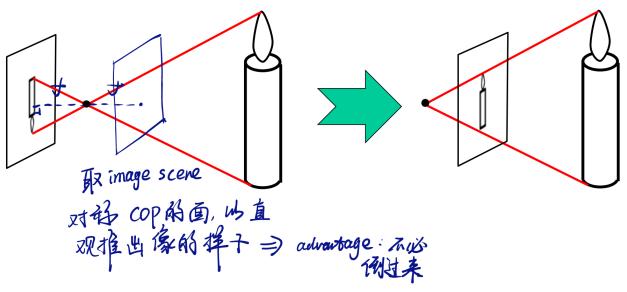


Pinhole model:

- Captures pencil of rays – all rays through a single point
- The point is called Center of Projection (COP) 投影中心
- The image is formed on the Image Plane
- Effective focal length f is distance from COP to Image Plane

In this model, using similar triangles, we can compute the coordinate mapping:

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



Also, we can use the trick shown on the left to avoid inversion:

We can (and most of time we do) pretend that the image plane is in front of camera center.

But our eyes, who generate images of 2D for us, can help us perceive 3D world. Why is that so?

3D感知来自于：双目立体视觉 or 移动+经验.
Why did evolution opt for such strange solution?

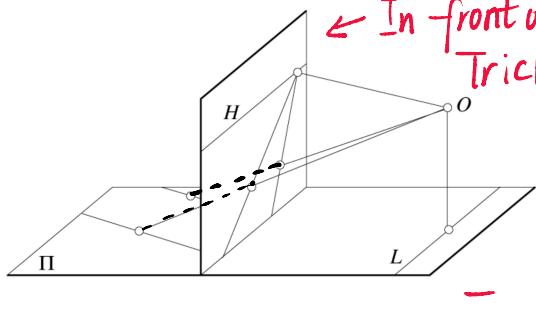
- Nice to have a passive, long-range sensor
- Can get 3D with stereo or by moving around, plus experience

Passive: We capture lights, instead of using sonar to detect'

Long-range: Light can be far away

Stereo: 两只眼从略微不同角度看世界，会产生两个略微不同的2D，大脑通过比较这个差异，可精确计算物体远近

But from 3D \rightarrow 2D, there are perspective projection:



\leftarrow In front of camera

Trick!

Parallel lines in real 3D world don't go parallel in 2D images

正交 平行投影
Orthographic or parallel projection

- What happens if we walk infinitely far away and zoom infinitely far in?

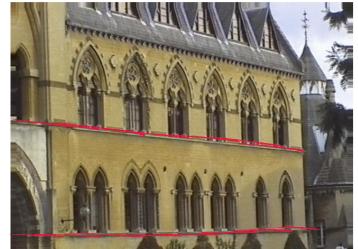
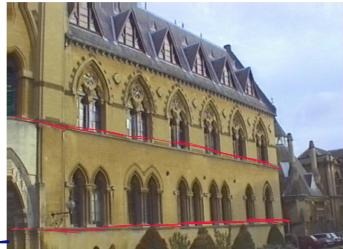
Left: Intersect



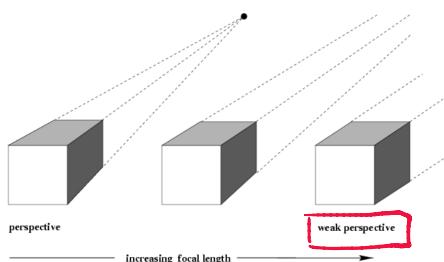
Right: Almost parallel

平行投影：
所有投影线平行

正交：不仅
平行，而且
垂直于投影平面。



同时，depth也开始不可信.....



For perspective, if we let COP far away from obj and let focal length also rather far?

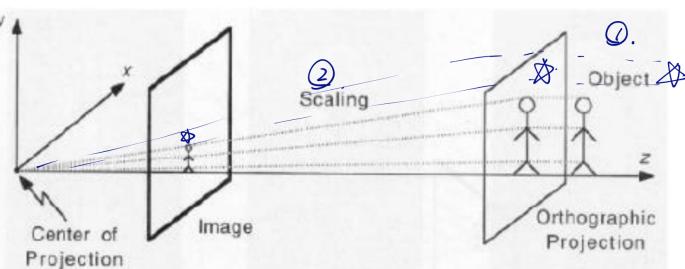
i.e., :

What happens if we walk infinitely far away and zoom infinitely far in?

\Rightarrow We are getting close to parallel perspective.

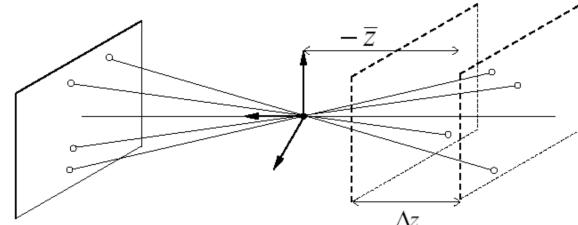
Or we can use another approximate modeling:

Scaled Orthographic or "Weak Perspective"



First orthographic one,
then normal scaling

If $\Delta z \ll -\bar{z}$, ortho first
then scale.



Orthographic or parallel projection

- Special case of perspective projection [透视投影]
 - Distance from center of projection to image plane is infinite
 - Projection equation: simply drop the z coordinate!
- ∞ distance
- COP (Center of Projection)
- image plane
- World
- $(x, y, z) \rightarrow (x, y)$

$$\text{If } \Delta z \ll -\bar{z} : \begin{aligned} x' &\approx -mx \\ y' &\approx -my \end{aligned} \quad m = -\frac{f}{\bar{z}}$$

Justified if scene depth is small relative to average distance from camera

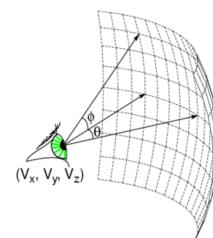
Wrap up.

Another Projection:

Spherical One

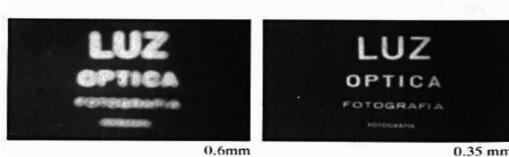
(Actually more similar

Spherical Projection



to how our eyes see !) $(\theta, \phi, z) \rightarrow (\theta, \phi)$

⊗ Doesn't depend on focal length



Now if wanting

to build a camera. a few things to calibrate Like Aperture.

⇒ How aperture affects image !

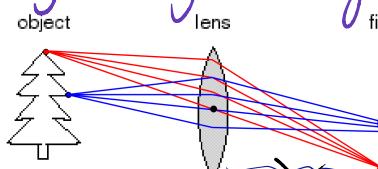
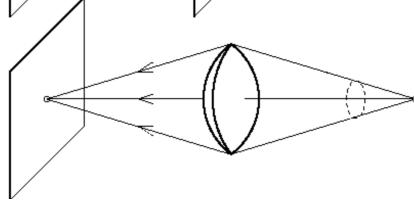
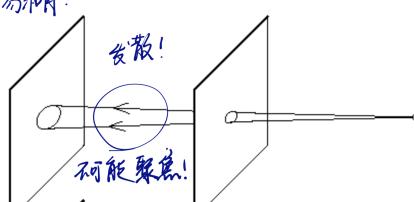
Too large & Too small ⇒ blurry

⇒ Easy to understand ↓ (衍射)

That's why we need a len: ① less light to get through ② Diffraction

可见 pin hole size (aperture) 很影响呈像大小与质量 !

很容易糊 !



Focal Point: 平行光
穿过 len 后 汇
聚的 焦点 .

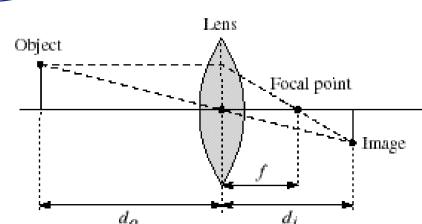
A lens focuses light onto the film

There is a specific distance at which objects are "in focus"
other points project to a "circle of confusion" in the image

• Changing the shape of the lens changes this distance

$$\frac{1}{d_o} + \frac{1}{d_i} = \frac{1}{f}$$

d_o : 物 → len d_i : "in focus" → len



$$\text{Thin lens equation: } \frac{1}{d_o} + \frac{1}{d_i} = \frac{1}{f}$$

Depth of Field 景深

看起来清晰可接受的范围或区域



DEPTH OF FIELD
DEPTH OF FIELD

光圈控制景深!!

Aperture ↑, DOF ↓

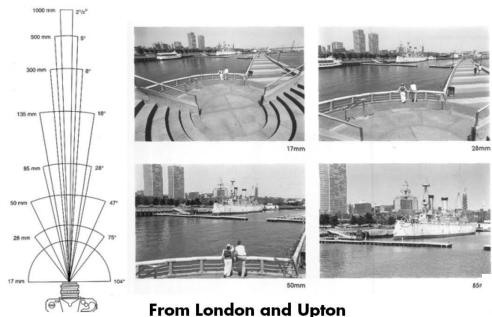
If Aperture ↓, DOF ↑, and
need more exposure

F-number: focal length / aperture diameter

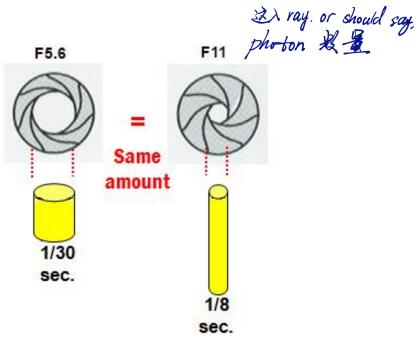
$$F = \frac{\text{focal length } (f)}{\text{aperture diameter}} \quad , \text{ 光圈越大, 在 f 不变下, } \text{ 光圈越小}$$

In camera community:

Field of View (Zoom)



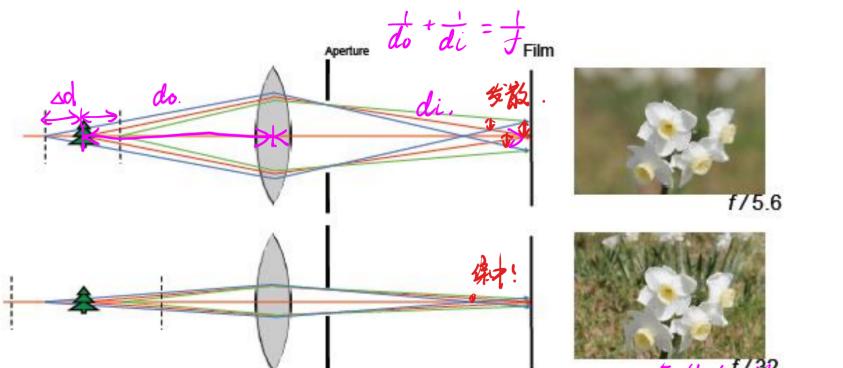
Exposure: shutter speed vs. aperture



But lens
aren't
flawless:

Another topic is Depth of Field (DOF)
Settled di & f => settled do ??

Aperture controls Depth of Field



光圈越大，景深越小。因为 DOF 深/浅物体进来的 ray 少且集中

Changing the aperture size affects depth of field

- A smaller aperture increases the range in which the object is approximately in focus
- But small aperture reduces amount of light - need to increase exposure 但同时, ray 少, 因此曝光要多些

$$\text{F-number} = \frac{\text{focal length}}{\text{aperture diameter}} \propto \text{DOF}$$

Another topic: FOV (Field of view).

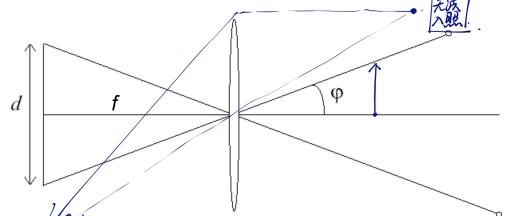
With fixed image

size, f ↑, FOV ↓

$$\varphi = \tan^{-1} \left(\frac{d}{2f} \right)$$

Field of View / Focal Length

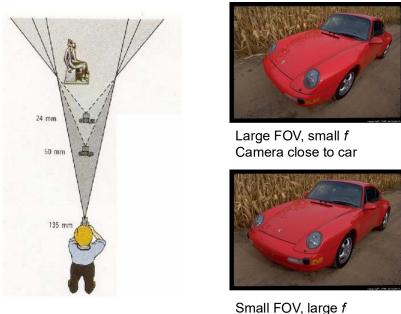
FOV depends of Focal Length



Size of field of view governed by size of the camera retina:

$$\varphi = \tan^{-1} \left(\frac{d}{2f} \right) \quad \varphi = \arctan \left(\frac{d}{2f} \right)$$

Smaller FOV = larger Focal Length

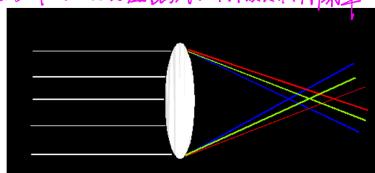


Finally: exposure:
shutter speed ↓ aperture.

How much photon / rays of light?

Lens Flaws: Chromatic Aberration
色散 Dispersion: wavelength-dependent refractive index
• enables prism to spread white light beam into rainbow

Modifies ray-bending and lens focal length: $f(\lambda)$
白光实像由多种波长的光组成而成: 不同波长不同折射率



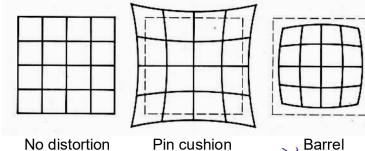
✿ color fringes near edges of image

Corrections: add 'doublet' lens of flint glass, etc.

① Rays of different
wave length → diff
ray-bending → Dispersion

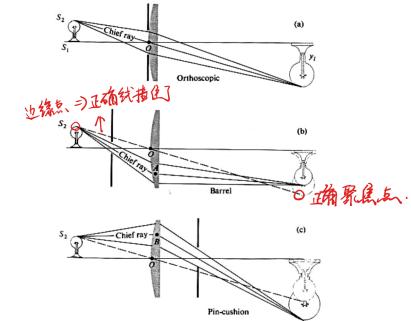
→ Chromatic Aberration.
易于 image obj 边缘处产生!

straight lines curve around the image center



Radial distortion of the image

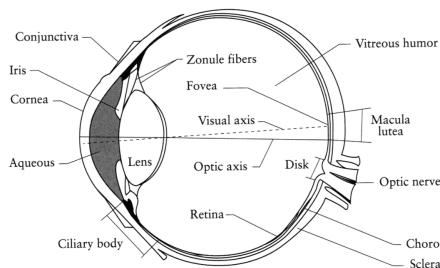
- Caused by imperfect lenses
 - Deviations are most noticeable for rays that pass through the edge of the lens
- 边缘最显著.



② Bending can be different at the center & fringe part of len

Lec2. Capturing Light

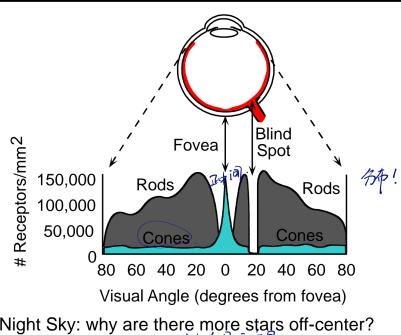
The Eye



The human eye is a camera!

- Iris - colored annulus with radial muscles 红膜→光圈
- Pupil - the hole (aperture) whose size is controlled by the iris
- What's the "film"? 瞳孔→光圈
- photoreceptor cells (rods and cones) in the retina 脉络膜

Distribution of Rods and Cones



Night Sky: why are there more stars off-center?
夜视星更亮

Our eye :

Iris control Aperture

Pupil → Aperture

Retina → Film

On retina:

cone & rod ← 光, gray
↑ 光, color

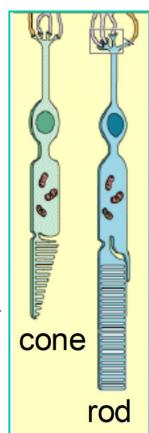
Cones
cone-shaped
less sensitive
operate in high light
color vision

白天
color vision

Rods

rod-shaped
highly sensitive
operate at night
gray-scale

夜晚
灰阶

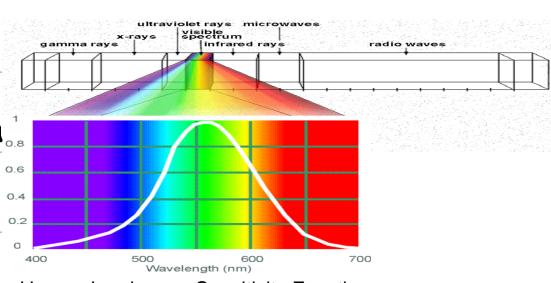


Retina 小角度范围内分布绝大部分 cone.

而其它大角度分布大量 cones

同时，光是一种 wave,

人可感知光的 wavelength
有限 (400-700nm) =>
这恰好也是太阳

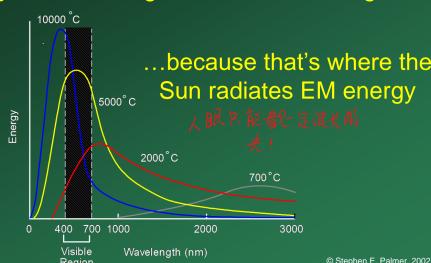


Human Luminance Sensitivity Function

温度发出光线的 wavelength 分布中占比亮的
wavelength zone !!

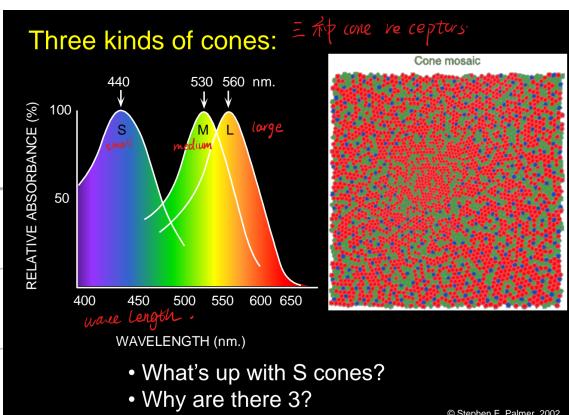
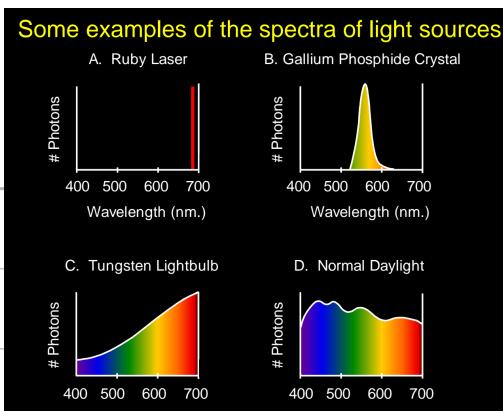
Any patch of light can be completely described physically by its spectrum: the number of photons (per time unit) at each wavelength 400 - 700 nm.

Why do we see light of these wavelengths?



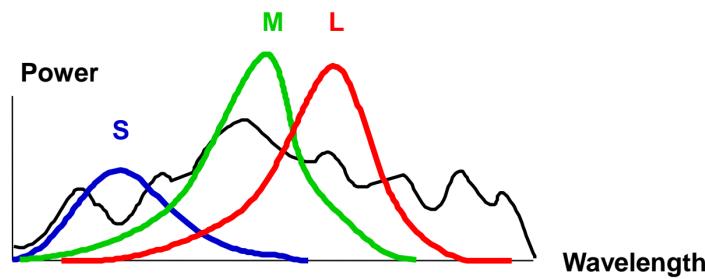
...because that's where the Sun radiates EM energy
人眼只能看到这个温度的光！

© Stephen E. Palmer, 2002



三种 cone.
吸收波长范围
如左图所示

Trichromacy



Rods and cones act as filters on the spectrum

- To get the output of a filter, multiply its response curve by the spectrum, integrate over all wavelengths
- Each cone yields one number
不能用三种就表达所有光谱
- How can we represent an entire spectrum with 3 numbers?
- We can't! A lot of the information is lost
 - As a result, two different spectra may appear indistinguishable
» such spectra are known as **metamers** 同色异谱物

Slide by Steve Seitz

每个 cone 视为 filter,
关心“真实光谱”中有多少
能量落在 cone 的敏感区域内
※: Each cone produces

one number

⇒ Represent spectrum with
3-dim data

Metamers: 两种光，物理
波长不同，但人眼看颜色相同

颜色恒常性 → 让我们感知物体真实颜色

Color Constancy: the ability to perceive the invariant color of a surface despite ecological variations in the conditions of observation.
而非反射到我们眼中的物理颜色

Another of these hard inverse problems:
Physics of light emission and surface reflection underdetermine perception of surface color

同时，我们对颜色的感知也可能受环境 /
心理影响。“看颜色”并非“所见即所得”
感知到物体表面不变的颜色尽管观察的
环境条件（如光照）在不断变化

Camera White Balancing



- Manual
 - Choose color-neutral object in the photos and normalize
- Automatic (AWB)
 - Grey World: force average color of scene to grey
 - White World: force brightest object to white

← (自动) 白平衡

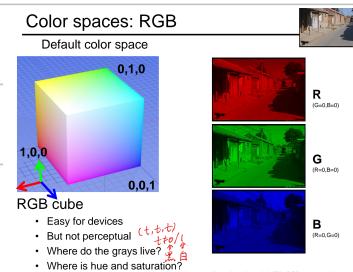
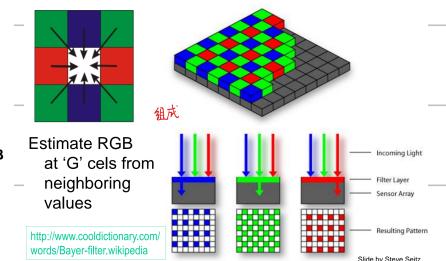
Image representation

- Images represented as a matrix
- Suppose we have an $N \times M$ RGB image called "im"
 - $im(1, 1)$ = top-left pixel value in R-channel
 - $im(y, x, b)$ = y pixels down, x pixels to right in the b^{th} channel
 - $im(N, M, 3)$ = bottom-right pixel in B-channel

row	column	R	G	B
0.92	0.93	0.94	0.97	0.98
0.89	0.90	0.91	0.96	0.97
0.86	0.88	0.85	0.93	0.94
0.83	0.85	0.82	0.90	0.91
0.80	0.82	0.79	0.87	0.88
0.77	0.74	0.58	0.51	0.52
0.74	0.69	0.55	0.48	0.49
0.71	0.66	0.52	0.45	0.46
0.68	0.63	0.49	0.42	0.43
0.65	0.60	0.46	0.39	0.40
0.62	0.57	0.43	0.36	0.37
0.59	0.54	0.40	0.33	0.34
0.56	0.51	0.37	0.29	0.30
0.53	0.48	0.34	0.26	0.27
0.50	0.45	0.31	0.23	0.24
0.47	0.42	0.28	0.19	0.20
0.44	0.39	0.25	0.16	0.17
0.41	0.36	0.22	0.13	0.14
0.38	0.33	0.19	0.10	0.11
0.35	0.30	0.16	0.07	0.08
0.32	0.27	0.13	0.04	0.05
0.29	0.24	0.10	0.01	0.02
0.26	0.21	0.07	0.00	0.01
0.23	0.18	0.04	0.00	0.00
0.20	0.15	0.01	0.00	0.00
0.17	0.12	0.00	0.00	0.00
0.14	0.09	0.00	0.00	0.00
0.11	0.06	0.00	0.00	0.00
0.08	0.03	0.00	0.00	0.00
0.05	0.00	0.00	0.00	0.00
0.02	0.00	0.00	0.00	0.00
0.00	0.00	0.00	0.00	0.00

RGB Bayer Grid

Practical Color Sensing: Bayer Grid



RGB color space: Gray: (t, t, t) (黑/白)
Hue: 绕着 Gray Diagonal 角度。
Saturation: 到 Gray Diagonal 的距离。

Lec3. Pixels and Images

We can think of an **image** as a function, f , from \mathbb{R}^2 to

R:

強度.

- $f(x, y)$ gives the **intensity** at position (x, y)
- Realistically, we expect the image only to be defined over a rectangle, with a finite range:
 $f: [a, b] \times [c, d] \rightarrow [0, 1]$

What is an image now?

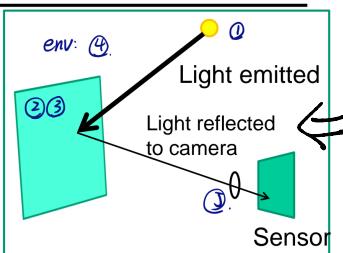
We use RGB! So for (x, y) .

we need intensity of RGB channels.

How does a pixel get its value?

Major factors

- ① Illumination strength and direction
- ② Surface geometry
- ③ Surface material
- ④ Nearby surfaces
- ⑤ Camera gain/exposure



So the intensity of light matters. How a pixel get its value? See left.

How does ② & ③ decide how intense the reflected light is?

Some light is absorbed (function of albedo ρ)
Remaining light is scattered (diffuse reflection)
Examples: soft cloth, concrete, matte paints

⇒ A simple and practical model:

Lambertian Reflectance Model

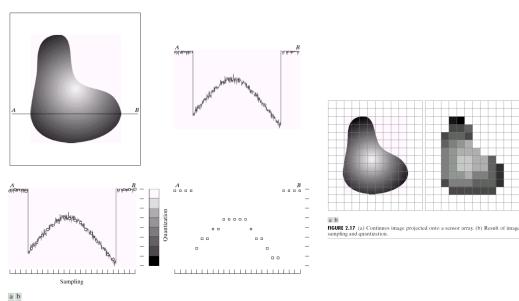
△ Suitable Material: soft cloth, concrete, matte paints

We know light can specular (鏡面反射) or Diffuse (漫反射), in Lambertian we suppose purely diffuse:

$$I = I_{\text{source}} \cdot \rho (\text{albedo}) \cdot (\vec{N} \cdot \vec{L})$$

\vec{L} : light direction vector, $\|\vec{N}\|_2 = \|\vec{L}\|_2 = 1$ $\cos \theta$

Sampling and Quantization

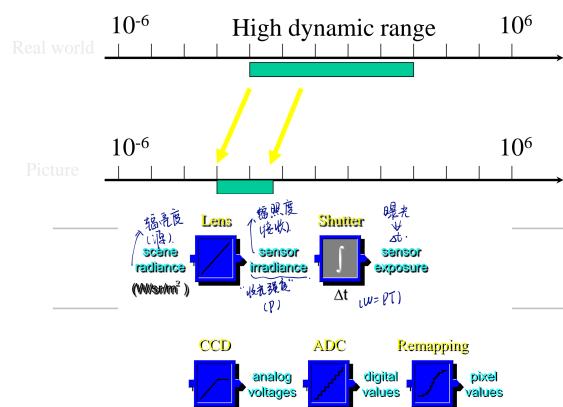


有了 Intensity 后, 用 pixel 来 sampling. 因为 real world 是连续的, 但 image, i.e., pixel space, 是离散的.

⇒ In fact, images can be regarded as signals! So sampling frequency (resolution) is important. We will talk about it later.

After getting intensity, have to map it to range $[0, 255]$.

And we are interested about how we can enhance it, like 'white balancing'.



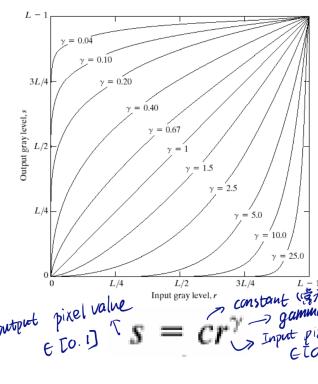


FIGURE 3.4 Plots of the function $s = cr^y$ for various values of y ($c = 1$ in all cases).

Method 1: Power - law transformations

$$S = cr^r, \quad r \in [0,1], \quad S \in [0,1]$$

Method 2: Negative : $S = 1 - r$
(Mostly Applied in Medical)

Method 3: Contrast Stretching

Method 4: Log : $S = c \cdot \log(1+r)$

Method 5: Image Histogram

$$S = 255 \times CDF(r)$$

$$r \in [0, 255], \quad r \in \mathbb{Z}$$

$$CDF(r) \in [0, 1]$$

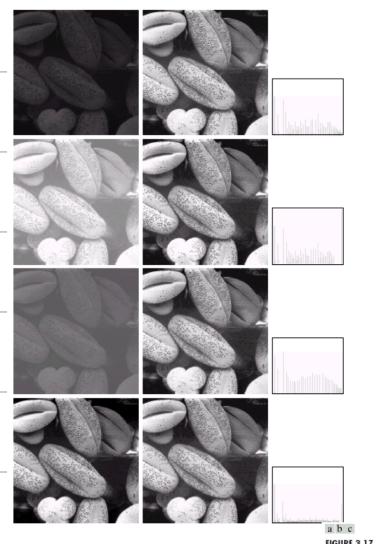


FIGURE 3.17

Now, back to sampling
stuff. According to
Nyquist Theorem

$$f_s > 2 \cdot f_{\max}$$

If failed:

Aliasing !!

How to ease aliasing ?



Sample more often

- Join the Mega-Pixel craze of the photo industry
- But this can't go on forever

Make the signal less "wiggly"

- Get rid of some high frequencies *
- Will loose information
- But it's better than aliasing

丢失部分信息

We can use Convolution

as filter to image ,

like a box

$$h[\cdot, \cdot]$$

filter :

$$\frac{1}{9} \begin{matrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{matrix}$$

⇒ We can get rid of high freq !

- Transformations on signals; e.g.:
 - bass/treble controls on stereo
 - blurring/sharpening operations in image editing
 - smoothing/noise reduction in tracking
- Key properties
 - linearity: $\text{filter}(f + g) = \text{filter}(f) + \text{filter}(g)$
 - shift invariance: behavior invariant to shifting the input
 - delaying an audio signal
 - sliding an image around
- Can be modeled mathematically by convolution

Cross-correlation

In Convolution, the operation sequence on entry is reversed
(Actually in DL, we don't flip)

Let F be the image, H be the kernel (of size $2k+1 \times 2k+1$), and G be the output image

$$G[i, j] = \sum_{u=-k}^k \sum_{v=-k}^k H[u, v]F[i + u, j + v]$$

This is called a **cross-correlation** operation:

$$G = H \otimes F$$

- Can think of as a "dot product" between local neighborhood and kernel for each pixel
= convolution if kernel is symmetric

Lec4: Convolution and Derivatives

Cross-correlation vs. Convolution

cross-correlation: $G = H \otimes F$

$$H \otimes F$$

$$G[i, j] = \sum_{u=-k}^k \sum_{v=-k}^k H[u, v]F[i + u, j + v]$$

A **convolution** operation is a cross-correlation where the filter is flipped both horizontally and vertically before being applied to the image:

(symmetric)

$$G[i, j] = \sum_{u=-k}^k \sum_{v=-k}^k H[u, v]F[i - u, j - v]$$

It is written:

$$G = H * F$$

$$H * F$$

A Gaussian kernel gives less weight to pixels further from the center of the window

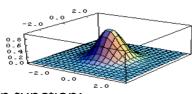
0	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0
0	0	0	90	90	90	90	90	90	0	0	0
0	0	0	90	90	90	90	90	90	0	0	0
0	0	0	90	90	90	90	90	90	0	0	0
0	0	0	90	90	90	90	90	90	0	0	0
0	0	0	90	90	90	90	90	90	0	0	0
0	0	0	90	90	90	90	90	90	0	0	0
0	0	0	90	90	90	90	90	90	0	0	0
0	0	0	90	90	90	90	90	90	0	0	0
0	0	0	90	90	90	90	90	90	0	0	0
0	0	0	90	90	90	90	90	90	0	0	0

$$F[x, y]$$

$$h(u, v) = \frac{1}{2\pi\sigma^2} e^{-\frac{u^2+v^2}{\sigma^2}}$$

$$\begin{matrix} 1 \\ 16 \end{matrix} \begin{matrix} 1 & 2 & 1 \\ 2 & 4 & 2 \\ 1 & 2 & 1 \end{matrix}$$

$$H[u, v]$$



This kernel is an approximation of a Gaussian function:

31

Slide by Steve Seitz

Removes "high-frequency" components from the image (low-pass filter)

Convolution with self is another Gaussian

$$\begin{matrix} \bullet & \bullet & \bullet \end{matrix} * \begin{matrix} \bullet & \bullet & \bullet \end{matrix} = \begin{matrix} \bullet & \bullet & \bullet \end{matrix}$$

Convolving twice with Gaussian kernel of width σ = convolving once with kernel of width $\sigma\sqrt{2}$

Gaussian (lowpass) pre-filtering



Gaussian 1/2

Solution: filter the image, then subsample

Filter size should double for each $\frac{1}{2}$ size reduction. Why?

Slide by Steve Seitz

Notation: $b = c * a$

Convolution is a multiplication-like operation

- commutative $a * b = b * a$
- associative $a * (b * c) = (a * b) * c$
- distributes over addition $a * (b + c) = a * b + a * c$
- scalars factor out $\alpha * a * b = a * \alpha b = \alpha(a * b)$
- identity: unit impulse $e = [..., 0, 0, 1, 0, 0, ...]$
- $a * e = a$

Conceptually no distinction between filter and signal

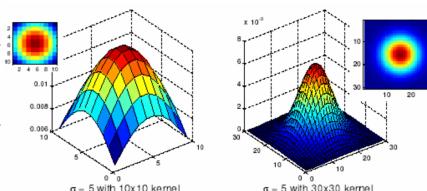
Usefulness of associativity

- often apply several filters one after another: $((a * b_1) * b_2) * b_3$
- this is equivalent to applying one filter: $a * (b_1 * b_2 * b_3)$

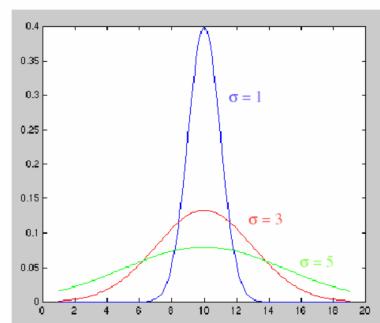
Instead of box filter, we can use

Gaussian to blur!

The Gaussian function has infinite support, but discrete filters use finite kernels



Effect of σ

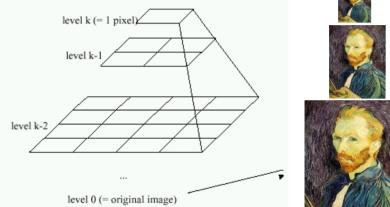


Consider: If we want to reduce the size of an image, we will see aliasing!

So we can blur it with Gaussian first, then down sample to reduce size.

Image Pyramids

Idea: Represent $N \times N$ image as a "pyramid" of $1 \times 1, 2 \times 2, 4 \times 4, \dots, 2^k \times 2^k$ images (assuming $N = 2^k$)



Known as a **Gaussian Pyramid** [Burt and Adelson, 1983]

- In computer graphics, a **mip map** [Williams, 1983]
- A precursor to **wavelet transform**

This forms an Image Pyramid

- The Fourier transform of the convolution of two functions is the product of their Fourier transforms

$$F[g * h] = F[g]F[h]$$

- The inverse Fourier transform of the product of two Fourier transforms is the convolution of the two inverse Fourier transforms

$$F^{-1}[gh] = F^{-1}[g]*F^{-1}[h]$$

- Convolution in spatial domain is equivalent to multiplication in frequency domain!

读懂 frequency domain 图 —— Spectrum

中心点 (DC 分量): 图正中心, 代表图像平均亮度

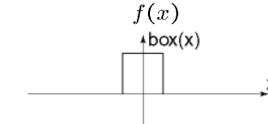
离中心距离: 越近, 代表频率越低

点亮度: 代表该频率在原始图中能量

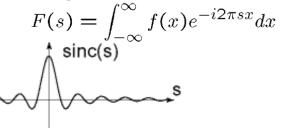
纹理: spectrum 中 水平亮线, 对应 figure 中 垂直方向边缘.

竖直, 对应 水平边缘

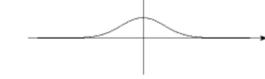
Spatial domain



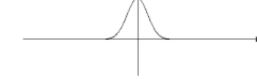
Frequency domain



gauss(x; σ)



gauss(s; 1/σ)



sinc(s)



box(x)

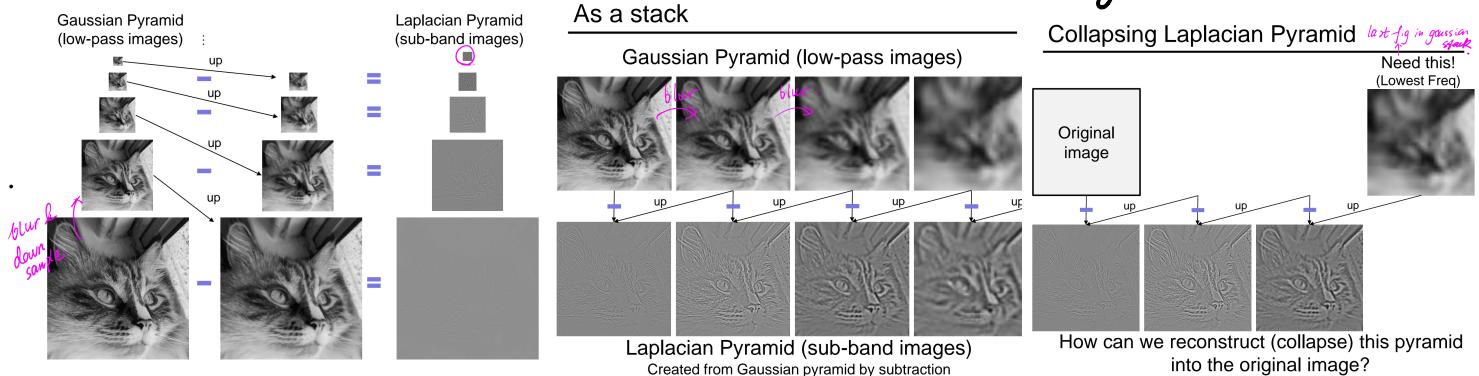


Why does the Gaussian give a nice smooth image, but the square filter give edgy artifacts?

This can be used to explain this
Box-filter Spatial Domain 上规整,
但 frequency domain 上不好, 故老积
时 Frequency domain 乘法(点乘)
会有很多 artifact ; Gaussian 反之 ; 甚至有 sincs) 这种低通滤波

Lec6 Pyramid Blending, Templates, NL Filters

下面将介绍如何完成上述三个任务 : ① Blending



Intro to how Gaussian & Laplacian Pyramids & Stack are built.

Then: Pipeline:

Note: Laplacian Stack Collapse
把所有的加起来即可 (迭代)

而 Pyramid 中还需要不断地:

$I_i = L_i$ it Up sample (I_{i-1}),
(“滚雪球”式地重建回原始图像)

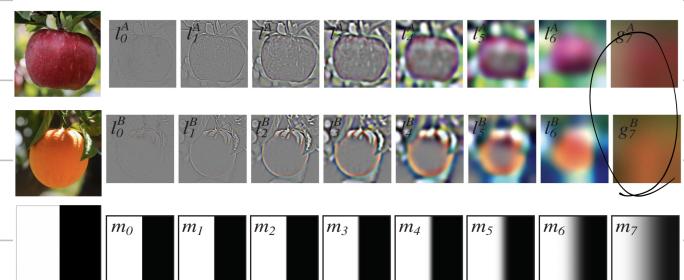
Image Blending with the Laplacian Pyramid

Build Laplacian pyramid for both images: LA, LB

Build Gaussian pyramid for mask: G

Build a combined Laplacian pyramid L

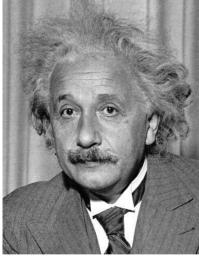
Collapse L to obtain the blended image



Template matching

Goal: find in image

- Main challenge: What is a good similarity or distance measure between two patches?
- Correlation
 - Zero-mean correlation
 - L2 distance
 - Normalized Cross Correlation



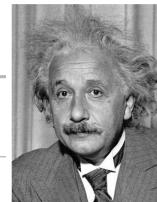
Goal: find in image

Method 1: filter the image with zero-mean eye

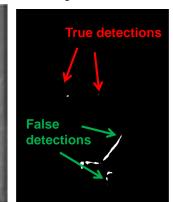
$$h[m, n] = \sum_{k, l} (g[k, l] - \bar{g}) f[m + k, n + l]$$

mean of g

Input



Filtered Image (scaled)



Thresholded Image

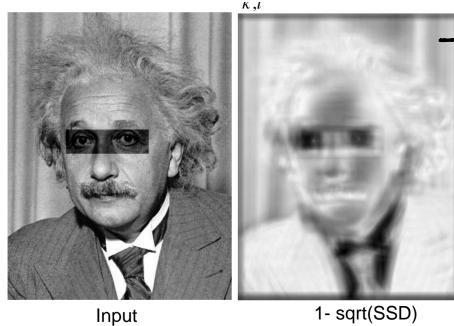
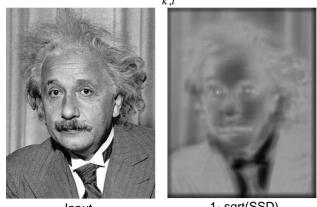
$(-\bar{g})$ is for
high contrast
increasing
(Kernel sum = 0)

② Template

Goal: find in image

Method 2: L2 distance (sum of squared diffs)

$$h[m, n] = \sum_{k, l} (g[k, l] - f[m + k, n + l])^2$$



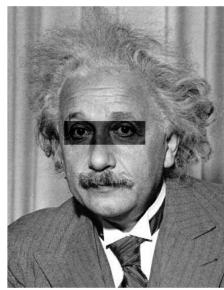
→ less sensitive to intensity

sensitive
⇒ Inspire correlation

Goal: find in image

Method 3: Normalized cross-correlation

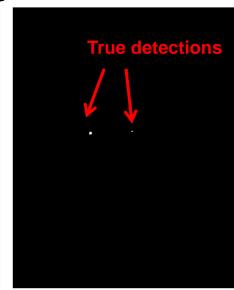
$$h[m, n] = \frac{\sum_{k, l} (g[k, l] - \bar{g})(f[m + k, n + l] - \bar{f}_{m, n})}{\left(\sum_{k, l} (g[k, l] - \bar{g})^2 \sum_{k, l} (f[m + k, n + l] - \bar{f}_{m, n})^2 \right)^{0.5}}$$



Input

Normalized X-Correlation

Good!



Thresholded Image

Zero-mean filter: fastest but not a great matcher

L2/SSD: next fastest, sensitive to overall intensity

Normalized cross-correlation: slowest, invariant to local average intensity and contrast



Summary to different methods.

Reducing Gaussian noise

③ Denoise.

Gaussian can mitigate, but can't eliminate.

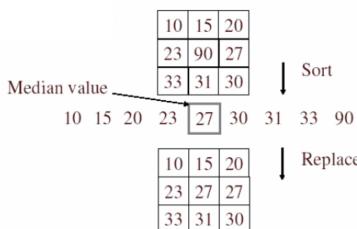
Also it can blur image at the same time!

⇒ Solution: "median" filter:

↓ Median vs. Gaussian filtering

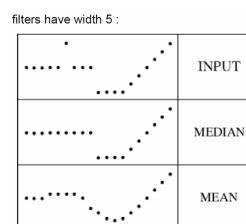
Alternative idea: Median filtering

A median filter operates over a window by selecting the median intensity in the window



What advantage does median filtering have over Gaussian filtering?

Robustness to outliers



(NL)

Gaussian

(High contrast)

Take lots of time,

not linear

Median

but good !!



3x3



5x5



7x7



3x3



5x5



7x7

Smoothing with larger standard deviations suppresses noise, but also blurs the image



Lec 7. Image Transformation + Warp

Previous image filtering center on change range of image.
now warp focus on change domain of image

What a 2×2 matrix can do ...

2-D Rotation

This is easy to capture in matrix form:

$$\begin{bmatrix} x' \\ y' \end{bmatrix} = \underbrace{\begin{bmatrix} \cos(\theta) & -\sin(\theta) \\ \sin(\theta) & \cos(\theta) \end{bmatrix}}_{R} \begin{bmatrix} x \\ y \end{bmatrix}$$

Is this a linear transformation? Yes

Even though $\sin(\theta)$ and $\cos(\theta)$ are nonlinear functions of θ ,

- x' is a linear combination of x and y
- y' is a linear combination of x and y

What is the inverse transformation?

- Rotation by $-\theta$ $R^{-1} = R^T$
- For rotation matrices

旋转矩阵

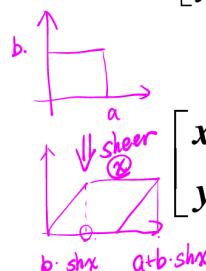
2D Rotate around (0,0)?

$$\begin{aligned} x' &= \cos \Theta * x - \sin \Theta * y \\ y' &= \sin \Theta * x + \cos \Theta * y \end{aligned}$$

$$\begin{bmatrix} x' \\ y' \end{bmatrix} = \begin{bmatrix} \cos \Theta & -\sin \Theta \\ \sin \Theta & \cos \Theta \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix}$$

2D Shear?

$$\begin{aligned} x' &= x + sh_x * y \\ y' &= sh_y * x + y \end{aligned}$$



$$\begin{bmatrix} x' \\ y' \end{bmatrix} = \begin{bmatrix} 1 & sh_x \\ sh_y & 1 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix}$$

Linear transformations are combinations of ...

- Scale,
- Rotation,
- Shear, and
- Mirror/Reflection

$$\begin{bmatrix} x' \\ y' \end{bmatrix} = \begin{bmatrix} a & b \\ c & d \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix}$$

Properties of linear transformations:

- Origin maps to origin
- Lines map to lines
- Parallel lines remain parallel
- Ratios are preserved
- Closed under composition

$$\begin{bmatrix} x' \\ y' \end{bmatrix} = \begin{bmatrix} a & b & e & f & i & j \\ c & d & g & h & k & l \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix}$$

Affine transformations are combinations of ...

- Linear transformations, and
- Translations

$$\begin{bmatrix} x' \\ y' \\ w \end{bmatrix} = \begin{bmatrix} a & b & c \\ d & e & f \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x \\ y \\ w \end{bmatrix}$$

exactly the
def of affine
transformation

Properties of affine transformations:

- Origin does not necessarily map to origin (translation!)
- Lines map to lines
- Parallel lines remain parallel
- Ratios are preserved
- Closed under composition
- Models change of basis

Will the last coordinate w always be 1? Yes!

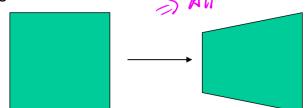
Projective transformations ...

- Affine transformations, and
- Projective warps

$$\begin{bmatrix} x' \\ y' \\ w' \end{bmatrix} = \begin{bmatrix} a & b & c \\ d & e & f \\ g & h & i \end{bmatrix} \begin{bmatrix} x \\ y \\ w \end{bmatrix}$$

have more freedom
to allow perspectives

- Parallel lines do not necessarily remain parallel
- Ratios are not preserved
- Closed under composition
- Models change of basis



⇒ Affine can't do this.

Only linear 2D transformations
can be represented with a 2×2 matrix

But can't bring in Translation

⇒ 3×3 ? ⇒ with additional 1

if last row: 0 0 1

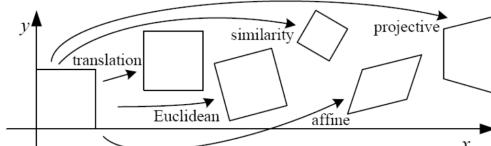
Then it is affine transformation



But if it's not:

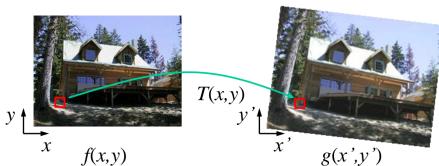
its perspective Transformations!

2D image transformations



Name	Matrix	# D.O.F.	Preserves:	Icon
translation	$[I t]_{2 \times 3}$	2		□
rigid (Euclidean)	$[R t]_{2 \times 3}$	3		◇
similarity	$[sR t]_{2 \times 3}$	4		△
affine	$[A]_{2 \times 3}$	6		□
projective	$[\tilde{H}]_{3 \times 3}$	8	straight lines only	□

So how to warp? First: $f(\text{source}) = \text{target} \dots$



Send each pixel $f(x,y)$ to its corresponding location

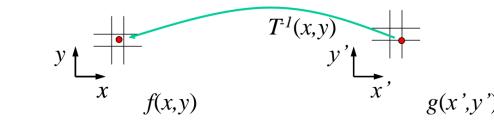
$$(x',y') = T(x,y) \text{ in the second image}$$

Q: what if pixel lands "between" two pixels?

A: distribute color among neighboring pixels (x',y')

- Known as "splatting" (Check out griddata in Matlab)
- Generally, a very bad idea. Why? Loss of detail

Maybe we can try: $f^{-1}(\text{target}) = ? \text{ source} \dots$



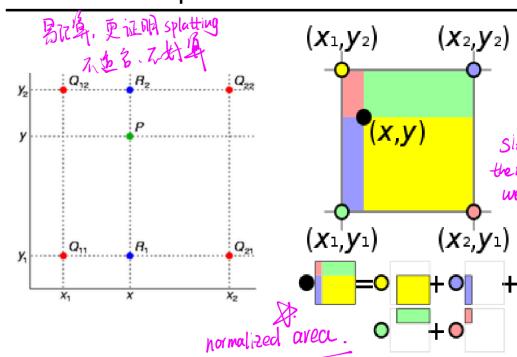
Get each pixel $g(x',y')$ from its corresponding location
 $(x,y) = T^{-1}(x',y')$ in the first image

Q: what if pixel comes from "between" two pixels?

A: Interpolate color value from neighbors

- nearest neighbor, bilinear, Gaussian, bicubic
- Check out interp2 in Matlab / Python

Bilinear Interpolation

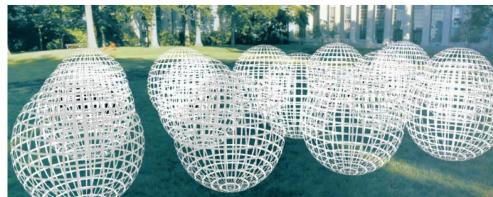


Use bilinear

Interpolation
for deciding
numerical value.

Lec 8 Mosaic

The Plenoptic Function



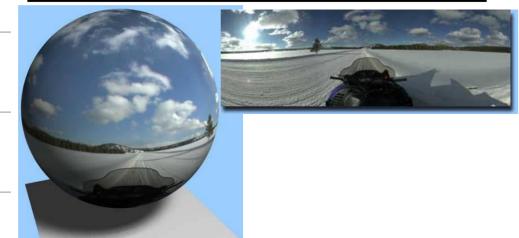
$$P(\theta, \phi, \lambda, t, V_x, V_y, V_z)$$

- Can reconstruct every possible view, at every moment, from every position, at every wavelength
- Contains every photograph, every movie, everything that anyone has ever seen! it completely captures our visual reality! Not bad for a function...

Plenoptic Function (全光函数) 描述了从一个特定观察点，在任意时间和方向上能看到的任意波长光的强度

Let's think our
vision plane is a ball.
receiving light at every

Spherical Panorama



All light rays through a point form a panorama
Totally captured in a 2D array -- $P(\theta, \phi)$
Where is the geometry???

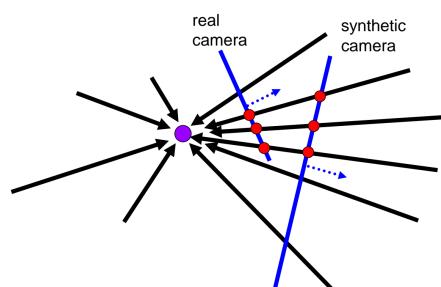
A pencil of rays contains all views

angle \Rightarrow Spherical Panorama!

以该视角审视照片 image 的形成，深度信息
都可抛去。^x

As long as COP don't move, we can
synthesize any camera view whose
COP is identical

* (深度信息需要多角度视角以获取)

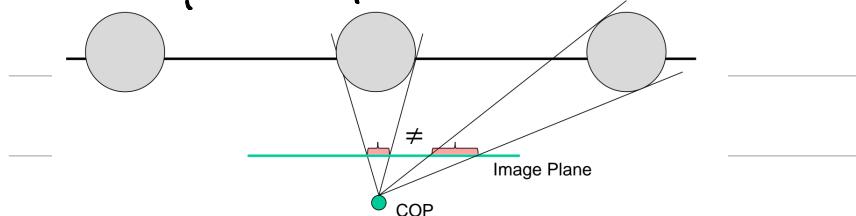


Can generate any synthetic camera view
as long as it has the same center of projection!

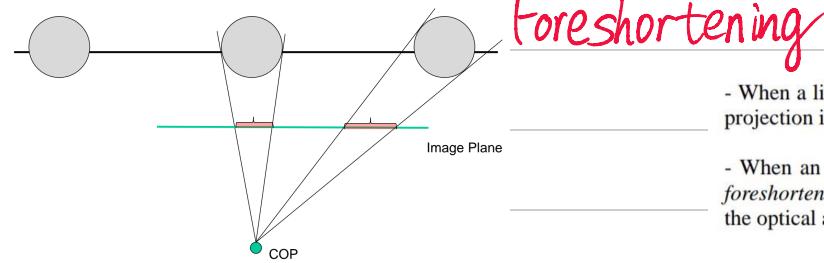
How to relate two images with the same COP? Find corresponding points and warp!

⇒ Homography

Another topic: Perspective Distortion:



Less noticeable with long focal length (i.e. you see distortion more with wide-angle camera)

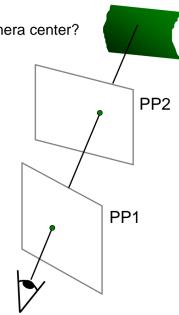


Scale & Foresighting

Image reprojection

Basic question

- How to relate two images from the same camera center?
- how to map a pixel from PP1 to PP2



Answer

- Cast a ray through each pixel in PP1
- Draw the pixel where that ray intersects PP2

But don't we need to know the geometry of the two planes in respect to the eye?

Observation:

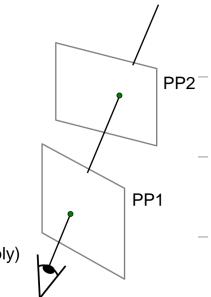
Rather than thinking of this as a 3D reprojection, think of it as a 2D **image warp** from one image to another

A: Projective – mapping between any two PPs with the same center of projection

- rectangle should map to arbitrary quadrilateral
- parallel lines aren't
- but must preserve straight lines
- same as: unproject, rotate, reproject

called Homography

$$\begin{bmatrix} wx' \\ wy' \\ w \\ p' \end{bmatrix} = \begin{bmatrix} * & * & * \\ * & * & * \\ * & * & * \\ H \end{bmatrix} \begin{bmatrix} x \\ y \\ 1 \\ p \end{bmatrix}$$

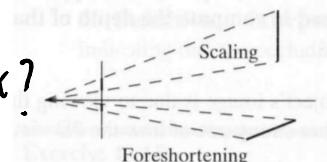


To apply a homography H

- Compute $p' = Hp$ (regular matrix multiply)
- Convert p' from homogeneous to image coordinates

- When a line (or surface) is parallel to the image plane, the effect of perspective projection is **scaling**.

- When a line (or surface) is not parallel to the image plane, we use the term **foreshortening** to describe the projective distortion (i.e., the dimension parallel to the optical axis is compressed relative to the frontal dimension).



At last: How to work out Homography Matrix?

Image rectification



- Need Point Pairs

- Find Optimal matrix.

H 's DOF = 8

How to Compute? Suppose $h_{33}=1$:

For one point: $(x_p, y_p, 1) \xrightarrow{H} (x, y, 1)$

$$A = \begin{bmatrix} 1 & y & 1 & 0 & 0 & 0 & -x \cdot x_p & -y \cdot x_p \\ 0 & 0 & 0 & x & y & 1 & -x \cdot y_p & -y \cdot y_p \end{bmatrix}$$

可见 1 point $\Rightarrow 2$ data $\therefore 3$ points fully determine

≥ 4 points overconstrained: $\min_H \|A\bar{H} - b\|_2^2$

$$H^* = (A^T A)^{-1} A^T b \quad (\text{close form solution})$$

$$\begin{bmatrix} p' \\ w \\ 1 \end{bmatrix} = \begin{bmatrix} a & b & c \\ d & e & f \\ g & h & i \end{bmatrix} \begin{bmatrix} x \\ y \\ 1 \end{bmatrix}$$

Can set scale factor $i=1$. So, there are 8 unknowns.
Set up a system of linear equations:

$$Ah = b$$

where vector of unknowns $h = [a, b, c, d, e, f, g, h]^T$

Need at least 8 eqs, but the more the better...

Solve for h . If overconstrained, solve using least-squares:

$$\min \|Ah - b\|^2$$

$$h = \begin{bmatrix} h_{11} \\ h_{12} \\ h_{13} \\ \vdots \\ h_{31} \\ h_{32} \end{bmatrix} \quad b = \begin{bmatrix} x_p \\ y_p \end{bmatrix}$$

Lec 9 & 10: How to auto - find matching points?

Not realistic

- Search in $O(N^8)$ is problematic
- Not clear how to set starting/stopping value and step

What can we do?

- Use pyramid search to limit starting/stopping/step values

Alternative: gradient decent on the error function

- i.e. how do I tweak my current estimate to make the SSD error go down?
- Can do sub-pixel accuracy
- BIG assumption?
 - Images are already almost aligned (<2 pixels difference!)
 - Can improve with pyramid
- Same tool as in motion estimation

~~Iterate through them and find best H.~~

But:

Can be improved with pyramid & GI

Feature-based alignment

- Feature Detection:** find a few important features (aka Interest Points) in each image separately
- Feature Matching:** match them across two images
- Compute image transformation:** as per Project 4, P

Suitable way: Find feature points, and then match. How?

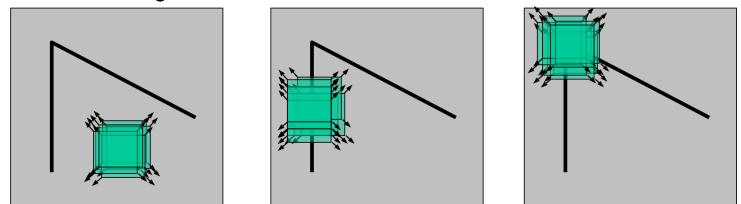
How do we match the features between the images?

- Need a way to describe a region around each feature
 - e.g. image patch around each feature
- Use successful matches to estimate homography
 - Need to do something to get rid of outliers

Issues:

- What if the image patches for several interest points look similar?
 - Make patch size bigger
- What if the image patches for the same feature look different due to scale, rotation, etc.
 - Need an invariant descriptor

$$E(u, v) = \sum_{(x, y) \in W} [I(x+u, y+v) - I(x, y)]^2$$



If u, v are small (small shift):)

- First-order Taylor approximation for small motions $[u, v]$:

$$\begin{aligned} I(x+u, y+v) &= I(x, y) + I_x u + I_y v + \text{higher order terms} \\ &\approx I(x, y) + I_x u + I_y v \end{aligned}$$

plug in

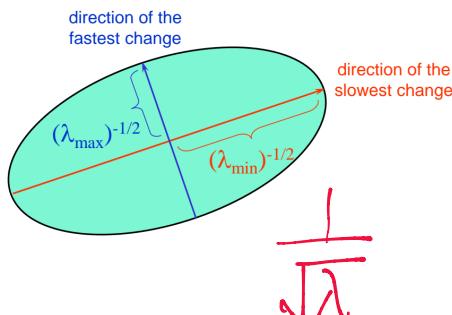
$$= I(x, y) + \begin{bmatrix} I_x & I_y \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix}$$

Consider a horizontal "slice" of $E(u, v)$: $[u \ v] M \begin{bmatrix} u \\ v \end{bmatrix} = \text{const}$

This is the equation of an ellipse.

$$\text{Diagonalization of } M: M = R^{-1} \begin{bmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{bmatrix} R$$

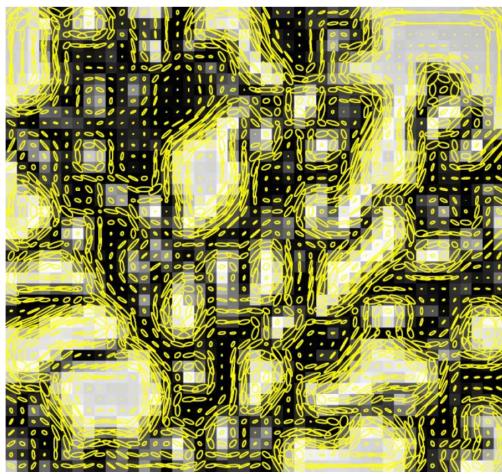
The axis lengths of the ellipse are determined by the eigenvalues and the orientation is determined by R



$$\begin{aligned} E(u, v) &= \sum_{(x, y) \in W} [I(x+u, y+v) - I(x, y)]^2 \\ &\approx \sum_{(x, y) \in W} [I(x, y) + \begin{bmatrix} I_x & I_y \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix} - I(x, y)]^2 \\ &= \sum_{(x, y) \in W} \left(\begin{bmatrix} I_x & I_y \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix} \right)^2 M \\ &= \sum_{(x, y) \in W} [u \ v] \begin{bmatrix} I_x^2 & I_x I_y \\ I_x I_y & I_y^2 \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix} \end{aligned}$$

Diagonalize M & get eigenvectors.

Can use ellipse to visualize!



\Rightarrow Gradient \uparrow , axis \downarrow
 \Rightarrow Good to reflect "change intensity"

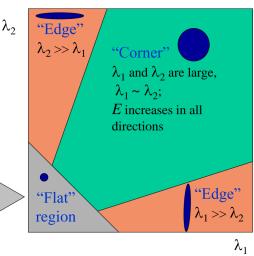
Measure of corner response:

$$R = \frac{\det M}{\text{Trace } M}$$

$$\det M = \lambda_1 \lambda_2$$

$$\text{trace } M = \lambda_1 + \lambda_2$$

Interpreting the eigenvalues
 Classification of image points using eigenvalues of M :



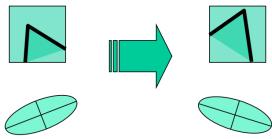
$$M = \sum_{(x,y) \in W} \begin{bmatrix} I_x^2 & I_x I_y \\ I_y I_x & I_y^2 \end{bmatrix}$$

Pipeline shown on the right:

① Compute each point's I_x, I_y

② For each (\hat{x}, \hat{y}) , draw window W . Use results in ① to compute M

Some property of harris detector :

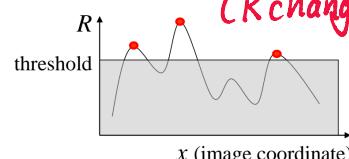


Ellipse rotates but its shape (i.e. eigenvalues) remains the same

Corner response R is invariant to image rotation

✓ Only derivatives are used \Rightarrow invariance to intensity shift $I \rightarrow I + b$

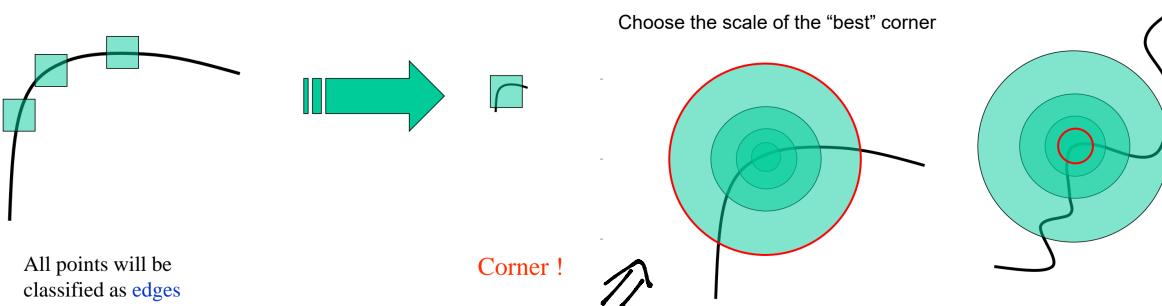
✓ Intensity scale: $I \rightarrow aI$



\Rightarrow Partial invariance to affine intensity change.

The problem: how do we choose corresponding circles independently in each image?

Choose the scale of the "best" corner



All points will be classified as edges

在所有可能尺度中，使角点响应函数 R 达到局部最大的尺度

Also, we want points are equally distributed ! \Rightarrow ANMS.

For point set P , for $p_i \in P$, compute r_i :

$$r_i = \min_{P_j \in P, R_j > c \cdot R_i} \|P_j - P_i\|_2, \text{ then sort using } r_i \text{ and take top-} p \text{ points,}$$

② Match Harris Corners \Rightarrow Invariant & Distinctive

Solution: (In project) Take 40×40 patch \rightarrow Blur

\rightarrow Downsample to $8 \times 8 \rightarrow$ Normalize: $z = \frac{x - u}{\sigma}$

Then to match them: Similarity between two descriptors: Euclidean

And want: $\frac{1-NN}{2-NN} < \varepsilon$, to ensure match is correct

Finally, to get rid of outlier ...

RANSAC Algorithm

*: Use all inliers! Not its original H ...

RANSAC loop:

1. Select four feature pairs (at random)
2. Compute homography H (exact)
3. Compute inliers where $dist(p_i', H p_i) < \varepsilon$
4. Keep largest set of inliers
5. Re-compute least-squares H estimate on all of the inliers

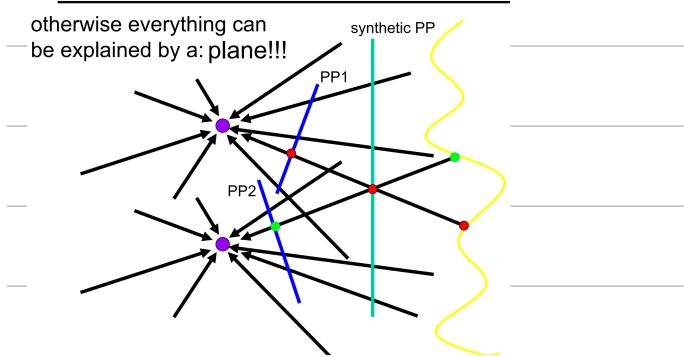
Lec 11 & 12: 3D Intro & Stereo

- is 3D = depth from a single image?

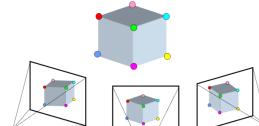
Structure and depth are inherently ambiguous from single views.

- 2.5D = per-pixel depth from a single image

Need two different camera centers



3D Points (Structure)



Correspondences



Camera (Motion)



There is a world coordinate frame and camera looking at the world

How can we model the geometry of a camera?

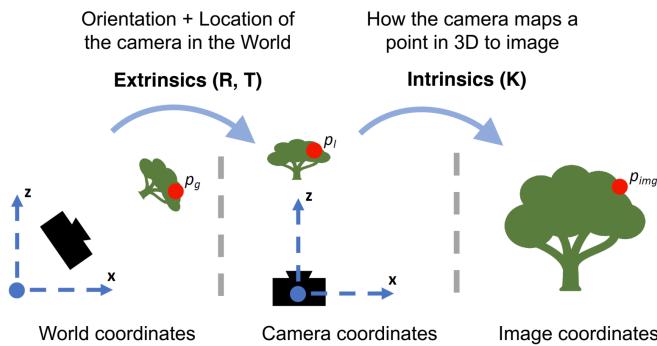


Three important coordinate systems:

1. World coordinates
2. Camera coordinates
3. Image coordinates

To go from pixels to 3D location in the **world coordinates**, we need to know two things about the camera:

1. Position & Orientation of the camera with respect to the world (extrinsics)
2. How the camera maps a point in the world to image (intrinsic)



Want:

Image Coordinates

$$\mathbf{x}_i = \begin{bmatrix} x_i \\ y_i \end{bmatrix} \quad \xleftarrow{\text{Perspective Projection (3D to 2D)}} \quad \mathbf{x}_c = \begin{bmatrix} x_c \\ y_c \\ z_c \end{bmatrix}$$

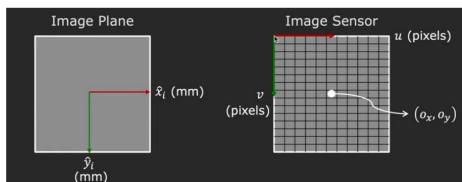
Camera Coordinates

World Coordinates

$$\mathbf{x}_w = \begin{bmatrix} x_w \\ y_w \\ z_w \end{bmatrix} \quad \xleftarrow{\text{Coordinate Transformation (3D to 3D)}} \quad \mathbf{x}_c$$

Coordinate Transformation (3D to 3D)

Image Plane to Image Sensor Mapping



1. Account for pixel density (pixel/mm) & aspect ratio by scalars: $[m_x, m_y]$
 $m_x x_i, m_y y_i$

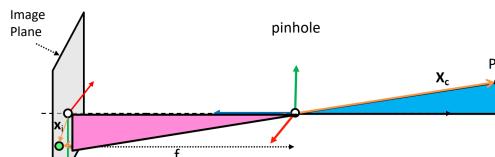
2. Usually the top left corner is the origin. But in the image plane, the origin is where the optical axis pierces the plane! Need to shift by: (o_x, o_y)

$$u_i = \alpha_x x_i + o_x = \alpha_x f \frac{x_c}{z_c} + o_x$$

where $[f_x, f_y] = [m_x f, m_y f]$

Pixel Coordinates:

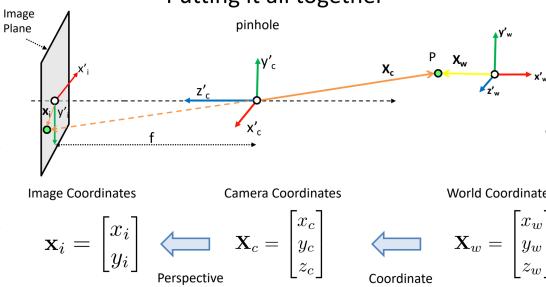
$$u_i = f_x \frac{x_c}{z_c} + o_x \quad v_i = f_y \frac{y_c}{z_c} + o_y$$



$$\mathbf{x}_i = \begin{bmatrix} x_i \\ y_i \end{bmatrix} \quad \frac{x_i}{f} = \frac{x_c}{z_c} \quad \mathbf{x}_c = \begin{bmatrix} x_c \\ y_c \\ z_c \end{bmatrix}$$

$$\text{Image Coordinates} \quad x_i = f \frac{x_c}{z_c}$$

Putting it all together



$$M = \begin{bmatrix} f_x & 0 & o_x & 0 \\ 0 & f_y & o_y & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix} \times \begin{bmatrix} R_{3 \times 3} & \mathbf{t} \\ 0_{1 \times 3} & 1 \end{bmatrix}$$

mx, my : mm \rightarrow pixel

Image plane: origin : center!

Image Sensor : Top-Left point!

Intrinsic

(with homogeneous)

Perspective projection + Transformation to Pixel Coordinates:

$$u_i = f_x \frac{x_c}{z_c} + o_x \quad v_i = f_y \frac{y_c}{z_c} + o_y$$

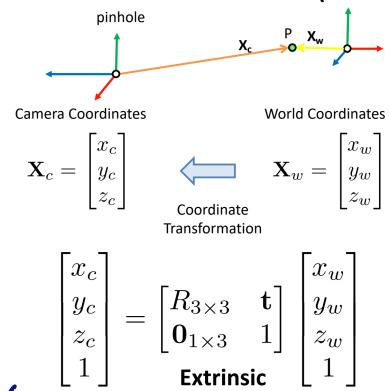
$$\begin{bmatrix} u \\ v \\ 1 \end{bmatrix} \equiv \begin{bmatrix} \tilde{u} \\ \tilde{v} \\ \tilde{w} \end{bmatrix} = \begin{bmatrix} f_x & 0 & o_x & 0 \\ 0 & f_y & o_y & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} x_c \\ y_c \\ z_c \end{bmatrix}$$

Intrinsic Matrix

Extrinsic

(need homogeneous
for translation).

Camera Transformation (3D-to-3D)



M's DOF:

$$5 + 6 = 11$$

Intrinsic (4+1)

Extrinsic (6)

* Skew: Some camera's

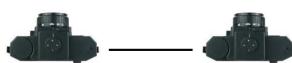
pixel is parallelogram. Need skew.

Now we can use multi-view to sense depth.

From now, suppose cameras are calibrated (In/Extrinsic Matrix Unknown)

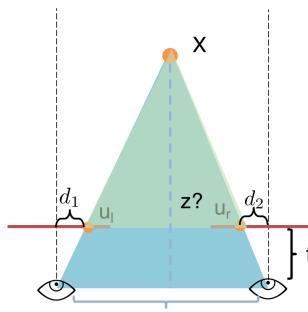


- Assume parallel optical axes
- Two cameras are calibrated
- Find relative depth



Key Idea: difference in corresponding points to understand shape

Solving for Depth in Simple Stereo



Do we have enough to know what is Z?

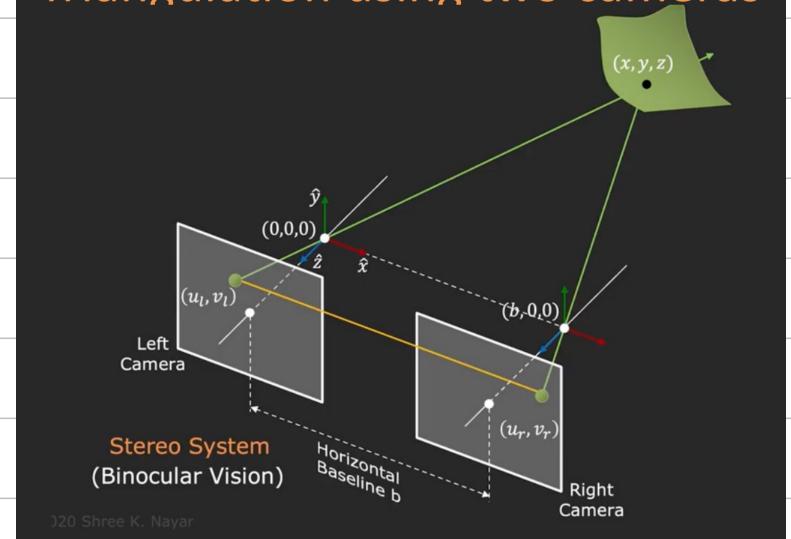
Yes, similar triangles!

$$\triangle \frac{B - (u_l - u_r)}{z - f} = \frac{B}{z}$$

$$z = \frac{fB}{u_l - u_r}$$

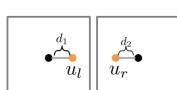
disparity (how much corresp. pixels move)

Triangulation using two cameras

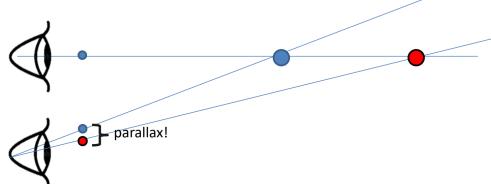


$Z(\text{depth}) \propto \frac{1}{u_l - u_r}$ ← disparity

Base of \triangle : $B - (d_1 + d_2)$
in image coordinates: $= B - (u_l - u_r)$



Parallax



Parallax = from ancient Greek *parállaxis*
= *Para* (side by side) + *allássō*, (to alter)
= Change in position from different view point

Two eyes give you parallax, you can also move to see more
parallax = "Motion Parallax"

视差 : disparity caused due to two eyes.

Consider task:

Stereo Matching: Finding Disparities

Goal: Find the disparity between left and right stereo pairs.



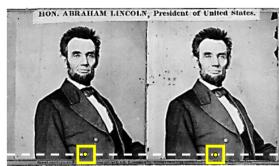
Left/Right Camera Images



Disparity Map (Ground Truth)

So at least we want to
find pairs of points who has same disparity:

Your basic stereo algorithm



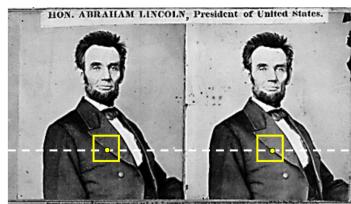
For every epipolar line:

For each pixel in the left image

- compare with every pixel on same epipolar line in right image
- pick pixel with minimum match cost

Improvement: match *windows*, + clearly lots of matching strategies

Dense correspondence search



For each epipolar line

For each pixel / window in the left image

- compare with every pixel / window on same epipolar line in right image
- pick position with minimum match cost (e.g., SSD, correlation)

Epipolar line:

对于左边一个点，我们知道必
在右图中一条特定直线上
有对应点

(parallel movements)

Issue: Effect of window size:

Want window large enough to have sufficient intensity variation, yet small enough to contain only pixels with about the same disparity.

Issues with Stereo

- Surface must have non-repetitive texture



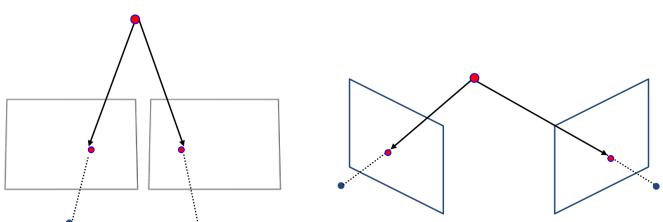
- Foreshortening effect makes matching a challenge



Slide Credit: Shn

More general case: Cameras are calibrated
but doesn't have to have parallel optical axes

- The two cameras need not have parallel optical axes.



Option 2

1. Solve for correspondences

2. Estimate camera

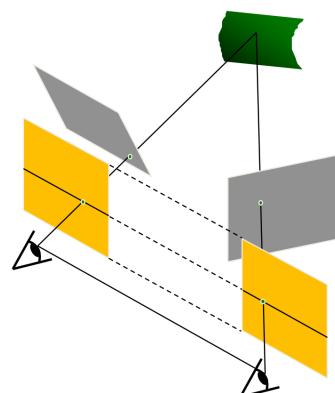
- What is the relationship between the camera + correspondences?

3. Triangulate

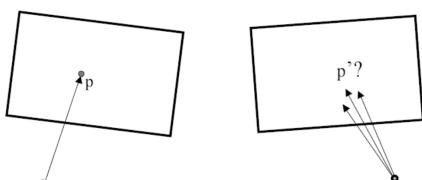
Option 1: Rectify via homography Quite Familiar

- reproject image planes onto a common plane
 - plane parallel to the line between optical centers
- pixel motion is horizontal after this transformation
- two homographies, one for each input image reprojection

– C. Loop and Z. Zhang. [Computing Rectifying Homographies for Stereo Vision](#). CVPR 1999.



Issue:



- Given p in left image, where can corresponding point p' be?

To be continued in future class