

# Computer Vision

- Goal: enables machine to understand images & videos
  - { measurement: compute properties of 3D World (distance, shape...)
  - { perception & interpretation: recognize objects, people, activities..
- Course outline:
  - 1) Image Processing Basic
  - 2) Segmentation
  - 3) Local features & Matching
  - 4) Object Recognition & Categorization
  - 5, 3D Reconstruction
  - 6, Deep Learning

  
Ng Huu Duc

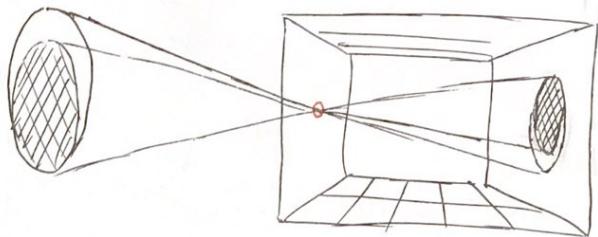
# Image Formation

Log

+ Camera obscura:

"Dark Chamber,"

Leonardo Da Vinci 1544



+ Pinhole Camera:

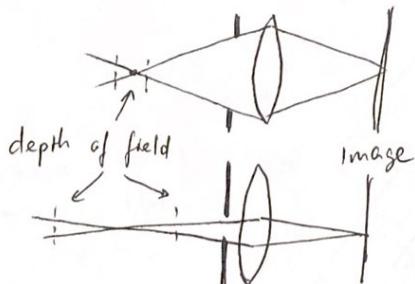
- Pinhole size = Aperture

{ too big - blurring  
too small - also blur, but because of diffraction  
but then, image is dark

⇒ Lenses: keep image sharp while capture more light

- The thin lens

- Focus & Depth of Field:



Large Aperture: small depth of field  
(only object with the correct distance will be at focus, while background is blurry)

Small Aperture: large depth of field  
but need more light

-  $f \propto$  field of view:  $f$  gets smaller, image ⇒ wide range  
 $f$  gets larger, — telescopic

+ Digital image: Discretize the image into pixels

Quantize light intensities ⇒ pixel values

+ Resolution: No. of pixel (most commonly understand)

+ Color Sensing: Bayer grid - demosaicing

Color image: RGB

just 1 of many color space

Luv  
xyz ..

Grey-scale Image

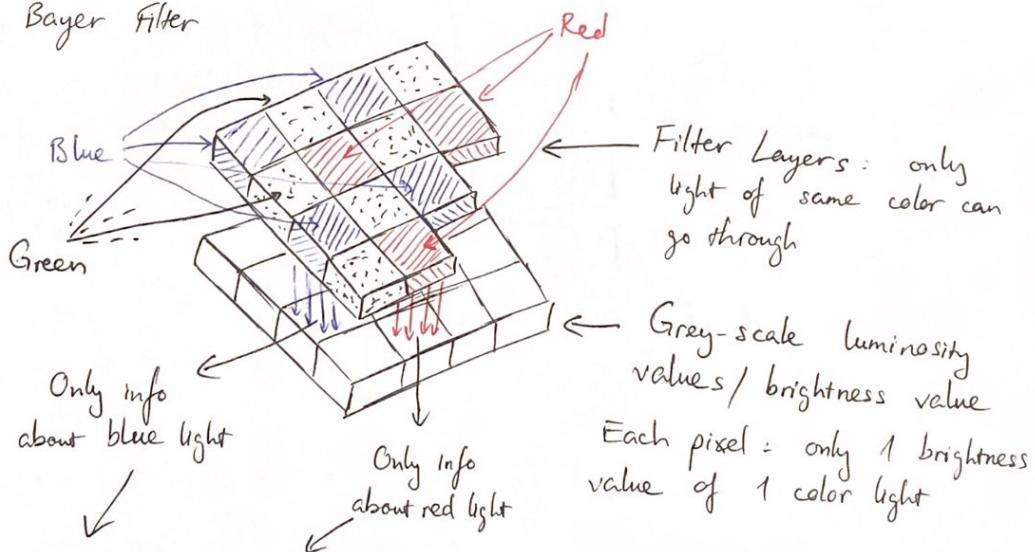
### ④ Demosaicing

Digital camera takes in ~~image~~ light through a filter (Bayer or Xtrans)

⇒ we get a grey-scale image

We need to do Demosaicing based on the filter's pattern to get the color image from the raw image

Example: Bayer Filter



⇒ Use brightness values of surrounding pixels about other 2 colors (a.k.a. demosaicing)

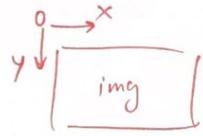
④ Raw image has a green-cast

Twice many green as red & blue cause human eyes are twice as sensitive to the green part to other red / blue part.

# Image Processing

Locy

## 1. Linear Filters



- Type of noise :
  - Salt & pepper noise
  - Impulse noise
  - Gaussian noise

$$\text{noise} = \text{randn}(\text{size}(im)) * \sigma$$

$$\text{output} = im + \text{noise}$$

Basic Assumption: i.i.d - independent & identically distributed

- Correlation Filter:  $G[i,j] = \frac{1}{(2k+1)^2} \sum_{u=-k}^k \sum_{v=-k}^k F[i+u, j+v]$

Matlab

filter2  
imfilter

$$G = H \otimes F$$

$k$  is from the window size different weights:  $G[i,j] = \sum_{u=-k}^k \sum_{v=-k}^k H[u,v] \underbrace{F[i+u, j+v]}_{\text{non-uniform weights}}$

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

$$\boxed{0.0} F_{NN}$$

If  $H[u,v] = H[-u,-v]$   $\Rightarrow$  correlation  $\equiv$  convolution

conv2

$$G = H * F$$

Averaging Filter: Ringing Artifacts ?

- Gaussian Filter :  $\frac{1}{\sqrt{2\pi} \sigma} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$

more noise  $\Rightarrow$  ↑ $\sigma$

Rule of thumb: set filter width to 6 $\sigma$

$\Rightarrow$  blurring effect

Efficient Implementation: if Filter is separable  $\Rightarrow$  do 1D 2 times to have 20 computational cost: from  $O(K^2)$  to  $O(2K)$

K: Kernel size

- Boundary issues :
  - full
  - same
  - valid

$$\boxed{\text{output size}} = f + g$$

$$\boxed{\quad} = f$$

$$\boxed{\quad} = f - g$$

$$\boxed{\quad} = g$$

Pixel near boundary : 
 

- Clip filter (black)  $\Rightarrow$  dark border
- Wrap around
- Copy edge  $\rightarrow$  strong edge response
- Reflect across edge

## 2. Back-ground

- Taken Fourier Transform of a signal  $\Rightarrow$  Frequency coefficients  
 $\Rightarrow$  Frequency Spectrum

④ Duality: The better a func is localized in 1 domain the worse it is localized in the other

- Effect of Convolution:  $f * g \rightarrow F \cdot G$

A Gaussian has compact support in both domain

$\Rightarrow$  convenient choice for low-pass filter

[low pass filter  $\Rightarrow$  smooth]

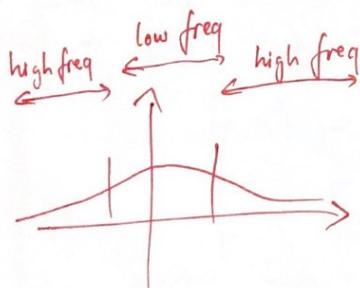
[high pass filter  $\Rightarrow$  edge, line]

- Sharpening filter: emphasize noise as well

## 3. Non - Linear Filter

- Median Filter: [remove spikes (good for impulse, salt & pepper noise)  
 edge preserving (unlike mean filter)]

If we increase filter size of Median Filter  $\Rightarrow$  reducing structure loose details



freq domain  
(Fourier)

### ④ Correlation vs Convolution

Conv is better, it has additional nice properties

- commutative:  $f * g = g * f$
- associative  $(f * g) * h = f * (g * h)$
- Fourier transform  $f * g \xrightarrow{\text{FT}} F \cdot G$

Both are linear shift invariant (LSI)

$$h \circ (f_0 + f_1) = h \circ f_1 + h \circ f_0$$

$$f \cdot h \xrightarrow{\text{FT}} F \cdot H$$

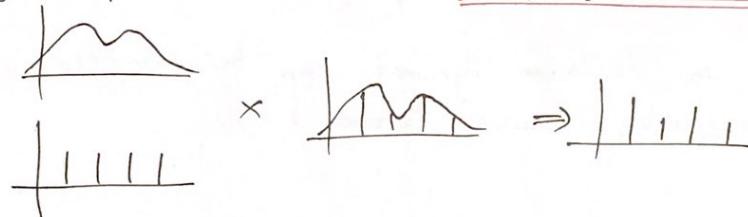
L03,

## 4. Multi-Scale representations

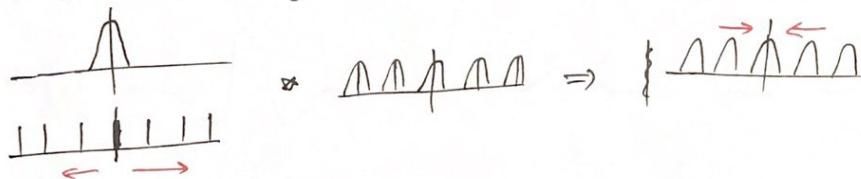
- Image Pyramid: very little over-head (in term of computational cost)

### ② Fourier Interpretation: Discrete Sampling:

- Sampling in spatial domain is like multiplying with a spike func



- ⇒ Sampling in frequency domain is like convolving with a spike func



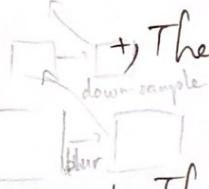
⇒ When we sampling with lower frequency, the magnitude spectrum will be overlaped (due to duality) ⇒ will not be able to reconstruct signal

+)Nyquist theorem and limit: to recover a certain freq  $f$ , at least sample with  $2f$

/by lighting/ ⇒ Aliasing artifacts in Graphics: overlapped (because sampling with edges? too low freq)

③ We can not recover high frequencies, but we can avoid artifacts by prior smoothing before resampling

with impulse img  
conv reproduces itself  
corr reflect itself

+ The Gaussian Pyramid: perform blurring  $\Rightarrow$  then down-sampling  


+ The Laplacian Pyramid:  
 $L_i = G_i - \text{expand}(G_{i+1})$   
 $G_i = L_i + \text{expand}(G_{i+1})$   
 $L_n = G_n$   
 $\Rightarrow L_{0 \rightarrow n-1}$  contain high frequency information

⑧ Images in Laplacian Pyramid can be compressed further than the corresponding Gaussian Pyramid images

Laplace Laplacian  $\sim$  Difference of Gaussians

  
 the mask  $\Leftrightarrow$  ?? a combinations of 2nd derivatives  $\approx$  edges  
 Laplace & from

Laplacian : 
$$\nabla^2 f = \frac{\partial^2 f}{\partial x^2} + \frac{\partial^2 f}{\partial y^2}$$

$$\begin{aligned} \frac{\partial^2 f}{\partial x^2} &= [f(x+1, y) - f(x, y)] - [f(x, y) - f(x-1, y)] \\ &= f(x+1, y) + f(x-1, y) - 2f(x, y) \end{aligned}$$

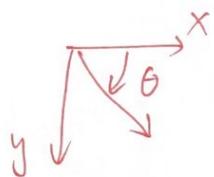
$$\Rightarrow \nabla^2 f = f(x \pm 1, y) + f(x, y \pm 1) - 4f(x, y)$$

$\Rightarrow$  Laplacian filter:

0	1	0
1	-4	1
0	1	0

## 5. Filters as templates:

Correlation filtering as Template Matching



## 6. Image Gradients

- Differentiation & Convolution:  $\frac{\partial f(x,y)}{\partial x} \approx \frac{f(x+1,y) - f(x,y)}{1}$

$$\Rightarrow [1] \boxed{-1}$$

Problem: it shifts the image

$\Rightarrow$  Prewitt, Sobel, Roberts Filters

$$\begin{bmatrix} -1 & 0 & 1 \\ -1 & 0 & 1 \\ -1 & 0 & 1 \end{bmatrix}; \begin{bmatrix} 1 & 0 & 1 \\ 0 & 0 & 0 \\ -1 & 0 & -1 \end{bmatrix}; \begin{bmatrix} 1 & 2 & 1 \\ 0 & 0 & 0 \\ -1 & -2 & -1 \end{bmatrix}$$

$$M_x = \begin{bmatrix} -1 & 0 & 1 \\ -2 & 0 & 2 \\ -1 & 0 & 1 \end{bmatrix} \quad M_y = \begin{bmatrix} 1 & 2 & 1 \\ 0 & 0 & 0 \\ -1 & -2 & -1 \end{bmatrix}$$

\* Gradient direction  
 $\theta = \tan^{-1} \left( \frac{\partial f}{\partial y} / \frac{\partial f}{\partial x} \right)$

Edge strength  
 $\| \nabla f \| = \sqrt{\left( \frac{\partial f}{\partial x} \right)^2 + \left( \frac{\partial f}{\partial y} \right)^2}$

- With noise, we need to smooth the image first

- Derivative Theorem of Convolution:

$$\frac{\partial}{\partial x} (h * f) = \left( \frac{\partial h}{\partial x} \right) * f$$

$$\Rightarrow g * (h * I) = (g * h) * I$$

$\Rightarrow$  Use "Derivative of Gaussian" filter  
 Only need to do 1 filter.. (?)

Derivative filter  
Gaussian filter

## 7) Edge Detection:

- Criteria for an "optimal" edge detector: robust to noise, good localization, single response

- Primary steps:
  - smoothing (suppress noise)
  - edge enhancement (filter for contrast)
  - edge localization
- Effect of  $\sigma$  on derivatives:
  - large  $\sigma \Rightarrow$  larger scale edges
  - small  $\sigma \Rightarrow$  finer features also
- Threshold:
  - high  $\Rightarrow$  less weak details (either relevant or not)
  - low  $\Rightarrow$  more

## ⊗ The Canny Edge Detector:

→ Gradient magnitude

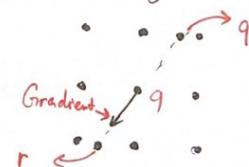
→ Thresholding

→ Non-Maximum Suppression:

Check if pixel is local maximum along gradient direction

(Select single max across width of the edge)

Keep  $q$  iff  $\text{Mag}(q) > \text{Mag}(p)$  and  $\text{Mag}(q) > \text{Mag}(r)$



→ Hysteresis Thresholding: 2 thresholds  $k_{\text{high}}$  &  $k_{\text{low}}$   
 $k_{\text{high}}$  starts edges,  $k_{\text{low}}$  continues edge chain (typical:  $\frac{k_{\text{high}}}{k_{\text{low}}} = 2$ )

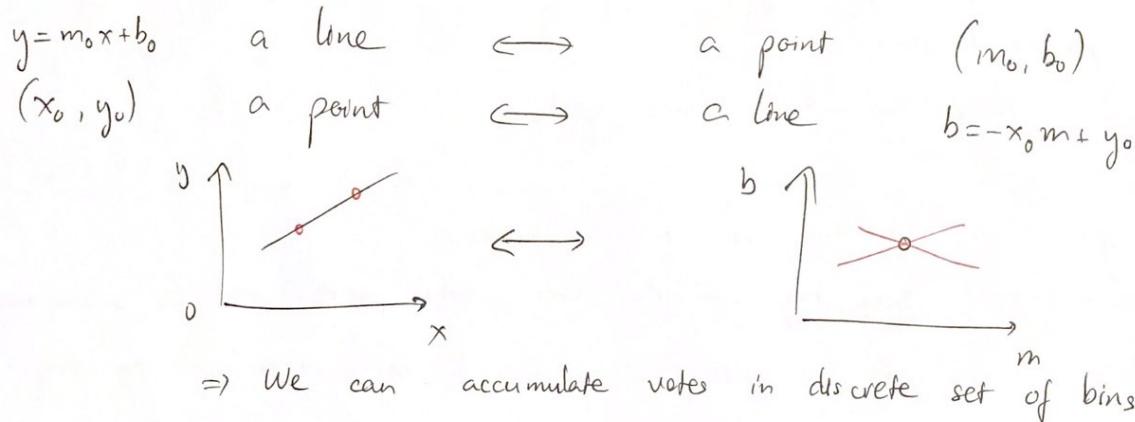
⊗ Edges  $\neq$  Object boundaries

Ex: complex texture, shadow ...

## 8, Structure Extraction:

- Hough transform:  $\begin{cases} \text{given points} \Rightarrow \text{what is the line} \\ \text{how many lines?} \\ \text{which point belongs to which lines} \end{cases}$

- Image  $\leftrightarrow$  Hough space

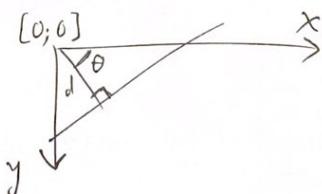


Ⓐ Problem: vertical lines

$\Rightarrow$  Polar Representation for lines:

$$x \cdot \cos \theta + y \cdot \sin \theta = d$$

$\Rightarrow$  Sinusoid segment in Hough space



Ⓑ Hough Transform Algorithm:

1, Initialize  $H(d, \theta) = 0$

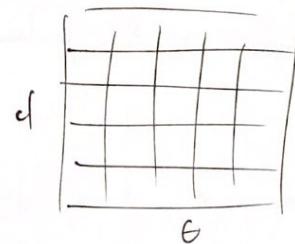
2, For each edge point  $(x, y)$

for  $\theta = 0 \rightarrow 180$  // stride

$$d = x \cos \theta + y \sin \theta$$

3, Final values of  $(\hat{d}, \hat{\theta})$  where  $H(d, \theta)$  is maximal

4, The detected line given by  $\hat{d} = x \cos \hat{\theta} + y \sin \hat{\theta}$



In practice,

range of

$$d: [-D; D]$$

$$\theta: \left[-\frac{\pi}{2}; \frac{\pi}{2}\right]$$

F

## Extensions:

- 1) Use Image gradient:  $\theta = \tan^{-1} \left( \frac{\partial f}{\partial y} / \frac{\partial f}{\partial x} \right)$  (for lines)
- 2) Give more votes for stronger edges (gradient magnitude)
- 3) Change the sampling of  $(d, \theta)$   $\Rightarrow$  more or less resolution  
less resolution  $\Rightarrow$  better to deal with noise, lose some details  
more  $\_\_\_ \Rightarrow \_\_\_$
- 4) Same procedure for circles, squares, .. (any shape can be parameterized)

circle  $\Rightarrow$  3D Hough space

- Tips:
- + Smoothing  $\Rightarrow$  vote for neighbors
  - + Keep tags on the votes, what points vote for which lines..
  - ↓ Cons: complexity search time  $\uparrow$  exponentially with No. model parameters
  - + Generalized Hough Transform: use define reference point

## + Hough transform for circle:

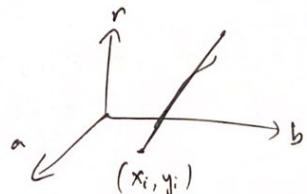
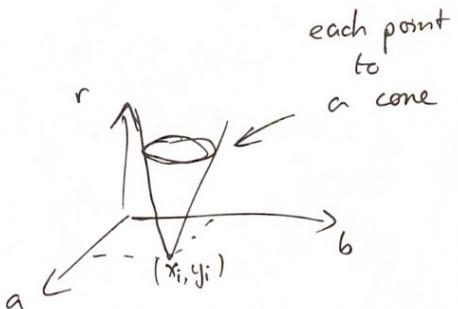
$$(x_i - a)^2 + (y_i - b)^2 = r^2$$

- Unknown  $r \propto$  gradient direction

- Unknown  $r$ , know grad\_dir

$$a = x - r \cdot c\theta$$

$$b = y + r \cdot s\theta$$



# Segmentation

L5,

## ④ Segmentation as Grouping Problem

### - Gestalt Theory:

The whole is greater than the sum of its parts.

Relationships among parts can yield new properties/features.

Gestalt Factors: proximity, similarity, parallelism, closure..

⇒ These factors make intuitive sense, but are very difficult to translate to algorithms.

## ⑤ Goals: identify groups of pixels that go together or find boundaries that separate objects

### ⑥ Segmentation as clustering:

#### - K-Means Clustering: each pixels as a data point in the feature space NP-hard even with $k=2$

#### K-Means++: as an initialization for centroids $c_i$

1) Random first center

2) Pick next center with the probability proportional to the distance to previous centers  $\frac{1}{\|p - c_i\|^2}$  ⇒ kind of prioritize further centers

3) Repeat till  $k$  centers

- A variety of feature space choice ⇒ we can group pixels in different ways. ↗ intensity, color, texture position

Problems  
with  
K-Means ?

## ② Probabilistic Clustering:

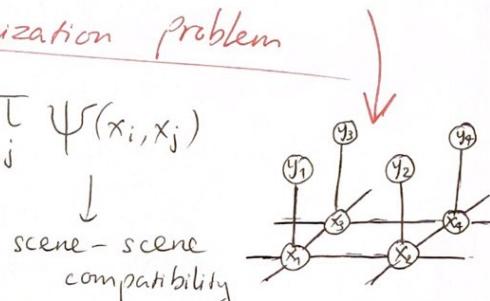
- Expectation - Maximization : compare with K-Means
- ② Model free Clustering : no shape, no number of cluster.
- Mean-Shift Algorithm : only specify window size
  - ↳ Mode: local maximum of the density of given distribution
  - ↳ Cluster: all data points in the attraction basin of a mode
  - ↳ Attraction basin: region for which all trajectories lead to same mode
- Tricks: 1) Assign all points within radius  $r$  around final mode  
2) within radius  $r/c$  of search path to the mode
- Cons: depends on window size  
computationally expensive, does not scale well with dimension of feature space

Markov  
Random  
Fields

## ③ Segmentation as Energy Minimization problem

$$\text{Maximize } p(x, y) = \prod_i \phi(x_i, y_i) \cdot \prod_{i,j} \psi(x_i, x_j)$$

↓  
image-scene compatibility



↳ Minimize energy formulation:

$$E(x, y) = \sum_i \phi(x_i, y_i) + \sum_{i,j} \psi(x_i, x_j)$$

single-node unary potentials      pairwise potentials

$$\begin{cases} \phi(x_i, y_i) = -\log \phi(x_i, y_i) \\ \psi(x_i, x_j) = -\log \psi(x_i, x_j) \end{cases}$$

## ④ Many inference algorithms:

- Gibbs sampling, simulated annealing
- Iterated conditional modes (ICM)
- Variational methods
- Belief propagation
- Graph cuts

L6, ④ Graph cut: ones of many inference algorithm  
 only for binary energy functions ??  
 but if we can fit it to our problem, can be implemented very fast

④ Unary terms: just normal classifier

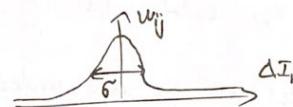
If we have given color model for initial classes:  $I^s, I^t$

$$\Rightarrow \begin{cases} \phi_i(s) \propto \exp(-\|I_i - I^s\|^2 / 2\sigma^2) \\ \phi_i(t) \propto \exp(-\|I_i - I^t\|^2 / 2\sigma^2) \end{cases}$$

④ Pairwise terms:

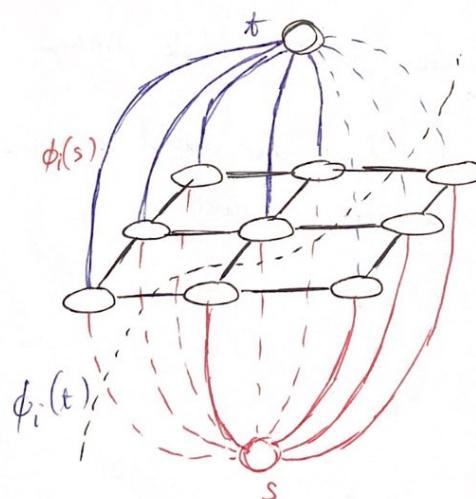
Ex:  $\psi(x_i, x_j) = w_{ij} \cdot \delta(x_i \neq x_j) \rightarrow$  only active when we cut that weight

$$w_{ij} = \exp \left\{ \frac{-\Delta I_{ij}^2}{2\sigma^2} \right\}$$



$\Rightarrow$  if pixels (in neighboring) is similar  $\Rightarrow$  great weight  
 not similar  $\Rightarrow$  small weight

Cause only when we cut that link between 2 pixels of different classes, that we will loss that weight



⊗ Example:

→ Color potentials: - as unary potentials.

e.g. modeled with MoGs

$$\phi(x_i, y_i; \theta_\phi) = -\log \sum_k \Theta_\phi(x_i; k) \cdot p(k|x_i) \mathcal{N}(y_i; \mu_k, \Sigma_k)$$

→ Edge potentials: - as pairwise potentials.

$$\varphi(x_i, x_j, g_{ij}(y); \theta_\psi) = \Theta_\psi \cdot g_{ij}(y) \cdot \delta(x_i \neq x_j)$$

$$g_{ij}(y) = e^{-\beta \|y_i - y_j\|^2} ; \quad \beta = \frac{1}{2} \left[ \text{avg}(\|y_i - y_j\|^2) \right]^{-1}$$

⇒ If there is no label change:  $x_i = x_j \Rightarrow \text{cost } \varphi = 0$

If there is label change:  $x_i \neq x_j$

If pixel values are indeed different  $y_i \neq y_j \Leftrightarrow y_i - y_j \uparrow$   
 $\Rightarrow g_{ij}(y) \rightarrow 0$  very small

If not  $y_i \approx y_j \Leftrightarrow g_{ij}$  will be great

-  $\beta$  is pre-calculated on the whole image (for each img or all, depends on each situation ..)

-  $\theta_\phi, \theta_\psi$  decide how our model is sensitive to color, or edge ..  
 need to be learned as well.

## ⑧ s - t - Min cut Problem:

- An strut divides the nodes between source & sink
  - The cost of the cut: sum of costs of all edges going from S to T
- $\Rightarrow$  st - mincut is st-cut with minimum cost
- edges that are cutted

## + Max flow Algorithms:

- Constraints:
  - Edges: Flow  $<$  Capacity
  - Nodes: Flow in = Flow out
- Algorithms:
  - Find path from source  $\rightarrow$  sink with positive capacity
  - Push max flow through
  - Adjust capacity
  - Repeat till no new path possible
- s-t Graph Cuts can only be applied?  
Can only minimize binary energies, that are sub-modular  
(discrete equivalent to convexity)
- Augmenting Path                          vs                          Push - Relabel  
all edge weights must  $> 0$                            $\forall w_{ij}$

### + $\alpha$ -Expansion Moves:

- Each move add new label & decrease energy
- Choose move that minimize the energy the most

### + $\alpha$ - $\beta$ swap

## ④ Application:

+ Interactive Image Segmentation

+ Iterated Graph Cut

## ⑤ Summary Graph-Cut:

- Applicable for a wide range of problems
- Very efficient algorithm for vision problems
- Become standard

- 
- Can only solve limited class of models
  - Only optimal for sub-modular problem  
For multi-label case, only approximate

# Object Detection

47,

Classification + Localization = Object Detection < Instance Segmentation

- + Object Recognition Challenges?
- + Appearance-based recognition: 3D objects can be represented by a set of images appearances, for recognition, it's sufficient to just compare the global representation
- + Identification vs Categorization

same object seen find particular object which we have seen before	new object, in only same class
---	--------------------------------

- + Sliding-window based Object Detection brute force
  - (Slide window around (1 pixel to the right, down...))
  - (Over different scale (scale the image, not the window))

- + Detection via Classification: Main Idea
  - 1, Obtain training data (Ex: Car & non-car data)
  - 2, Define features
  - 3, Define classifier

- + Feature Extraction:

- Global Appearance such as pixel-based appearance; (color or grayscale)
  - not robust
  - ↳ sensitive to small shifts, illumination, intra-class variation ..
- ⇒ Gradient-based Representations invariance to small shifts & rotations
  - localized histograms offer more spatial info than single global hist.

L<sub>8</sub>

## HOG (Histogram of Oriented Gradients)

If we want to stick with simple linear classifier, then our feature representation has to be very good. <sup>SVM</sup>

Inputs: prepared images with same shape  
divide to grid cell ( $8 \times 8$  pixels)

Training phase  $\rightarrow$  HOG Descriptor Processing Chain:

1, Gamma Compression:  $x \mapsto \sqrt{x}$  small performance improvement

2, Compute gradients (both magnitude & direction)

3, Weighted vote in spatial & orientation cells

From grad-direction  $\Rightarrow$  index of bin

Note: NOT '+1', but '+grad-mag'

4, Contrast normalize over overlapping spatial cells

Each block has 4 cells ( $2 \times 2$ )

Each block overlapped 50% with the next block.

$\Rightarrow$  L2 normalization within 4cell, Clipping (0,2), L2 norm again

Again, just like gamma suppression, the purpose is to reduce the effect of very strong gradients, that could overshadow others.

5, Convert 3D array of histogram into 1D feature vector  $\Rightarrow$  SVM

⊗ The SVM will deal with only 1 type of dimension vector thus, all training inputs, and then the window-shape in testing phase must be the same

## ⊕ HoG in Testing Phase:

- Is like convolving HoG feature map of test image with the learned template  $w$
- Non-Maximum Suppression: to get single detection
  - ⊕ Can break ties with noise
  - ⊕ Also clip the score ( $> 0.8 \dots$ )
  - ⊕ Is clustering problem.. could use meanshift..
- With Sliding-window approach:

↳ Move window by 1 cell, not 1 pixel

For different scale detection, down/upscale the image, while fixed everything else (cell size..)

Image Pyramid

## Boosting:

N data  
M iterations

Final decisions:

$$H(x) = \text{sign} \left( \sum_{m=1}^M \alpha_m h_m(x) \right)$$

$$I(A) = \begin{cases} 1, & \text{if } A \text{ true} \\ 0, & \text{else} \end{cases}$$

Example: Viola Jones face detector

As example for Adaboost

1) Initialize:  $w_n^{(1)} = \frac{1}{N}$  for  $n = 1, \dots, N$

2) For  $m = 1, \dots, M$  iterations

- Train weak classifier  $h_m(x)$  to minimize  $J_m = \sum_{n=1}^N w_n^{(m)} I(h_m(x_n) \neq t_n)$

- Estimate weighted error of  $h_m$  on  $X$ :  $E_m = \frac{J_m}{\sum_{n=1}^N w_n^{(m)}}$

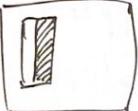
- Calculate a weighting coefficient for  $h_m(x)$ :  $\alpha_m = \ln \left( \frac{1-E_m}{E_m} \right)$

- Update weighting coefficients:  $w_n^{(m+1)} = w_n^{(m)} \exp \left\{ \alpha_m I(h_m(x_n) \neq t_n) \right\}$

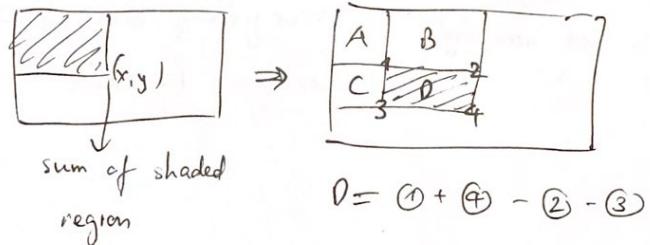
## ⊗ Limitations of Sliding - Window approach:

- Not all objects are "box" shape
  - Window doesn't consider context: Ex: pedestrian is on the ground only, not in the sky
- ⇒ That's why we call it brute-force

## ⊗ Niela-Jones Face Detector:

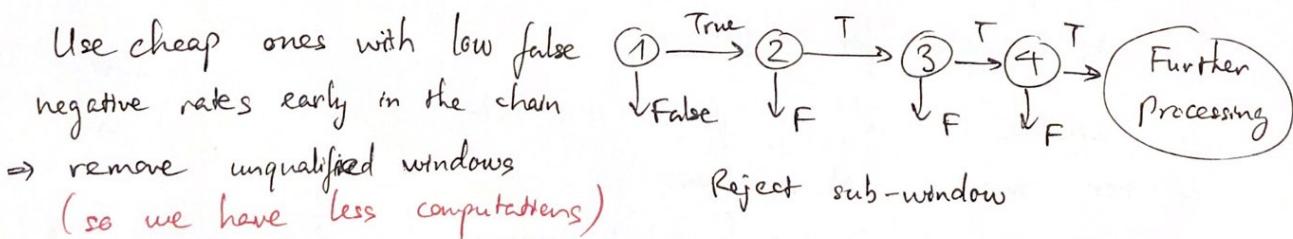
- "Rectangular" filters:  
 just rectangular regions with specific position in a patch

- Efficiently computable with Integral image



- Simple classifier for each filters: output > θ ?  
just like decision tree

- Cascading Classifiers for Detection:



- Train with 5k positives, 350M negatives

Test with 38 layer cascade

(1st layer 2 features filters → last layer 60 61 filters...)

# Local feature

- ↳ - Motivation: ↗ Global representations have major limitations  
 (pedestrian: head, arm, leg..)
- ↗ Local feature: more robust to Occlusions  
Articulation  
Intra-class variations
- (half the face is hidden,  
 find the eye  $\Rightarrow$  face;  
 find leg angle; different koala .. but same nose)
- Application: ↗ Image matching  
 ↗ Image stitching
- Procedure:
- 1, Find a set of distinctive key points
  - 2, Define region around each key points
  - 3, Extract & normalize region content
  - 4, Compute a local descriptor from normalized region
  - 5, Match local descriptors / Final corresponding pairs

- ⊗ Detector Requirements:
- Invariant to Geometric Invariance (translation, rotation, scale, ...)
  - Robust to lighting, noise, blur..
  - Local: not global feature  $\Rightarrow$  robust to occlusion clutter
  - Distinctive.
  - Efficient: real time
  - Quantitative

⊗ Many Existing Detectors available:

Hessian & Harris; Laplacian, DoG; Harris-/Hessian-Laplace;  
Harris-/Hessian-Affine; EBR & IBR.

- Corners are repeatable & distinctive (significant changes in all directions)

†) Harris Detector:

2<sup>nd</sup> Moment Matrix

$$M = \sum_{x,y} w(x,y) \begin{pmatrix} I_x^2 & I_x I_y \\ I_x I_y & I_y^2 \end{pmatrix}$$

window function

$$M = R^{-1} \cdot \begin{bmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{bmatrix} \cdot R \Rightarrow \begin{cases} \text{both } \lambda_1 \text{ & } \lambda_2 \text{ is great} \Rightarrow \text{corner} \\ \begin{cases} \lambda_1 \text{ is great} \\ \lambda_2 \approx 0 \end{cases} \Rightarrow \text{edge} \end{cases}$$

Eigenvalue decomposition

- Instead of find  $\lambda_i$ , we can do a cheap approximation:

$$R = \det(M) - \alpha \cdot \text{trace}(M)^2 = \lambda_1 \lambda_2 - \alpha (\lambda_1 + \lambda_2)^2$$

$$\alpha = \text{const} (0,09 \rightarrow 0,06)$$

Just need M

$$\begin{cases} R > 0 \Rightarrow \text{corner} \\ R < 0 \Rightarrow \text{edge} \end{cases}$$

- Window function  $w(x,y)$ : uniform - square window: not rotation invariant

⇒ Gaussian window (rotation invariant)

$$\Rightarrow M = g(\sigma) * \begin{pmatrix} I_x^2 & I_x I_y \\ I_x I_y & I_y^2 \end{pmatrix}$$

the window function is to consider the window around a pixel

Because the Harris Detector doesn't say that there is a corner at that pixel  
but that there is a corner in the window around that pixel

## $\Rightarrow$ Harris Detector algorithm summary:

Use derivative  
of Gaussian

1, Image derivatives:  $I_x, I_y$

1st derivative

2, Square img-der:  $I_x^2, I_x I_y, I_y^2$

$\sigma_1$

3, Gaussian filter  $g(\sigma)$ :  $g(I_x^2), g(I_x I_y), g(I_y^2) \quad \sigma_2 = 2\sigma_1$

4, Cornerness func:  $R = \det[M(\sigma_x, \sigma_y)] - \alpha \cdot [\text{trace}(M(\sigma_x, \sigma_y))]^2$

$\alpha = \text{const}$

(0.04  $\rightarrow$  0.06)

$$= g(I_x)^2 g(I_y)^2 - g(I_x I_y)^2 - \alpha g(I_x)^2$$

$$= g(I_x^2) \cdot g(I_y^2) - [g(I_x I_y)]^2 - \alpha [g(I_x^2) + g(I_y^2)]^2$$

5, Non-Maximum Suppression & Threshold

⊖  $\Rightarrow$  Problem: Harris Detector is rotation invariant  
but is not scale invariant

⊖ Note: with window func., Harris Detector takes into account larger image neighborhood, unlike Hessian Detector

## + Hessian Detector: 2nd derivative

Intuition: Search for strong derivatives in 2 orthogonal direction  
will focus on mainly corners & strong textured areas

Also use 2nd Gaussian  
derivative  
with  $\sigma_1$

$$- \text{Hessian}(I) = \begin{bmatrix} I_{xx} & I_{xy} \\ I_{yx} & I_{yy} \end{bmatrix} \quad \text{Harris: just corners} \Rightarrow \det(\text{Hessian}(I)) = I_{xx} I_{yy} - I_{xy}^2$$

- Apply NMS, threshold..

-  $\det(H(I)) > \text{threshold} \Rightarrow$  corners

$$\sigma^4 \det(H) = \sigma^4 [I_{xx}(\sigma) \cdot I_{yy}(\sigma) - I_{xy}^2(\sigma)] > +$$

⊖ In practice, actually combine both of these 2

If want precise locations  $\Rightarrow$  Harris  
additional locations  $\Rightarrow$  Hessian

$$\frac{\text{Tr}(H)^2}{\det(H)} = \frac{(r+1)^2}{r} \quad r = \frac{\sigma_1}{\sigma_2}$$

$r > 10$  eliminate (edges)

$r = 1$  good

\*  $\sigma^4$  because the peaks of der-Gaussian is reduced by  $\sigma^2$

$$G = \frac{1}{2\pi\sigma^2} \dots \text{and we have } I_{xx} \cdot I_{yy} \Rightarrow \sigma^4 \text{ so that } t \text{ is independent with } \sigma$$

## 40) + Scale Invariant Region Selection:

- Exhaustive search over different scale: is stupid
- Scale signature function
- Laplacian of Gaussian - "blob" detector  
Define characteristic scale as the scale produces peak of LoG response  

$$L_{xx}(\sigma) + L_{yy}(\sigma) \leq \frac{\sigma^2}{\sigma^2} \Rightarrow \text{list of } (x, y, \sigma)$$
- Efficient Approximation LoG by DoG  

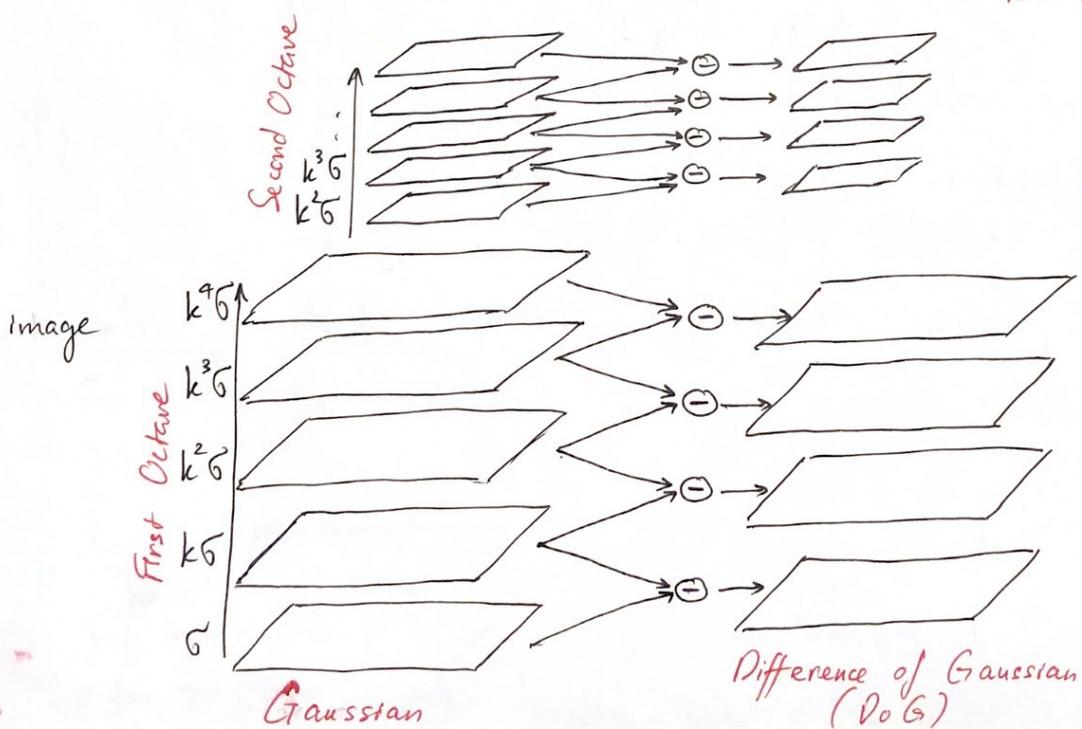
$$L = \sigma^2 (G_{xx}(x, y, \sigma) + G_{yy}(x, y, \sigma))$$

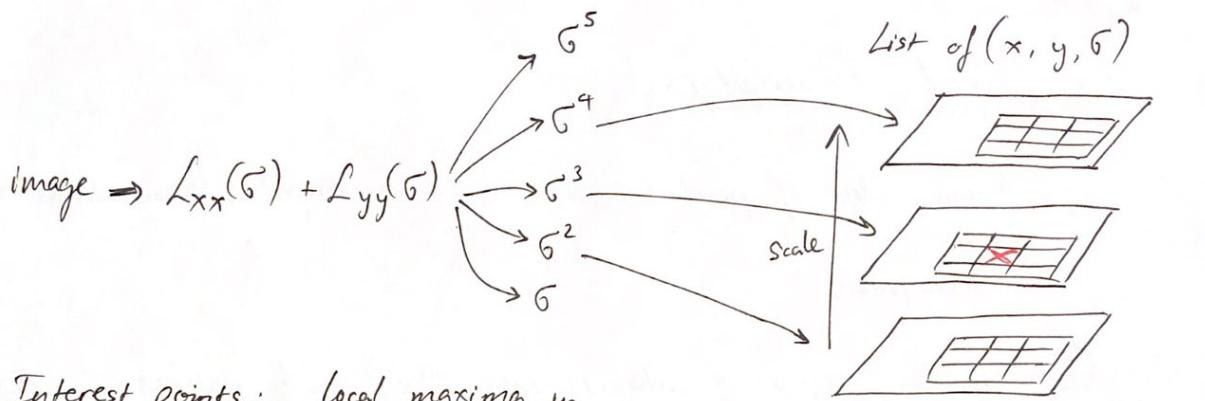
$$\text{DoG} = G(x, y, k\sigma) - G(x, y, \sigma)$$

## Scale Pyramid:

$$\sigma = 1,6$$

$$k = \sqrt{2}$$

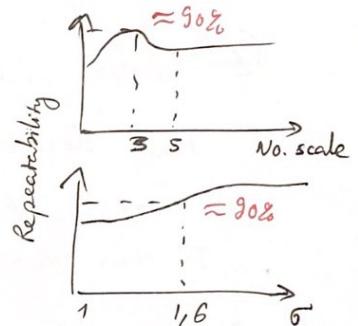




Interest points: local maxima in  
scale space of  $\text{Lc } G$  (min or max among  $27 - 1 = 26$  neighbors)

### Tips

- Should choose 3-5 scales sampled per octave
- Init  $\sigma = 1, 6$



		k scale						.
		1st	$1/\sqrt{2}$	1	$\sqrt{2}$	2	$2\sqrt{2}$	.
Octave	2nd	$\sqrt{2}$	2	$2\sqrt{2}$	4	$4\sqrt{2}$	.	
	3rd	$2\sqrt{2}$	4	$4\sqrt{2}$	8	$8\sqrt{2}$	.	

Note: The next octave is 4 times smaller than the previous  
 (2 times in each axis)

⇒ Thus applying  $\sigma_i$  at one octave  
 will be just equal to apply  $2\sigma_i$  at the previous octave?

## +)Local Descriptors:

- Again, list of pixel intensities . . . is stupid, translational variant

$\Rightarrow$  Histograms

① Divide region of interest into  $4 \times 4$  sub-patches /  $8$  bins

$$\Rightarrow \underline{4 \times 4 \times 8 = 128\text{-dim}} \quad \text{SIFT descriptor}$$

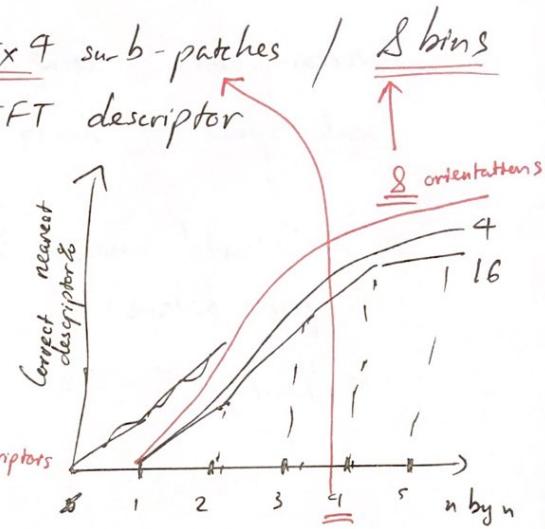
② To be rotational invariant:

Find dominant grad

$\Rightarrow$  rotate to that direction

(If there are 2 dominant grad  $\rightarrow$  create 2 descriptors)

③  $\otimes$  { Normalize  
clip  $(0,2)$   
Re-normalize



$\otimes$  Minimum size of the region should be  $\xrightarrow{32 \times 32}$  at worst  $16 \times 16$

+)Key point matching: find nearest neighbor

If the ratio of distance between best & 2nd best is very close ( $0,8$ )

$\Rightarrow$  don't pick it

SIFT: Scale Invariant Feature Transform

+ ) Find consistent configurations:

- Warping: given  $\begin{cases} \text{source} \\ \text{transformation} \end{cases}$   $\Rightarrow$  find output

2 types  
of  
problem

- Alignment: given  $\begin{cases} \text{source} \\ \text{output} \end{cases}$   $\Rightarrow$  find transformation

- 2D Affine Transformation:

$$\begin{matrix} \text{Rotation} & \left[ \begin{array}{c} x' \\ y' \\ w \end{array} \right] = \left[ \begin{array}{ccc} a & b & c \\ d & e & f \\ 0 & 0 & 1 \end{array} \right] \cdot \left[ \begin{array}{c} x \\ y \\ w \end{array} \right] & \text{Translation} \\ \left[ \begin{array}{ccc} \cos\theta & -\sin\theta & 0 \\ \sin\theta & \cos\theta & 0 \\ 0 & 0 & 1 \end{array} \right] & & \left[ \begin{array}{ccc} 1 & 0 & t_x \\ 0 & 1 & t_y \\ 0 & 0 & 1 \end{array} \right] \end{matrix}$$

$$\left[ \begin{array}{ccc} \cos\theta & -\sin\theta & 0 \\ \sin\theta & \cos\theta & 0 \\ 0 & 0 & 1 \end{array} \right]$$

$$\begin{matrix} \text{Scaling} \\ \left[ \begin{array}{ccc} s_x & 0 & 0 \\ 0 & s_y & 0 \\ 0 & 0 & 1 \end{array} \right] \end{matrix}$$

Shearing

$$\left[ \begin{array}{ccc} 1 & sh_x & 0 \\ sh_y & 1 & 0 \\ 0 & 0 & 1 \end{array} \right]$$

Trans + Rot = Euclidean trans.

Trans. + Rot. + Scaling = Similarity trans.

Trans + Rot + Scaling + Shear = Affine trans

+ Projective warp = Projective trans.

- Fitting a Affine Transformation: 6 variables  $\Rightarrow$  need 3 points

$$\left[ \begin{array}{c} x'_i \\ y'_i \end{array} \right] = \left[ \begin{array}{cc} m_1 & m_2 \\ m_3 & m_4 \end{array} \right] \left[ \begin{array}{c} x_i \\ y_i \end{array} \right] + \left[ \begin{array}{c} t_1 \\ t_2 \end{array} \right] \Leftrightarrow \left[ \begin{array}{cccccc} x_i & y_i & 0 & 0 & 1 & 0 \\ 0 & 0 & x_i & y_i & 0 & 1 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \end{array} \right] \left[ \begin{array}{c} m_1 \\ m_2 \\ m_3 \\ m_4 \\ t_1 \\ t_2 \end{array} \right] = \left[ \begin{array}{c} x'_i \\ y'_i \\ \vdots \end{array} \right]$$

$\Rightarrow$  Solve with least Squares Estimation

$$x = (\text{pseudo inverse of } A) \cdot B$$

Z

→ Homography: a projective transformation

$$\begin{bmatrix} x' \\ y' \\ z' \end{bmatrix} = \begin{bmatrix} h_{11} & h_{12} & h_{13} \\ - & - & - \\ h_{31} & h_{32} & 1 \end{bmatrix} \cdot \begin{bmatrix} x \\ y \\ 1 \end{bmatrix}; \quad \begin{bmatrix} x'' \\ y'' \\ 1 \end{bmatrix} = \frac{1}{z'} \begin{bmatrix} x' \\ y' \\ z' \end{bmatrix}$$

$$\Leftrightarrow x' = H \cdot x \quad ; \quad x'' = \frac{1}{z'} \cdot x'$$

$$\Leftrightarrow \begin{bmatrix} x_B, & y_B, & 1 & 0 & 0 & 0 & -x_A, x_B, & -x_A, y_B, & -x_A, \\ 0 & 0 & 0 & x_B, y_B, 1 & -y_A, x_B, & -y_A, y_B, & -y_A, \\ . & . & . & . & . & . & . \end{bmatrix} \begin{bmatrix} h_{11} \\ \vdots \\ h_{32} \\ 1 \end{bmatrix} = [0]$$

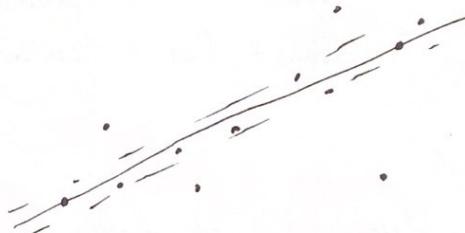
$$\Leftrightarrow A \cdot h = 0$$

Do some Singular Value Decomposition shit  $\Rightarrow h$

- When we have perpendicular camera axes,  $h_{33} \approx 0$
- If not, can set as  $A'h^* = b$ ,  $h^* = [h_{11} \dots h_{32}]^T$   
then use pseudo inverse to find  $H$

→ Dealing with Outliers: each transformation as a point (80)

- RANSAC:



- Alternatively: Generalized Hough transform

# Deep Learning

L11,12,13)

- Max pooling: robustness to translation
- If use Batch Norm  $\Rightarrow$  no need of biases  
Turn it off to have less params
- If next layer after Batch norm is linear, conv (might be a RELU in the middle), then don't need scale for  
Cause instead, we can multiply the weights later
- Instead of Max pooling, use stride in conv layer,  
it reduce running time with less calculation
- Global Average Pooling Layer to deal with input of different size, and also reduces params significantly  $(h \times w \times d) \rightarrow (1 \times 1 \times d)$   
average of each  $(h \times w)$
- Augmentation: better validation results
- VGGnet vs AlexNet: more conv layers with smaller filter size  
is better than a conv layer with larger size:  
2 @  $3 \times 3$  is better than 1 @  $5 \times 5$   
less params, while same receptive field
- $1 \times 1$  conv layer = "bottleneck" layer, for dimensionality reduction

- GoogleNet: Inception module  
No FC layer
- Residual Net: Skip connection  $\Rightarrow$  better grad propagation
- Why don't scale up the gradients when dealing with grad. vanishing?  
Cause the signal-noise ratio, we will also scale up the noise

- Top 1 vs Top 5 accuracy rate?

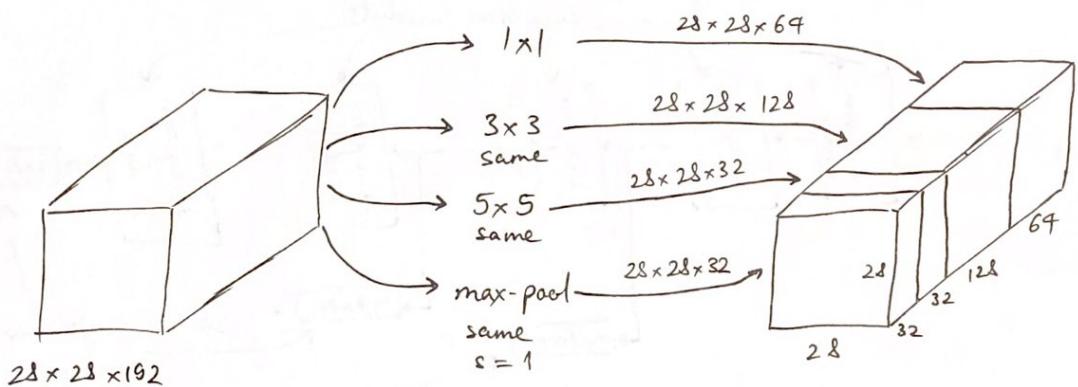
Example: output:  $\begin{cases} \text{Dog} & 70\% \\ \text{Cat} & 20\% \\ \text{Pig} & 5\% \\ \vdots & \end{cases}$   $\Rightarrow$   $\begin{cases} \text{Top 1: False} \\ \text{Top 5: Correct} \end{cases}$

True Label is "Cat".

- Transfer learning:
  - [ small data set: retrain last classifier layers ]
  - [ medium-sized dataset: retrain more layers (fine-tune) ]

- residual Net from exercise

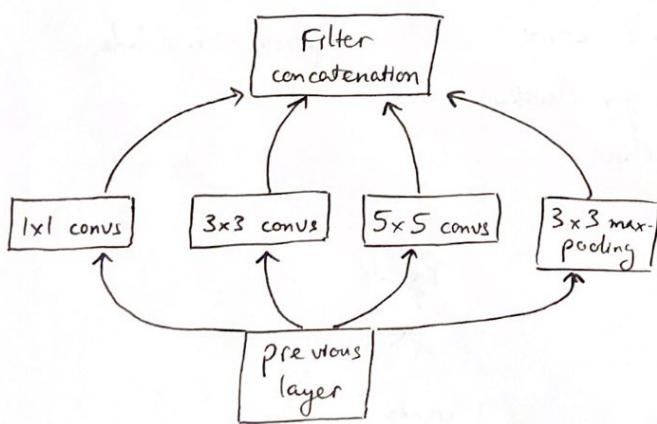
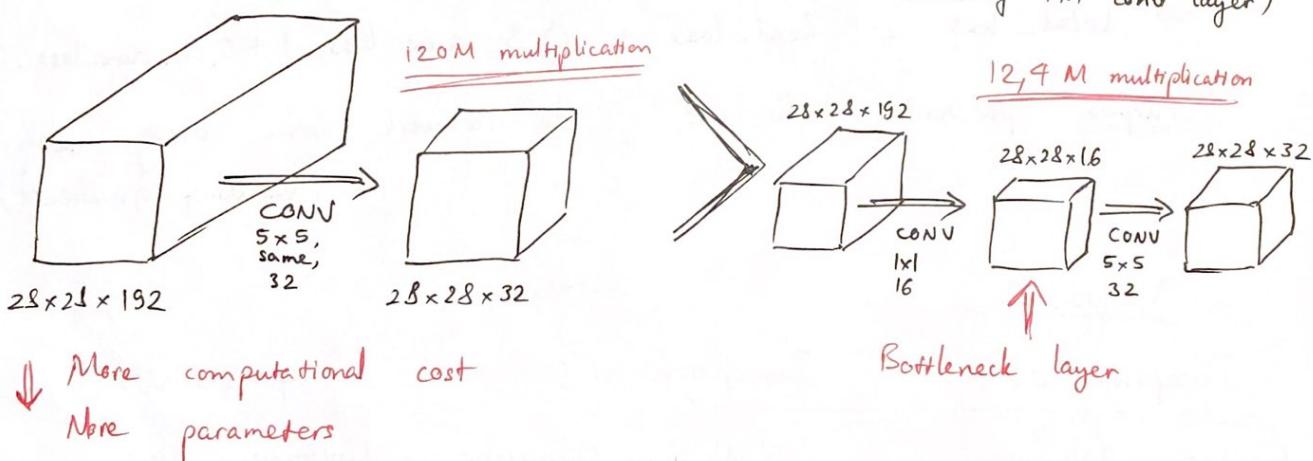
# + ) Google Net : Inception module



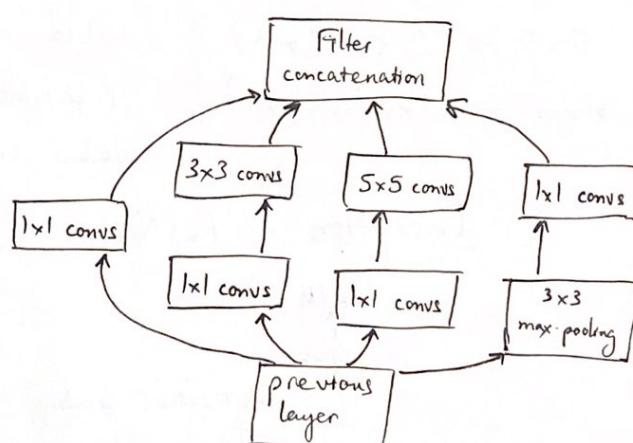
- Naive version

vs

with dimension reductions  
(use of  $1 \times 1$  conv layer)

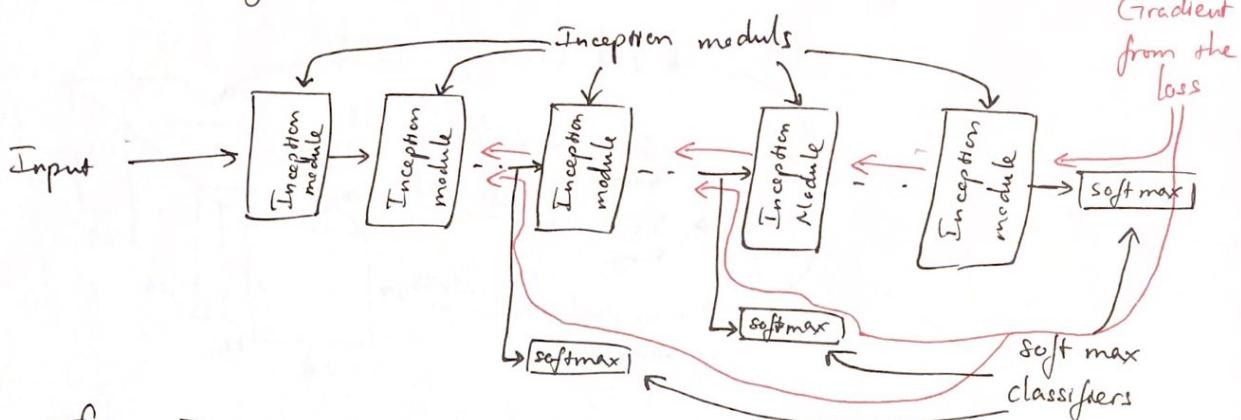


Inception module  
(Naive version)



Inception module  
with dimension reductions

## - Auxiliary classifiers:

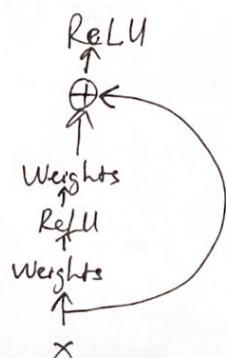
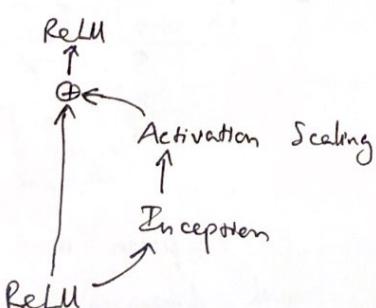


- The Inception network is simply a stack of 9 inception modules
- With 2 extra branches, each branch is applied a softmax to calculate auxiliary loss  
 $\Rightarrow \text{Total\_loss} = \text{Real\_loss} + 0.3 \cdot \text{aux\_loss\_1} + 0.3 \cdot \text{aux\_loss\_2}$
- ~~② Purpose:~~ prevent middle part of the network from "dying" out (vanishing gradients)

## - Versions:

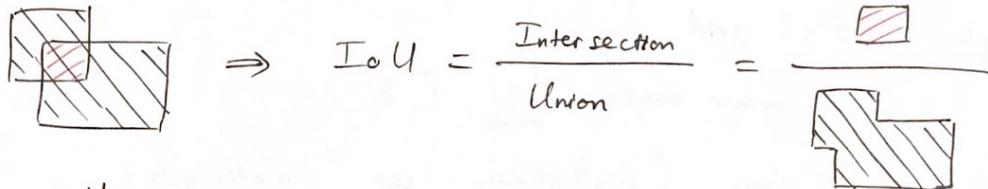
Inception .v2	Inception .v3	Inception .v4
Factorize filter size $5 \times 5 \Rightarrow 2 \times (3 \times 3)$	RMS Prop Optimizer Add $7 \times 7$ conv Batch norm for Auxiliary Label smoothing	Uniformize the Inception modules
$n \times n \Rightarrow (1 \times n) \times (n \times 1)$		

## - Inception - ResNet:



## +), IoU (Intersection over Union)

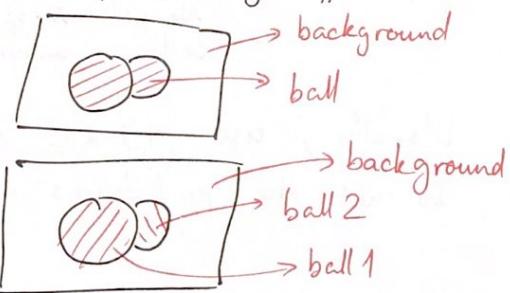
2 boxes



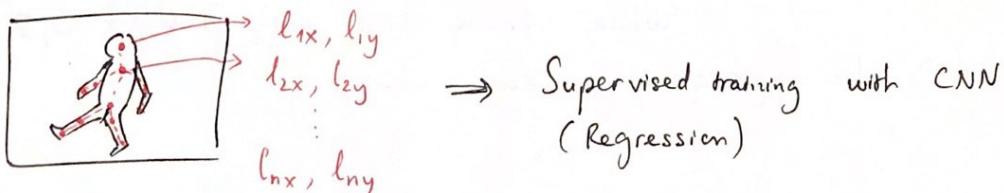
- Use to assess accuracy: correct if  $\text{IoU} \geq \text{threshold}$   
(commonly threshold = 0,5, or greater)
- Also used in Non max Suppression in YOLO

## +), Definition:

- Classification
- Classification & Localization } 1 object
- Detection : multiple objects, of possibly different classes
- Semantic Segmentation
- Instance Segmentation



## +), Landmark detection:



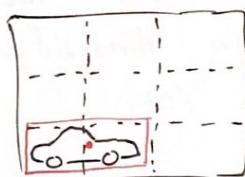
Label key landmark

⇒ Usage: { estimate pose  
          { recognize emotion

## ④ Anchor box based detectors:

+ YOLO: You only look once

Example:  $\begin{cases} 3 \times 3 \text{ grid} \\ 2 \text{ anchor boxes} \\ 3 \text{ class (pedestrian, car, motorcycle)} \end{cases}$



For each grid cell, the output will be:

$$y = \left[ p_c \underbrace{b_x \ b_y \ b_w \ b_h}_{\substack{\text{prob. that} \\ \text{there is object}}} \underbrace{c_1 \ c_2 \ c_3}_{\substack{\text{bounding box} \\ \text{coordinate}}} : \underbrace{p_c \ b_x \ b_y \ b_w \ b_h \ c_1 \ c_2 \ c_3}_{\substack{\text{prob. of} \\ \text{each class}}} \right]^T$$

=> Total output:

$$\underbrace{3 \times 3}_{\substack{\text{No. grid} \\ \text{cell}}} \times \underbrace{2 \times 8}_{\substack{\text{No. of} \\ \text{anchor box}}} \rightarrow 5 + \text{No. of classes}$$

- Usually, use  $19 \times 19$  grid cell  
so that the probabilities of 2 same class objects in the same grid cell is small

- Filter result with Non-max suppression:

↳ Threshold: remove boxes with low  $p_c$  confident. ( $p_c \leq 0,6$ )

↳ Non-max suppression: keep the box with highest confidence  
while remove boxes with IoU ( $> 0,5$ )

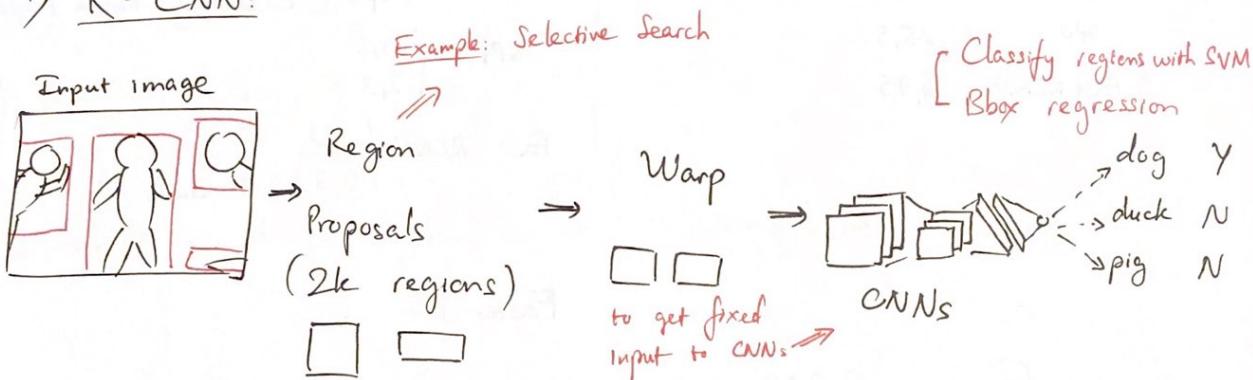
(Do this separately for each class)

+ SSD: Single shot detector

## \* CNNs for Object Detection:

### \* Regional-based proposal based Detectors:

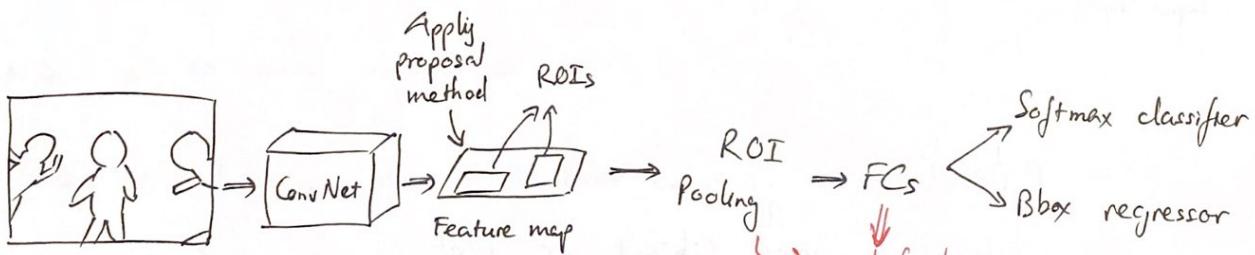
#### + R-CNN:



It works .. but there are so much to improve.

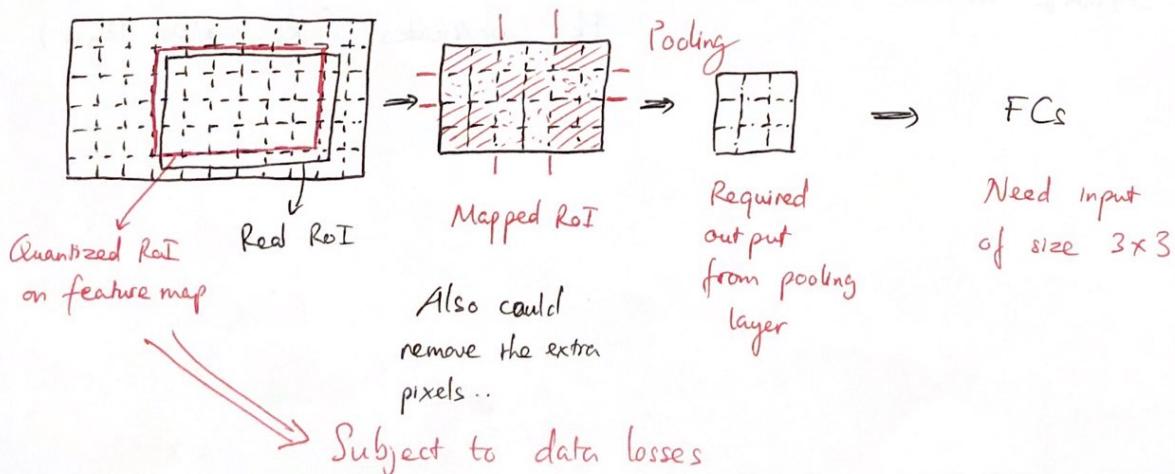
Training 3 days & testing 47s per image are slow

#### + Fast R-CNN:



\* Still use external proposal method (e.g. Selective Search) to get ROIs

#### - ROI Pooling:



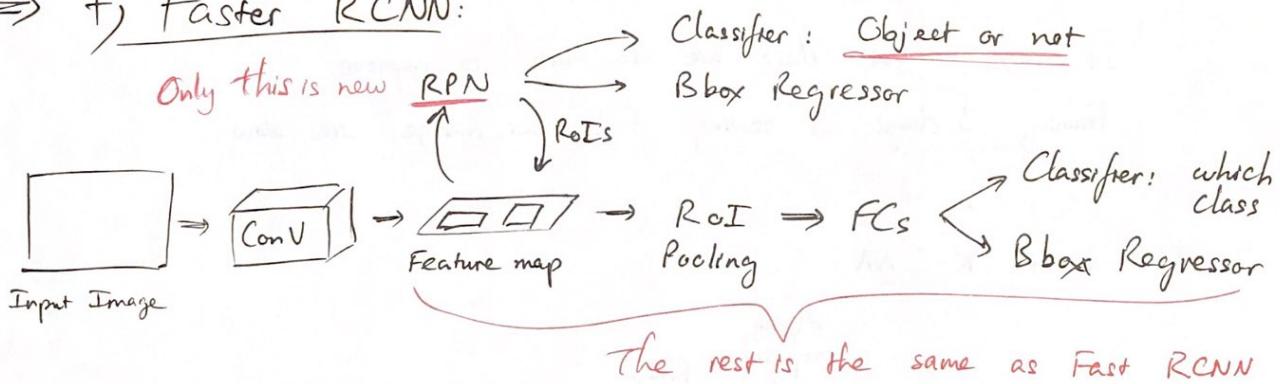
+ ) Compared R-CNN vs SPP vs Fast R-CNN:

	Training Time (hours)	Test time (secs)
R-CNN	84	49
SPP	25,5	47 (Excluding Region Proposal)
Fast R-CNN	8,75	4,3
		2,3 (_____)
		Fast-R-CNN { 2,3 0,32 (Exclude ..) }

→ Problem of Fast R-CNN: runtime dominated by Region Proposals!

Faster R-CNN: 0,2

⇒ + ) Faster R-CNN:



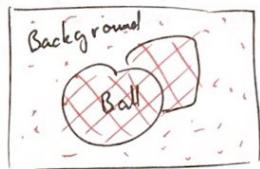
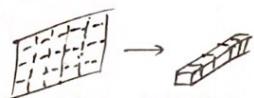
- RPN (Region Proposal Network): is another Conv-net

- [ Classify error: Object or Not
- Bbox error

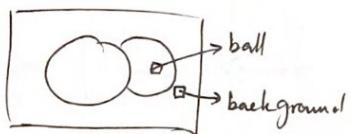
⇒ With  $k$  anchor boxes  $\Rightarrow \begin{cases} 2k \text{ scores } (\text{Classifier layer}) \\ 4k \text{ coordinates } (\text{Regressor layer}) \end{cases}$

→ Semantic Segmentation: label each pixel with class label  
but don't care how many objects there are

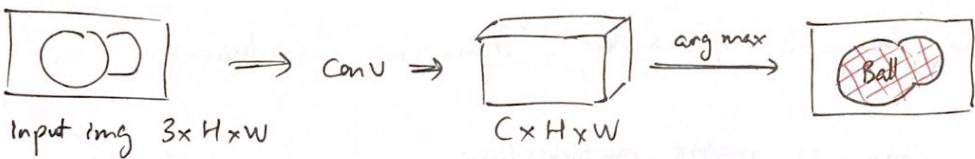
- Ideas: - Can't use normal FCs on the whole image, cause we will loose pixel location



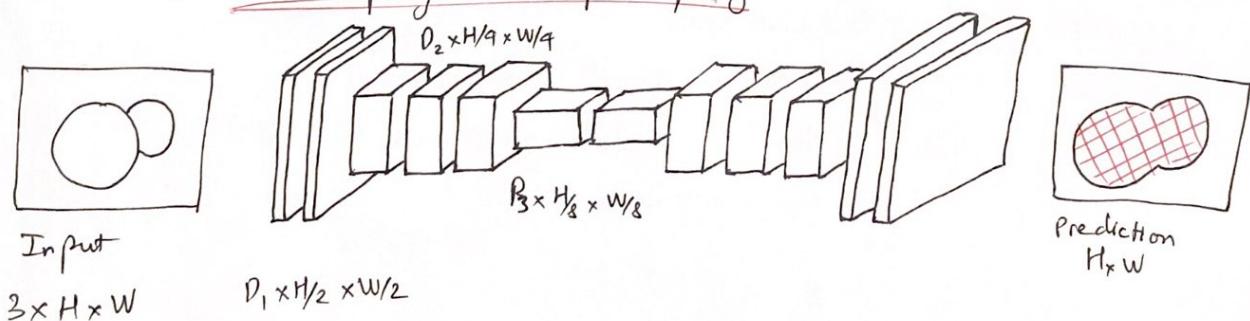
- Sliding window: run small patch through CNN  
very inefficient  
no reuse of features between shared patches



→ FCN (Fully Convolutional Net): can process arbitrary image size

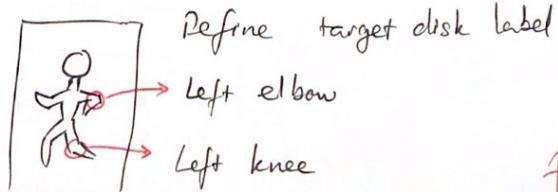


- Problem: conv at original image resolution will be expensive  
⇒ downsampling and upsampling inside the network



⊖ Application: Human Pose Estimation

- Example: SegNet



## + Up sampling : "Unpacking"

① Nearest neighbor:

$$\begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \Rightarrow \begin{bmatrix} 1 & 1 & 2 & 2 \\ 1 & 2 & 2 & 2 \\ 3 & 3 & 4 & 4 \\ 3 & 3 & 4 & 4 \end{bmatrix}$$

Simplest version

Problem: blocky output structure

② "Bed of Nails":

$$\begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \Rightarrow \begin{bmatrix} 1 & 0 & 2 & 0 \\ 0 & 0 & 0 & 0 \\ 3 & 0 & 4 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

Problem: fixed location for upsampled stimuli

③ Max unpooling: remember which element was max

$$\begin{bmatrix} 1 & 2 & 6 & 3 \\ 3 & 5 & 2 & 1 \\ 1 & 2 & 2 & 1 \\ 7 & 3 & 4 & 8 \end{bmatrix}$$

Max Pooling

$$\begin{bmatrix} 5 & 6 \\ 7 & 8 \end{bmatrix}$$

→ - . →

$$\begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}$$

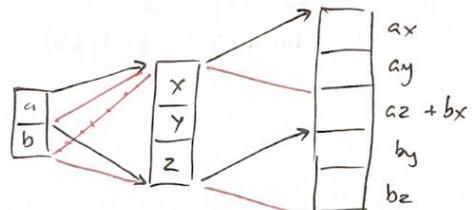
$$\begin{bmatrix} 0 & 0 & 2 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 3 & 0 & 0 & 4 \end{bmatrix}$$

Max Unpooling

(use corresponding pairs)

④ Learnable up sampling: Transpose Convolution

- Conv as matrix multiplication:

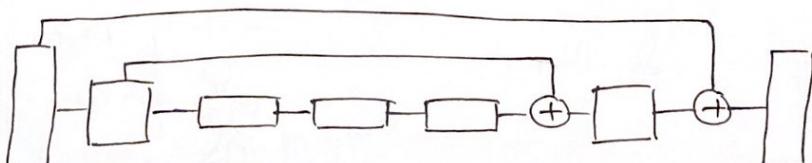


- Also called:

[ Up convolution  
Fractionally strided convolution  
Backward strided convolution ]

But not ~~deconvolution~~

+ Extension: Skip Connections  $\Rightarrow$  Example: U-Net



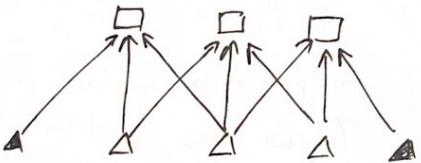
## +)Mask R-CNN:

- Object Detection + Semantic Segmentation = Instance Segmentation  
⊕ Faster-R-CNN + FCN = Masked R-CNN

- RoI Align vs RoI Pooling (Fast RNN)  
 Solve with bilinear interpolation ← Problem: loss of data misalignment of pixels

## +)Atrous Convolutions: (Dilated)

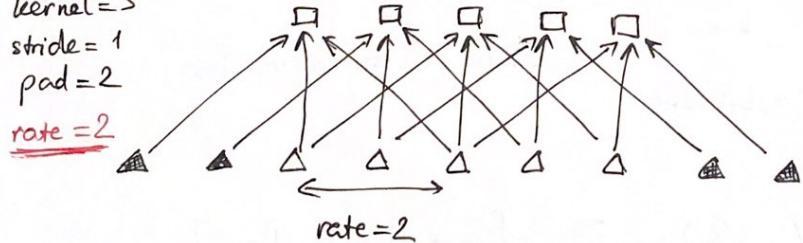
kernel=3  
stride=1  
pad=1



- increase receptive field without increasing computation
- like the effect of pooling without losing resolution

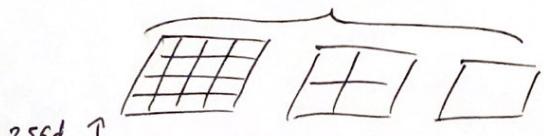
r as dilation factor

kernel=3  
stride=1  
pad=2  
rate=2



## +)Spatial Pyramid Pooling: sampling at many levels then stack all together.

$$16 \times 256 + 4 \times 256 + 1 \times 256 \Rightarrow \text{FCs}$$



⇒ lets the network decides which pooling scales to use, or even, how to combine different scales..

## $\Rightarrow \rightarrow$ Atrous Spatial Pyramid Pooling: (ASPP)

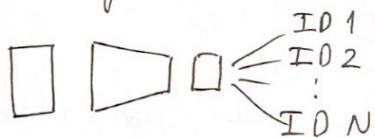
Apply Dilated Conv with  $r = 6, 12, 18, 24 \dots$  then stack the output

L16/

## ⊗ CNNs for Matching

3 Types of Models used for Matching Tasks

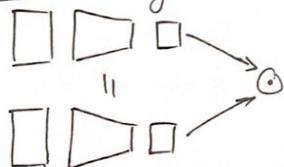
### - Identification:



### Training

Multi Class Classification Loss

### - Embedding

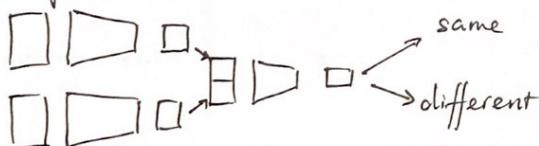


Siamese Network

Large-margin loss

or Triplet loss

### - Verification:



2-Class Classification loss

$$\rightarrow \text{Triplet loss} : L_{\text{tri}}(\theta) = \sum_{a,p,n} [m + D_{a,p} - D_{a,n}]_+$$

$$y_a = y_p \neq y_n$$

$D_{a,p}$  : distance between  $a \times n$

$[..]_+$  :  $\max(.., 0)$

- If we choose randomly, Anchor & Negative are probably already very different

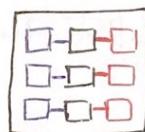
⇒ Need to : Mining Hard      | Triplets  
Medium-hard

### → Practical Implementation:

- Offline Hard Triplet Mining: (OHM)

- Embed data with  $f_\theta \Rightarrow$  for all data  $\Rightarrow$  very wasteful  
 ↗ Mine Hard triplets  
 - Update Embedding  $f_\theta$   
 high computational cost

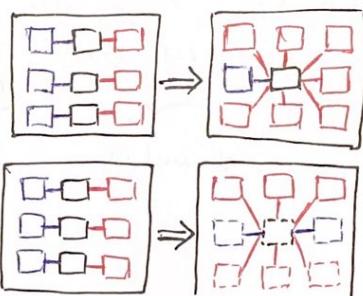
Then do mini-batch learning



- Online Hard Triplet Mining

Idea: 1) Other members of triplets as negative candidates

2) Triplets are constructed only within batch



- Batch All:  $d_{j,a,n}^{i,a,p} = D(f_\theta(x_a^i), f_\theta(x_p^i)) - D(f_\theta(x_a^i), f_\theta(x_n^j))$

$$L_{BA}(\theta; X) = \underbrace{\sum_{i=1}^P \sum_{a=1}^K}_{\text{all anchors}} \underbrace{\sum_{p=1}^K}_{\substack{p \neq a \\ p \neq a}} \underbrace{\sum_{j=1}^P \sum_{n=1}^K}_{\substack{j \neq i \\ j \neq i}} [m - d_{j,a,n}^{i,a,p}]_+$$

$P$  classes  
 each  $K$  images

$P \cdot K$        $K-1$        $(P-1) \cdot K$

Use all combinations of triplet within batch

- ### - Batch Hard:

For each anchor, within batch,  
select the hardest pos & neg.

- ## - Lifted Embedding:

consider all but the anchor-pos. pair  
as negatives

$$L_L(\theta; x) = \sum_{(a,p) \in X} \left[ D_{a,p} + \log \sum_{\substack{n \in X \\ n \neq a \\ n \neq p}} (e^{m - D_{a,n}} + e^{m - D_{p,n}}) \right]_+$$

- Generalized Lifted Embedding: batch all + lifted embedding

- (\*) Batch hard performs best  
additional improvements with soft margin

- Matching applications: Face Identification very good  $\Rightarrow$  privacy?  
Stereo Matching

# 3D Reconstruction

L17,

## (\*) Geometric Vision:

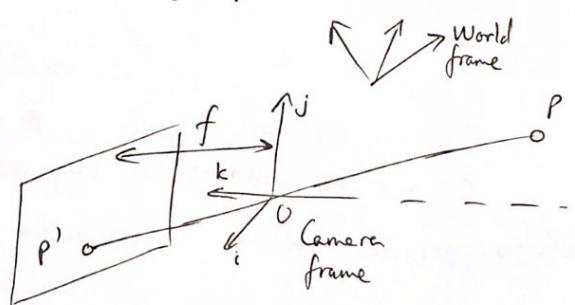
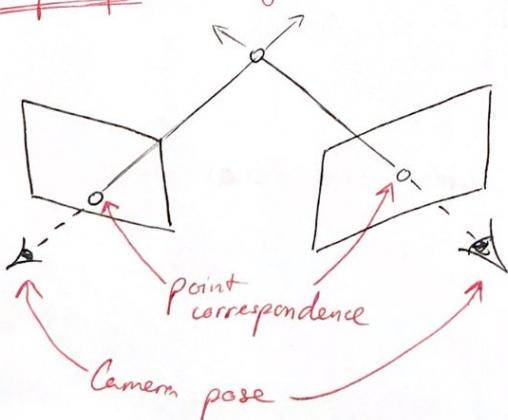
### + Visual Cues (Details)

- Shading
  - Texture
  - Focus
  - Perspective
  - Motion
- } The focus of lecture

### + Stereo Vision: process of extracting 3D information from multiple 2D views of a scene

## (\*) Epipolar geometry: the geometry of stereo vision

Basic principle: Triangulation



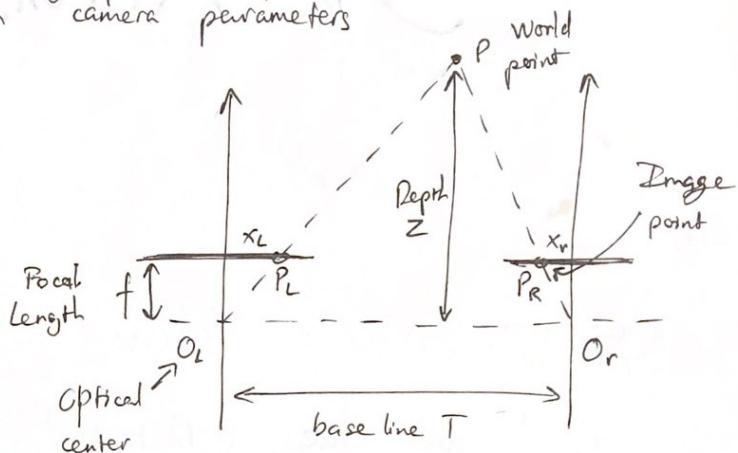
- Extrinsic params: camera frame  $\Leftrightarrow$  world frame rotation matrix, translation vector
- Intrinsic params: image coord.  $\Leftrightarrow$  pixel coord. focal length, pixel size, ...

- Simple case: [assuming parallel optical axes known camera parameters]

$$\frac{T - (x_r - x_l)}{2-f} = \frac{T}{Z}$$

$$\Rightarrow Z = f \frac{T}{x_r - x_l}$$

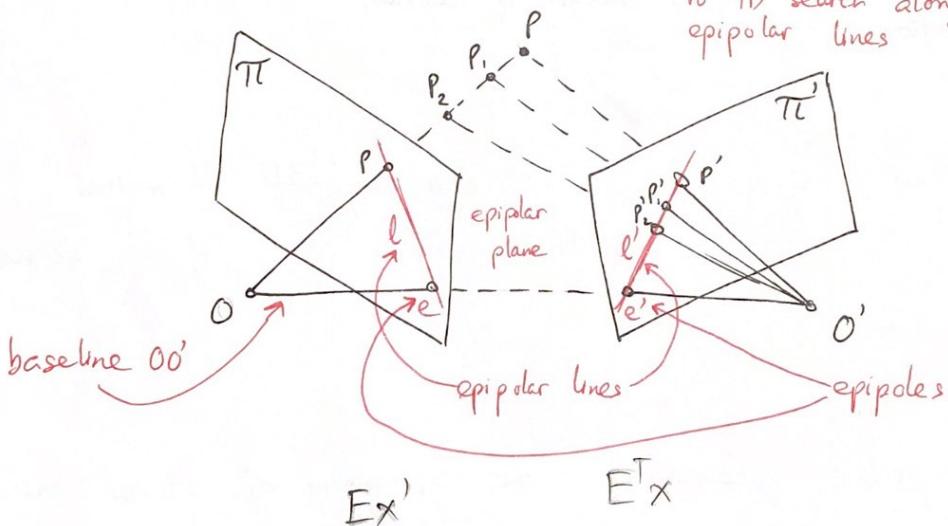
$x_r - x_l$ : disparity



$\Rightarrow$  We can estimate depth of image point from focal length f, base line T & disparity ( $x_r - x_l$ )

- Stereo correspondence constraints:

reduces correspondence problem to 1D search along conjugate epipolar lines



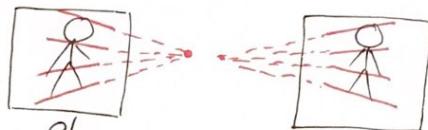
$$Ee' = 0 = E^T e$$

$E$ : as essential matrix :  $E = T_x \cdot R$

$e, e'$ : as positions of epipoles

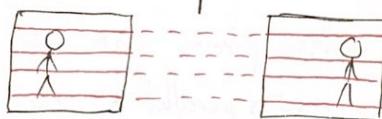
$\Leftrightarrow$  epipoles will be eigen vector of  $E$  with eigen values of 0

- Example:
  - ↳ Converging cameras: as position of 3D point varies, epipolar lines rotate about the baselines



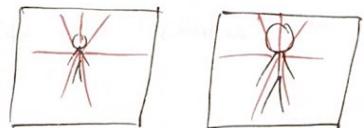
- 2, Motion Parallel with Image Plane:

the epipolar lines will be all parallel

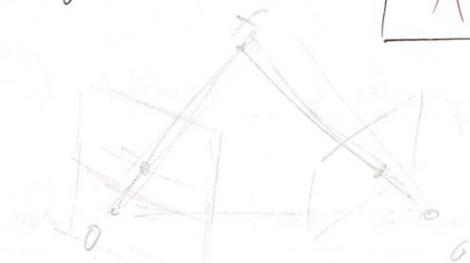


- 3, Forward Motion: epipoles are at the same coordinates in both images.

Epipolar lines converge like a star



### + Algebra stuffs:



$$X' = \underbrace{RX}_{\text{rotation matrix}} + \underbrace{T}_{\text{translational vector}} ; \quad \vec{a} \times \vec{b} = [a_x] \cdot \vec{b}$$

$$\downarrow \quad \quad \quad \downarrow \\ \text{skew matrix} \quad \text{corresponding to } \vec{a} \begin{bmatrix} 0 & -a_z & a_y \\ a_z & 0 & -a_x \\ -a_y & a_x & 0 \end{bmatrix}$$

$$\Rightarrow T \times X' = T \times RX + T \times T = T \times RX$$

$$\Rightarrow X' \cdot (T \times X') = X' \cdot (T \times RX)$$

$$\Rightarrow X' \cdot (T \times RX) = 0$$

$$\Rightarrow X' \cdot (T_x \cdot RX) = 0 \Rightarrow \text{Set } E = T_x \cdot R$$

$$\Rightarrow X'^T \cdot E \cdot X = 0 \text{ holds true for all rays } p \rightarrow p' \text{ in the same direction as } X \text{ & } X'$$

essential matrix  
(extrinsic infos)

$l' = E_p$  : epipolar line for point  $p$

$l = E^T p'$  :

$$\begin{bmatrix} x' \\ y' \\ 1 \end{bmatrix}^T - E \begin{bmatrix} x \\ y \\ 1 \end{bmatrix} = 0$$

$$\Leftrightarrow \begin{bmatrix} x' \\ y' \\ 1 \end{bmatrix} \cdot \begin{bmatrix} a & b & c \\ g & h & i \\ 1 & 1 & 1 \end{bmatrix} = 0$$

$$ax' + by' + c = 0$$

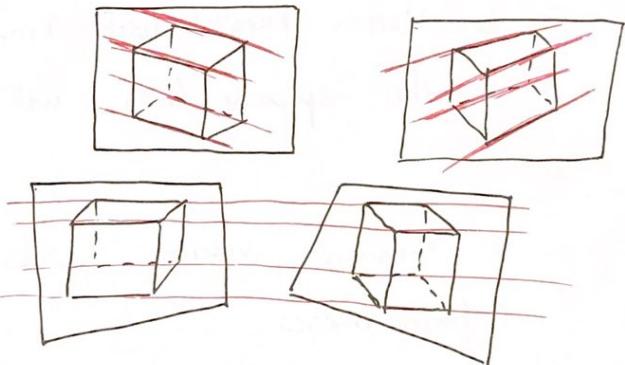
$\Rightarrow$  Epipolar constraint:

$$p'^T \cdot E \cdot p = 0$$

### + Stereo Image Rectification

Reproject images planes onto a common plane, that is parallel to the baseline

$\Rightarrow$  Scanlines are epipolar lines



### + Correspondence Search (Matching)

In practice,  
use both

#### Dense Correspondence Search

- For each pixel, find correspondence
- Easy when epipolar lines are scanlines (apply rectification)

#### Sparse Correspondence Search

- Only for a set of detected features
- Use feature description

(Harris ?)

#### Pros

- Simple process
- More depth  $\Rightarrow$  useful for surface reconstruction

- Efficiency
- Can have more reliable matches
- Less sensitive to illumination robust

#### Cons

Problem with textureless regions  
different view points

- Have to know enough to pick good features
- Sparse information

## ⊕ Stereo Reconstruction

- Main steps:
- Calibrate cameras
  - Rectify images
  - Compute disparity
  - Estimate depth

This is just ideal case.

- what if, how can we get extrinsic info from calibration?
- when triangulation failed?

⊕ Matrix stuff: Solve  $\boxed{Ax=0}$

$$A = U \cdot D \cdot V^T$$

$$A = U \begin{bmatrix} d_{11} & & \\ & \ddots & \\ & & d_{NN} \end{bmatrix} \begin{bmatrix} v_{11} \\ \vdots \\ v_{NI} \end{bmatrix} = \begin{bmatrix} v_{1N} \\ \vdots \\ v_{NN} \end{bmatrix}^T$$

- Solution of  $A \cdot x = 0$  is the null space vector of  $A$
- corresponds to the smallest (last) singular vector of  $A$

## ⊕ Camera Calibration:

⊕ Camera Parameters:

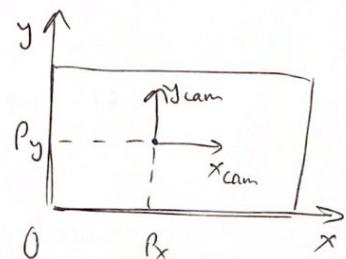
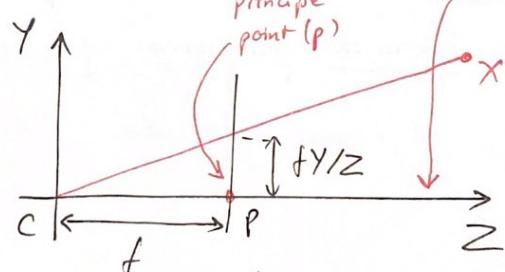
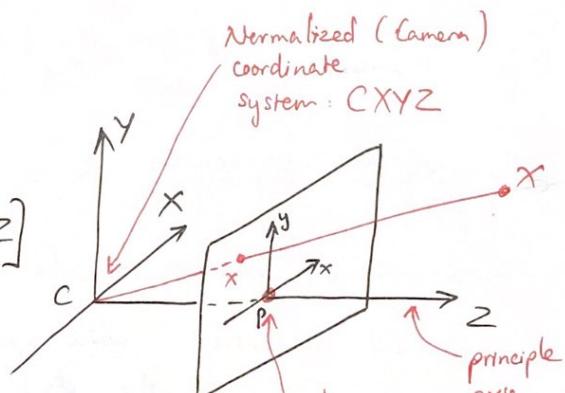
$$[x, y, z] \rightarrow [fx/z, fy/z]$$

$$\begin{bmatrix} fx \\ fy \\ z \end{bmatrix} = \begin{bmatrix} f & 0 \\ 0 & f \\ 1 & 0 \end{bmatrix} \begin{bmatrix} x \\ y \\ 1 \end{bmatrix}$$

$$\Leftrightarrow Z \cdot x = p_x - x_{cam}$$

$$\text{unknown normalization factor}$$

$$= \begin{bmatrix} f & 0 \\ 0 & f \\ 1 & 0 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \end{bmatrix} \cdot x_{cam}$$



- The principle point is not always in the middle of the image

$$\Rightarrow [x, y, z] \rightarrow [fx/z + px, fy/z + py]$$

### - Pixel Coordinates:

$$K = \begin{bmatrix} m_x & & \\ & m_y & \\ & & 1 \end{bmatrix} \cdot \begin{bmatrix} f & & p_x \\ & f & p_y \\ & & 1 \end{bmatrix} = \begin{bmatrix} \alpha_x & & x_0 \\ & \alpha_y & y_0 \\ & & 1 \end{bmatrix}$$

### - Camera Rotation & Translation:

$$\Rightarrow X_{cam} = \begin{bmatrix} \tilde{X}_{cam} \\ 1 \end{bmatrix} = \begin{bmatrix} R & -R\tilde{C} \\ 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} \tilde{X} \\ 1 \end{bmatrix} = \begin{bmatrix} R & -R\tilde{C} \\ 0 & 1 \end{bmatrix} \cdot X$$

$$\rightarrow Z_x = K[I|0] \cdot X_{cam} = K \cdot [R \quad | -R\tilde{C}] \cdot X \\ = K \cdot [R \quad | \quad t] \cdot X \\ = P \cdot X$$

Summary:

$$\Rightarrow \text{Intrinsic parameters: } K = \begin{bmatrix} m_x & & \\ & m_y & \\ & & 1 \end{bmatrix} \cdot \begin{bmatrix} f & s & p_x \\ & f & p_y \\ & & 1 \end{bmatrix} = \begin{bmatrix} \alpha_x & s & x_0 \\ \alpha_y & y_0 \\ 1 \end{bmatrix}$$

- S: skew, non rectangular pixels, in complicated settings..
  - Radial distortion: while opening angle is more than  $45^\circ$   
 $\Rightarrow$  the effect is no more negligible.

+ Extrinsic parameters:  $[R | t]$   $\Rightarrow$  Rotation + Translation

+), All: Camera Projection Matrix:  $P = K[R|t]$

- ⊗ DoF:
  - General pinhole camera: 9 DoF ( $3_{int} + 6_{ext}$ )  $p_x, p_y, f$
  - CCD Camera with square pixels: 10 DoF ( $m_x = m_y \Rightarrow 4_{int}$ )
  - \_\_\_\_\_ non \_\_\_\_\_. 11 DoF ( $m_x \neq m_y$ )
  - General camera: 11 DoF ( $m_x = m_y$ )

## Calibrating a Camera:

- Use known "calibration object"
- Algorithm for checkerboard pattern:
  - ① Perform Canny Edge Detector
  - ② Fit straight lines
  - ③ Intersect lines to obtain corners

Very good  
localization  
 $< 1/10$  pixel

- ④ Rule of thumb: number of constraints should exceed number of unknowns by a factor of 5

Ex: 11 params of  $P \Rightarrow$  should have 55 constraints

each point : 2 constraints  $\rightarrow$  should have  $55/2$  points

$$\lambda x_i = P x_i \Leftrightarrow x_i \times P x_i = 0 \Leftrightarrow \begin{bmatrix} 0 & -x_i^T & y_i x_i^T \\ x_i^T & 0 & -x_i x_i^T \\ -y_i x_i^T & x_i x_i^T & 0 \end{bmatrix} \begin{bmatrix} P_1 \\ P_2 \\ P_3 \end{bmatrix} = 0$$

$$\Leftrightarrow \begin{bmatrix} 0^T & x_1^T & -y_1 x_1^T \\ x_1^T & 0^T & -x_1 x_1^T \\ \dots & \dots & \dots \\ 0^T & x_n^T & -y_n x_n^T \\ x_n^T & 0^T & -x_n x_n^T \end{bmatrix} \begin{bmatrix} P_1 \\ P_2 \\ P_3 \end{bmatrix} = 0$$

need calibration points  
in more than one plane

or else, we will get  
degenerate solution

- After that, still need "matrix decomposition problem"  
to go from  $P \Rightarrow$  intrinsic, extrinsic parameters

+ Revisiting Triangulation: Find original  $\mathbf{x}$

→ Use Linear Algebraic Approach

$$\begin{bmatrix} \mathbf{x}_{1x} \end{bmatrix} \cdot P_1 \mathbf{x} = 0 \quad \Rightarrow \text{Stack}$$

$$\begin{bmatrix} \mathbf{x}_{2x} \end{bmatrix} \cdot P_2 \mathbf{x} = 0 \quad \Rightarrow \text{Solve with SVD}$$

2, Nonlinear Approach:  $\arg \min_{\mathbf{x}} [d^2(\mathbf{x}_1, P_1 \mathbf{x}) + d^2(\mathbf{x}_2, P_2 \mathbf{x})]$   
 { most accurate  
 but no closed-form solution, only iterative algorithm}

+ Revisiting Epipolar Geometry:

1, Calibrated Case:  $\mathbf{x}^T \cdot \mathbf{E} \cdot \mathbf{x}' = 0$ ;  $\mathbf{E} = [\mathbf{t}_x]^T \cdot \mathbf{R}$

- $\mathbf{E}$ : rank 2, singular, cause  $\mathbf{e}'$ ,  $\mathbf{e}$  is in the null space of  $\mathbf{E}$   
 $\mathbf{E}\mathbf{e}' = 0$ ,  $\mathbf{E}^T \mathbf{e} = 0$
- 5 dof, up to scale ( $\mathbf{p}_x, \mathbf{p}_y, m_x = m_y$ , skew, scale?)

2, Uncalibrated Case:  $\mathbf{x} = K\hat{\mathbf{x}}$ ,  $\mathbf{x}' = K'\hat{\mathbf{x}}'$

$$\Rightarrow \mathbf{x}^T \cdot K'^{-T} \cdot \mathbf{E} \cdot K^{-1} \cdot \mathbf{x}' = 0 \Leftrightarrow \mathbf{x}^T \cdot F \cdot \mathbf{x}' = 0$$

- Epipolar lines:  $F\mathbf{x}' \propto F^T \mathbf{x}$ .   
Fundamental matrix
- $F$ : singular, rank 2, 7 dof

⊗ Estimate F with Eight Point Algorithm

$$\begin{aligned} \mathbf{x} &= [u, v, 1]^T \\ \mathbf{x}' &= [u', v', 1]^T \end{aligned} \quad \Rightarrow \text{Stack 8 points}$$

$$\Rightarrow \text{Use SVD}$$

$$\Rightarrow \begin{bmatrix} u \\ v \\ 1 \end{bmatrix}^T \begin{bmatrix} F_{11} & F_{13} \\ \vdots & \vdots \\ F_{31} & F_{33} \end{bmatrix} \begin{bmatrix} u' \\ v' \\ 1 \end{bmatrix} = 0$$

If noisy data  
 $\Rightarrow F$  full rank 3, no epipoles  
 $\Rightarrow$  Use rank 2 estimation of  $F$  with SVD

$$\Rightarrow \begin{bmatrix} u, v, \dots, u, v, 1 \end{bmatrix} \begin{bmatrix} F_{11} \\ \vdots \\ F_{33} \end{bmatrix} = 0$$

~~Problem~~ Problem: Poor numerical conditioning  $u'u \gg v \gg 1$

⇒ The Normalized Eight-Point Algorithm:

1, Center image data at origin

Scale so mean squared distance = 2 pixels

2, Compute F

3, Enforce rank 2 constraint with SVD on F

4,

Set smallest singular value to 0

Nonlinear Least-Squared Approach:

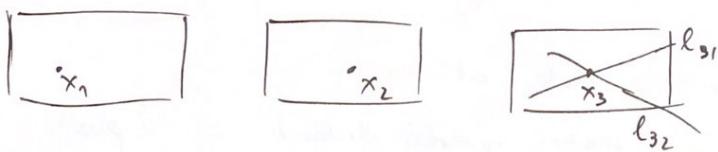
$$\text{minimize} \quad \sum_{i=1}^N [d^e(x_i, Fx'_i) + d^2(x'_i, F^T x_i)]$$

+ 3D Reconstruction with Weak Calibration:

to estimate world geometry without requiring calibrated camera  
only from a redundant set of correspondences

- Procedure:
  - Find interest point
  - Calculate correspondences
  - Compute epipolar geometry
  - Refine

## + Epipolar Transfer:



- Given projections of same point in 2 images:  $x_1$  &  $x_2$

- Failed when the 2 lines stack on each others