

# F12 solutions (2013)

RC<sub>1</sub>

Q1 (a)  $S(x, y) = I(x, y) * g_{\sigma}(x, y)$  where  $g_{\sigma}(x, y) = \frac{1}{2\pi\sigma^2} e^{-\frac{(x^2+y^2)}{2\sigma^2}}$

$$= I(x, y) * g_{\sigma}(x) * g_{\sigma}(y) \quad \text{where } g_{\sigma}(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{x^2}{2\sigma^2}}$$

$$= \sum_{-n}^n \sum_{-n}^n I(u, v) g_{\sigma}(x-u) g_{\sigma}(y-v)$$

$$= \sum_{-n}^n \sum_{-n}^n I(x-u, y-v) g_{\sigma}(u) g_{\sigma}(v)$$

Show how to compute  $g_{\sigma}(x)$  from  $\frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{x^2}{2\sigma^2}}$  curve eq.  $2n+1=7$  for  $\sigma=1$ . by sampling up to tails. (20%)

(b)  $S(x, y, \sigma^2)$ , family of blurred images with continuous  $\sigma$

(see backwork overleaf)

$$S(x, y, t) = g_{\sigma}(x, y) * I(x, y) \quad \text{for } t = \sigma^2$$

R2

## Q1 (b) Image pyramids and blob-detection (scale-space)

Need to generate a discrete set of images with difference amount of blur. We sample  $f(x, y, \sigma^2)$ , logarithmically spaced;

$$\sigma_i = 2^{\frac{i}{s}} \sigma_0, \quad \sigma_{i+1} = 2^{\frac{1}{s}} \sigma_i \quad (1)$$

with  $s$  images per octave (i.e. after  $s$  images,  $\sigma$  has doubled).

— apply incremental blur (gaussian  $\sigma_k$ ) between images in octave to get images with increasing amount of blur.

$$g(\sigma_{k+1}) = g(\sigma_i) * g(\sigma_k) \quad \sigma_{k+1}^2 = \sigma_i^2 + \sigma_k^2 \text{ at } \sigma_{i+1} = 2^{\frac{1}{s}} \sigma_i$$

$$\sigma_k = \sigma_i \sqrt{2^{\frac{2}{s}} - 1} \quad (2)$$

— Each blur is performed as 2 1D convolutions (see (a) ii).

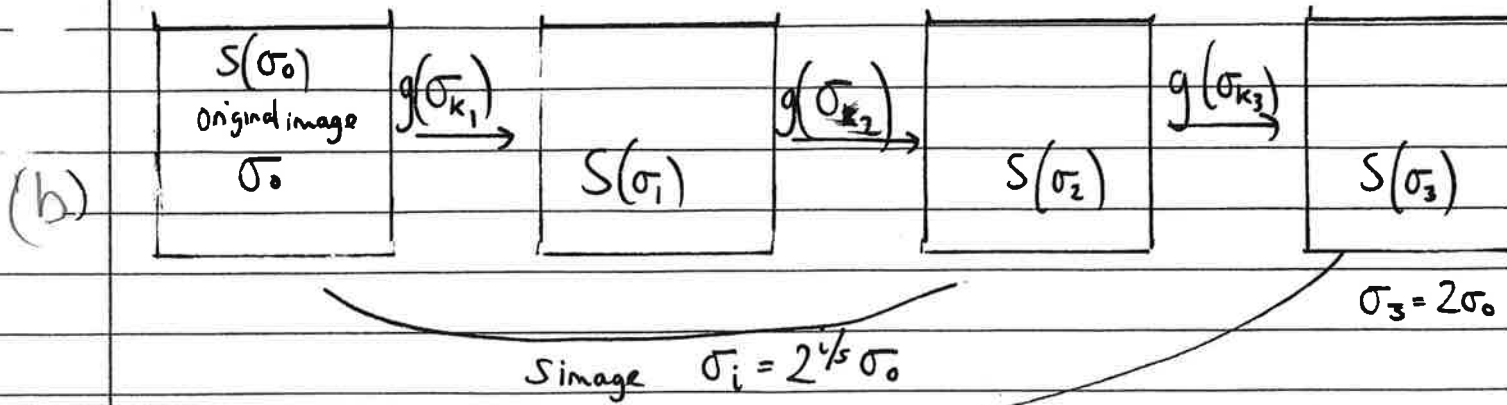
— After scale has doubled <sup>(no. of blurs)</sup>, resize image by subsampling by 2 (i.e.  $\frac{1}{2}$  size images). We can represent blurred images with fewer pixels without loss of information (Nyquist). [biggest saving] (3)

— some (small) incremental blur kernels used in each octave.

$\sigma_{k_1}, \sigma_{k_2}, \sigma_{k_3}$  etc. on sub-sampled images,

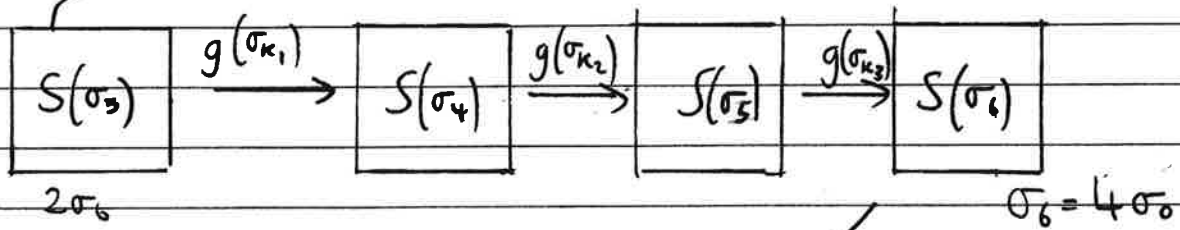
but really represent filtering with larger kernels,  $2\sigma_{k_1}, \dots, 4\sigma_{k_1}, \dots, 8\sigma_{k_1}$ .

1st octave  $\sigma_0 \rightarrow 2\sigma_0$

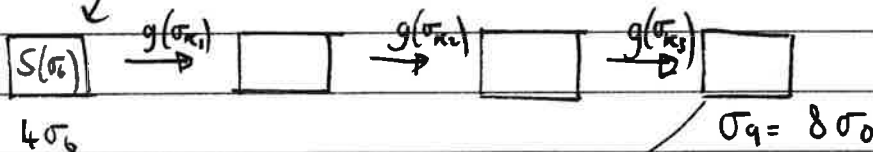


2nd octave  $2\sigma_0 \rightarrow 4\sigma_0$

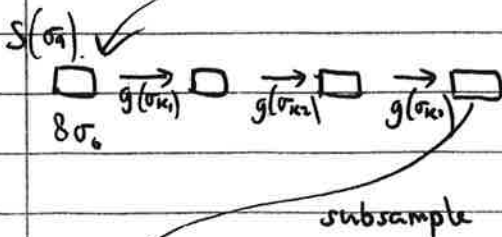
sub-sample to  $\frac{1}{4}$  size



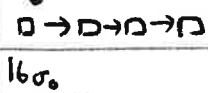
3rd octave  $4\sigma_0 \rightarrow 8\sigma_0$



4th octave  $8\sigma_0 \rightarrow 16\sigma_0$



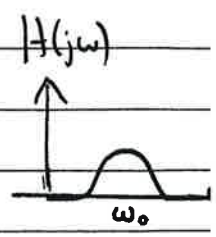
5th octave  $16\sigma_0 \rightarrow 32\sigma_0$



16 pixel images

# 11 (c) Band-pass filtering

- tuned to a small band of spatial frequencies  
(difference of high pass and low pass filtered o/p)



$$-\nabla^2 G_{\sigma_i} * I = \nabla^2 S(\sigma_i^2) \approx S(\sigma_{i+1}^2) - S(\sigma_i^2)$$

i.e. subtract neighbouring images in same octave.

- generate a pyramid of DOG images by subtraction; in same octave (2)  
(neighbours in "same" octave).

(c) Blobs are localised at max/min of  $\nabla^2(G_{\sigma} * I)$  response.  
Need to search over  $\nabla^2 S(x, y, \sigma_i^2)$  for (max/min) local

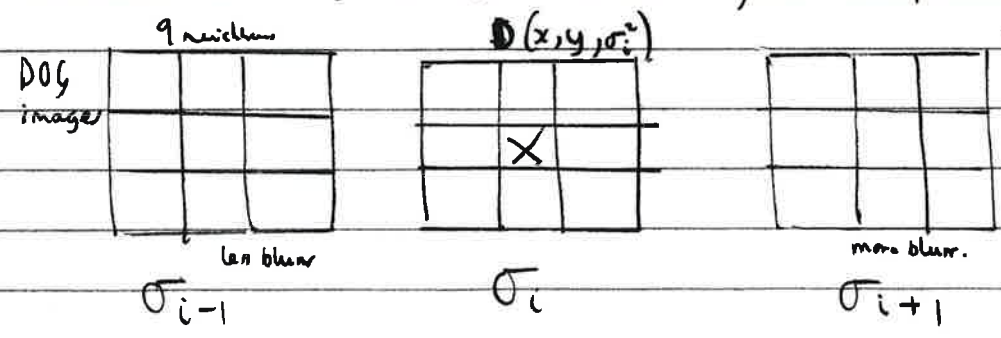
- More efficient to search over difference of gaussian pyramid (DOG)

$$\nabla^2 S \approx D(x, y, \sigma_i^2) \approx S(x, y, \sigma_{i+1}^2) - S(x, y, \sigma_i^2)$$

(10%)

for max/min in  $x, y$  and  $\sigma$ .

- Evaluate 26 neighbours of  $D(x, y, \sigma_i^2)$  to see if local max/min



(30%)

- Local max/min is blob location; <sup>(x, y)</sup> scale (size of feature) is  $\sigma_i$

- SIFT descriptor looks at 16x16 pixels sampled from  $S(x, y, \sigma_i^2)$  in correct octave.

(d) cont Orientation: Look at gradients in  $16 \times 16$  patch and bin (histogram at  $10^\circ$  intervals). Find dominant orientation after smoothing histogram.

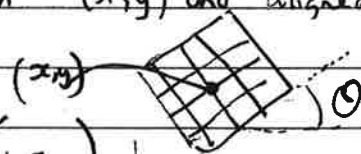
2b

## Q1e) SIFT, matching and classification

(a) - For each interest point at location  $(x, y)$  and scale  $\sigma$  estimate the dominant orientation,  $\theta$ , by looking at histogram of edge gradients from  $\nabla S(x, y, \sigma^2)$ .  
Bins  $10^\circ$  apart + smoothed by gaussian. (36 bins)

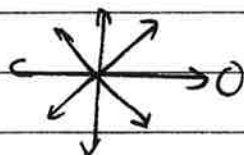


(b) - Sample  $16 \times 16$  gradients from  $(x, y)$  and aligned with  $\theta$  at image of scale,  $\sigma$



- Smooth these gradients with  $g(1.5\sigma)$  to emphasize gradients at interest point.

- produce orientation histogram <sup>(HOG)</sup> for each  $4 \times 4$  block (cell)  
Each bin records grad. magnitude interpolated over 8 dir's (quadrants) only.



- concatenate 16 histograms of gradients (HOGs) to vector of 128D

- normalize to unit length; truncate to 0.2 to avoid illumination effects, and normalize to 1.

- SIFT encodes 2D shape - invariant to lighting by using edges + normalization step + to exact 2D position by histogram/bin effects, [pooling]  
Edges taken around a blob in centre.

(20/4)

Q2 (a) Consider  $x_i = \frac{f X_i}{Z_i}$  and let  $X_i = a_i + \lambda_i b$

In lim. as  $Z_i \rightarrow \infty$

$$\begin{aligned} x_{\infty} &= \frac{f X_i}{Z_i} = \frac{f (a_i + \lambda_i b)}{a_i + \lambda_i b} \\ &= \frac{f (b_1, b_2)}{b_3} \quad \text{and indep. of } a_i \end{aligned}$$

i.e. All have common vanishing pt if //.

(20%)

Q2 (b)  $X_c = R X + T = \left[ \begin{array}{c|c|c} R & T & X \\ \hline 0 & 1 & 1 \end{array} \right]$  In this example  $R=I$   
 $T=0$

ii) perspective

$$\begin{bmatrix} su \\ sv \\ s \end{bmatrix} = \begin{bmatrix} f & 0 & 0 & 0 \\ 0 & f & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} X_c \\ Y_c \\ Z_c \\ 1 \end{bmatrix}$$

It is used to represent homogeneous co-ordinates

iii) CCD scaling

$$\begin{aligned} u &= u_0 + k_u \frac{f X_c}{Z_c} \\ v &= v_0 + k_v \frac{f Y_c}{Z_c} \end{aligned}$$

$$\begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = \begin{bmatrix} k_u & 0 & u_0 \\ 0 & k_v & v_0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} X_c \\ Y_c \\ Z_c \end{bmatrix}$$

Combine to get projection matrix.

Assumptions: pin-hole, no-non linear distortion, planar perspective

(20%)

$$2(c) \text{ (i)} \begin{bmatrix} s_u \\ s_v \\ s \end{bmatrix} = \begin{bmatrix} k_u & 0 & u_0 \\ 0 & k_v & v_0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} f & 0 & 0 & 0 \\ 0 & f & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} x_c \\ y_c \\ z_c \\ 1 \end{bmatrix}$$

$$u = \frac{k_u f x_c}{z_c} + u_0$$

$$v = \frac{k_v f y_c}{z_c} + v_0$$

$$\text{while } \begin{bmatrix} s_A u_A \\ s_A v_A \\ s_A \end{bmatrix} = \begin{bmatrix} k_u & 0 & u_0 \\ 0 & k_v & v_0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} f & 0 & 0 & 0 \\ 0 & f & 0 & 0 \\ 0 & 0 & 0 & z_A \end{bmatrix} \begin{bmatrix} x_c \\ y_c \\ z_c \\ 1 \end{bmatrix}$$

$$\therefore u_A = \frac{k_u f x_c}{z_A} + u_0$$

$$v_A = \frac{k_v f y_c}{z_A} + v_0 \quad (30\%)$$

$$\therefore (u - v_A) = \frac{k_u f x_c}{z_c} \left( 1 - \frac{z_c}{z_A} \right) = \frac{k_u f x_c}{z_c} \left( \frac{z_A - z_c}{z_A} \right)$$

$$v - v_A = \frac{k_v f y_c}{z_c} \left( 1 - \frac{z_c}{z_A} \right) = \frac{\Delta z}{z_A} (v - v_0) \quad (10\%)$$

$$= \frac{\Delta z}{z_A} (v - v_0)$$

(ii) good approximation at center of image of for small variation in depth,  $\Delta z$  (20%)  
 (iii) Calibrate from 4 pts by linear least-squares, (if more non-coplanar pts)

(3)

$$(a) \begin{matrix} 3 \times 1 \\ \begin{bmatrix} su \\ sv \\ s \end{bmatrix} \end{matrix} = \begin{matrix} 3 \times 3 \\ \begin{bmatrix} k \end{bmatrix} \end{matrix} \begin{matrix} 3 \times 4 \\ \begin{bmatrix} R & | & T \end{bmatrix} \end{matrix} \begin{matrix} 4 \times 1 \\ \begin{bmatrix} x \\ y \\ z \\ 1 \end{bmatrix} \end{matrix} \text{ in general}$$

Assume  $Z=0$ ,

$$\begin{matrix} 3 \times 1 \\ \begin{bmatrix} su \\ sv \\ s \end{bmatrix} \end{matrix} = \begin{matrix} 3 \times 3 \\ \begin{bmatrix} k \end{bmatrix} \end{matrix} \begin{matrix} 3 \times 3 \\ \begin{bmatrix} r_1 & r_2 & | & T \\ \vdots & \vdots & & \end{bmatrix} \end{matrix} \begin{matrix} 3 \times 1 \\ \begin{bmatrix} x \\ y \\ \cancel{z} \\ 1 \end{bmatrix} \end{matrix} \text{ for planar object}$$

$3 \times 3$

(b) (i) 4

(ii) RANSAC. — random sample of 4 pts

— estimate  $t_{ij}$ — check for inliers by Euclidean distance at  $u' = Au$ 

— accept or correct matches if large # inliers

(iii). Set up  $2N$  equations from  $N$  correspondences

$$u' = \frac{h_{11}u + h_{12}v + h_{13}}{h_{31}u + h_{32}v + h_{33}}$$

$$v' = \frac{h_{21}u + h_{22}v + h_{23}}{h_{31}u + h_{32}v + h_{33}}$$

$$h_{21}u + h_{22}v + h_{23} = v'(h_{31}u + h_{32}v + h_{33})$$

$$h_{31}u + h_{32}v + h_{33}$$

Normalize data to unit variance and zero mean.

Set up  $A\bar{h} = 0$  and solve by SVD.

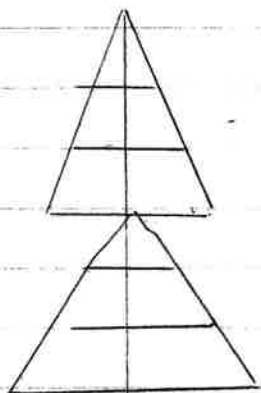


(c) Need to have  $K$  so as to recover  $R$  and  $t$ .

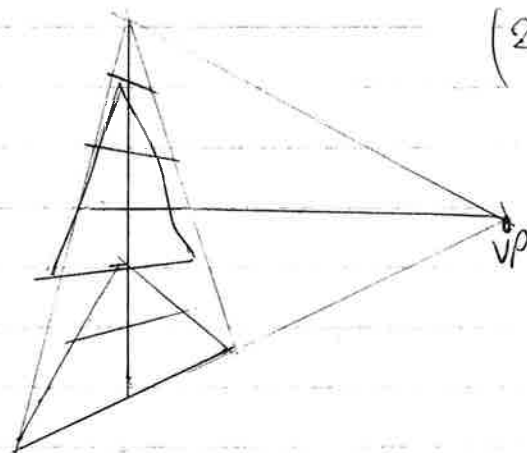
(d) Consider lines of symmetry as parallel lines eg.  $y=0, y=a^i e_k$   
Under perspective consider lines/pls at  $D$

$$\begin{bmatrix} u' \\ v' \\ 1 \end{bmatrix} = \begin{bmatrix} \quad \\ \quad \\ \quad \end{bmatrix} \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} = \begin{bmatrix} h_{11} \\ h_{21} \\ h_{31} \end{bmatrix}$$

$$\begin{aligned} u'_{VP} &= h_{11}/h_{31} \\ v'_{VP} &= h_{21}/h_{31} \end{aligned}$$



bilateral  
symmetry is lost



4a) The geometry of a stereo camera constrains each point-feature in one image to lie on a corresponding epipolar line in the other image. Epipolar lines meet at the epipole: this is the image of one camera's optical centre in the other camera's image plane. There are two epipoles - one for each image.

The fundamental matrix relates points in the left & right images of a stereo pair:

$$\underline{w}'^T \underline{F} \underline{w} = 0 \quad \text{where } \underline{w}'^T = [u, v, 1]$$

$\underline{w}$  = points pixel coordinates in the left image and

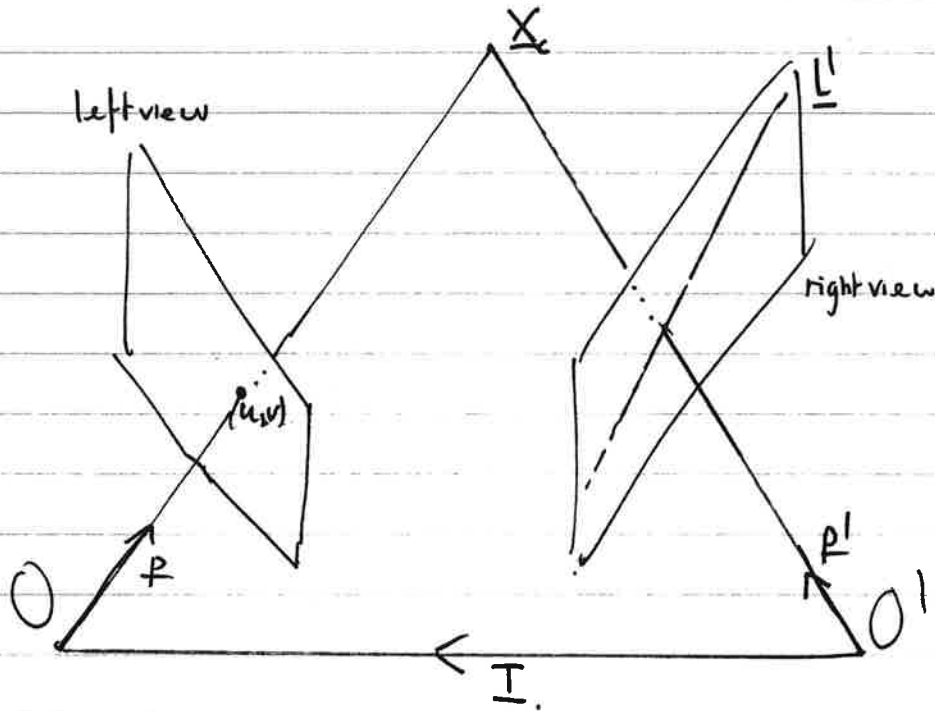
$\underline{w}'$  are the corresponding coordinates in the right image

$\underline{F}$  has zero determinant & can be determined only up to a scale factor.

Epipolar lines can be derived from the fundamental matrix.

4(b)

$$\underline{x}'_c = R \underline{x}_c + \underline{I}$$



$$\underline{\omega} = K p \quad \text{and} \quad \underline{\omega}' = K' p'$$

$$\underline{x}' = R \underline{x} + T$$

$$\underline{x}' \cdot (\underline{I}_x R \underline{x}) = 0 \quad \text{by coplanarity of } [\underline{x}', \underline{x}, \underline{I}]$$

$$p'^T \underline{I}_x R p = 0$$

$$\therefore \underline{\omega}'^T K'^{-T} \underline{I}_x R K^{-1} \underline{\omega} = 0$$

$$\underline{\omega}'^T F \underline{\omega} = 0 \quad \text{where} \quad \underline{F} = K'^{-T} \underline{I}_x R K^{-1}$$

$$\therefore \underline{L}' = \underline{K}'^{-T} \underline{I}_x R \underline{K}^{-1} \underline{\omega}$$

4(c). A pt on an epipolar line in right view corresponds to  $\underline{w}$  in left. lies on a line  $\underline{L}'$  such that

$\underline{w}' \cdot \underline{L}' = 0$  where line is represented by  $\begin{pmatrix} l_1 \\ l_2 \\ l_3 \end{pmatrix}$  and  $\underline{w}'$

But  $\underline{w}' \cdot F \underline{w} = 0$  and hence  $\underline{L}' = F \underline{w}$ .  
where  $\underline{L}' = K^{l-T} T_x R K^{-l} \underline{w}$  (20%)

(d) Each pair of correspondences gives 1 equation in unknown  $F$  element  $f_{ij}$

$\dots A \underline{f} = 0$  and solve by least-squares

Matrix  $F$  has max rank 2 ( $\det F = 0$ ). Enforce by finding best matrix  $F^*$  to estimate  $\hat{E}$  by minimizing Frobenius norm. (20%)

e)  $E$  can be decomposed into  $\underline{e}'_x(M)$  where  $\underline{e}'$  is an epipole in right vi or if internal cameras are known it can be decomposed into  $\underline{t}$  and  $\underline{R}$

$P_i = K [1 | 0]$  and  $\underline{P}_i = K [R | t]$  (20%)

need to know  $K$  and  $K'$  of cameras