



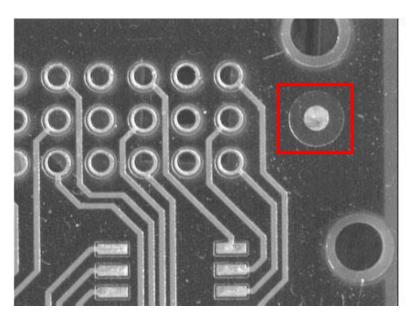
# **Object Detection**

Luigi Di Stefano (luigi.distefano@unibo.it)

#### Introduction (1)



 The object detection problem occurs in countless applications and can be formulated as follows. Given a reference image (aka model image) of the sought object, determine whether the object is present or not in the image under analysis (aka target image) and, in case of detection, estimate the pose of the object.





 Depending on the application, the pose may often be given by a translation, a roto-translation or a similarity (roto-translation plus scale).

#### Introduction (2)



- The problem has indeed a number of diverse facets, such as seeking to detect a single or multiple instances of the model as well as either a single or multiple models. For the sake of simplicity, and without much loss of generality, hereinafter we will refer mainly to the basic "singlemodel/single instance" setting. Generalization to any of the other variants turns out more often than not straightforward.
- Typical nuisances to be dealt with are intensity changes, occlusions and clutter. Moreover, computational efficiency is a major requirement in most practical applications. In particular, it is harder to achieve efficiency as long as the dimensionality of the pose space and/or the number of models gets larger.
- Many different object detection approaches have been proposed in the computer vision literature. Here, we will focus on template matching (aka pattern matching), shape-based matching, the Hough Transform and detection by local invariant features.

#### Template Matching



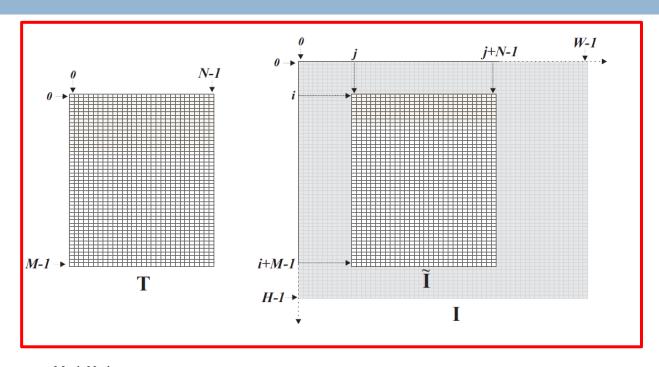
The model image is slid across the target image to be compared at each position



to an equally sized window by means of a suitable (dis)similarity function

## (Dis)Similarity Functions (1)





$$SSD(i,j) = \sum_{m=0}^{M-1} \sum_{n=0}^{N-1} \left( I(i+m,j+n) - T(m,n) \right)^2$$
 (Sum of Squared Differences)

Both  $\widetilde{I}(i, j)$  (the window at position (i, j) of the target image having the same size as T) as well as T can be thought of as  $M \cdot N$ -dimensional vectors. Accordingly, the SSD represents the squared L2 (Euclidean) norm of their difference.

#### (Dis)Similarity Functions (2)



$$SAD(i,j) = \sum_{m=0}^{M-1} \sum_{n=0}^{N-1} \left| I(i+m,j+n) - T(m,n) \right|$$
 (Sum of Absolute Differences)

The SAD represents the L1 norm of the difference between vectors  $\tilde{l}(i,j)$  e T.

$$NCC(i, j) = \frac{\sum_{m=0}^{M-1} \sum_{n=0}^{N-1} I(i+m, j+n) \cdot T(m, n)}{\sqrt{\sum_{m=0}^{M-1} \sum_{n=0}^{N-1} I(i+m, j+n)^{2}} \cdot \sqrt{\sum_{m=0}^{M-1} \sum_{n=0}^{N-1} T(m, n)^{2}}}$$
(Normalised Cross-Correlation)

$$NCC(i,j) = \frac{\tilde{\mathbf{I}}(i,j) \cdot \mathbf{T}}{\|\tilde{\mathbf{I}}(i,j)\| \cdot \|\mathbf{T}\|} = \frac{\|\tilde{\mathbf{I}}(i,j)\| \cdot \|\mathbf{T}\| \cdot \cos \theta}{\|\tilde{\mathbf{I}}(i,j)\| \cdot \|\mathbf{T}\|} = \cos \theta \qquad \qquad \qquad \tilde{\mathbf{I}}(i,j) = \alpha \cdot \mathbf{T}$$
(Invariant to linear intensity changes)

The NCC represents the cosine of the angle between vectors  $\tilde{I}(i,j)$  e T.

## (Dis)Similarity Functions (3)



#### (Zero-Mean Normalised Cross-Correlation, Correlation Coefficient)

$$\mu(\tilde{I}) = \frac{1}{MN} \sum_{m=0}^{M-1} \sum_{n=0}^{N-1} I(i+m, j+n)$$

$$\mu(T) = \frac{1}{MN} \sum_{m=0}^{M-1} \sum_{n=0}^{N-1} T(m,n)$$

#### La NCC is computed after subtraction of the means:

$$I(i+m, j+n) \to \left(I(i+m, j+n) - \mu(\tilde{I})\right) \qquad T(m,n) \to \left(T(m,n) - \mu(T)\right)$$

$$ZNCC(i, j) = \frac{\sum_{m=0}^{M-1} \sum_{n=0}^{N-1} \left(I(i+m, j+n) - \mu(\tilde{I})\right) \cdot \left(T(m,n) - \mu(T)\right)}{\sqrt{\sum_{m=0}^{M-1} \sum_{n=0}^{N-1} \left(I(i+m, j+n) - \mu(\tilde{I})\right)^{2}} \cdot \sqrt{\sum_{m=0}^{M-1} \sum_{n=0}^{N-1} \left(T(m,n) - \mu(T)\right)^{2}}}$$

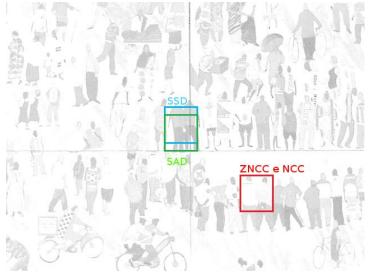
$$\widetilde{\mathbf{I}}(i,j) = \alpha \cdot \mathbf{T} + \beta$$

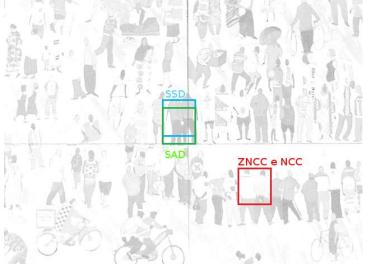
(Invariant to affine intensity changes)

## Exemplar comparison under significant intensity changes







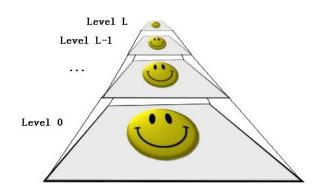


**ZNCC** turns out a similarity function very robust to intensity changes!

#### Fast Template Matching



- Template matching may turn out exceedingly slow whenever the model and/or target images have a large size (i.e. computational complexity is O(M×N×W×H)).
- Development of fast template matching algorithms is therefore an active research topic. A popular approach is to deploy an image pyramid:



Very fast approach, though the number of levels needs to be chosen carefully (and empirically) to avoid loss of information detrimental to detection.

- Smoothing and sub-sampling (typically 1/2 on both sides at each level).
- Full search at top level and then local refinements traversing back the pyramid down to original image.



# A taxonomy of fast template matching approaches



- Approximate Methods: the result is not guaranteed to be the same as would be provided by a full-search.
  - Image Pyramid
    Full-search at top level, then local refinement throughout
    successive levels down to the original full-size image.
  - Sub-Templates
    - A relatively small sub-template (e.g. a salient region or a set of salient points) is sought first in order to highlight a set of candidate positions, then the full-size comparison is carried out at these positions only.
- Exact (aka Exhaustive) Methods: the result is guaranteed to be the same as would be provided by a full-search.
  - Call-out with dissimilarity functions (e.g. SAD-SSD): computation of the function at the current position is terminated as soon as its value gets higher than the current minimum (or a threshold).
  - Matching in the Fourier domain (FFT).
  - Deployment of efficiently computable bounds of the (dis)similarity function to skip candidates which cannot improve the current degree of matching.

## Template Matching via FFT (1)



 Due to the Convolution Theorem, template matching can be carried out in the Fourier domain. Indeed, in case of the NCC function:
 Box-filtering

$$NCC(i,j) = \frac{1}{K(i,j)} \Big( T(i,j) \circ I(i,j) \Big), \quad K(i,j) = \left\| \tilde{I}(i,j) \right\| \cdot \left\| T \right\|$$

$$= \frac{1}{K(i,j)} \Big( T(i,j) \circ I(i,j) - \left\| \tilde{I}(i,j) \right\| \cdot \left\| T \right\|$$

$$= \frac{1}{K(i,j)} \left( T(i,j) \circ I(i,j) - \left\| \tilde{I}(i,j) \right\| \cdot \left\| T \right\|$$

$$= \frac{1}{K(i,j)} \left( T(i,j) \circ I(i,j) - \left\| \tilde{I}(i,j) \right\| \cdot \left\| T \right\|$$

$$= \frac{1}{K(i,j)} \left( T(i,j) \circ I(i,j) - \left\| \tilde{I}(i,j) \right\| \cdot \left\| T \right\| \right)$$

$$= \frac{1}{K(i,j)} \left( T(i,j) \circ I(i,j) - \left\| \tilde{I}(i,j) \right\| \cdot \left\| T \right\| \right)$$

$$= \frac{1}{K(i,j)} \left( T(i,j) \circ I(i,j) - \left\| \tilde{I}(i,j) \right\| \cdot \left\| T \right\| \right)$$

$$= \frac{1}{K(i,j)} \left( T(i,j) \circ I(i,j) - \left\| \tilde{I}(i,j) \right\| \cdot \left\| T \right\| \right)$$

$$= \frac{1}{K(i,j)} \left( T(i,j) \circ I(i,j) - \left\| \tilde{I}(i,j) \right\| \cdot \left\| T \right\| \right)$$

$$= \frac{1}{K(i,j)} \left( T(i,j) \circ I(i,j) - \left\| \tilde{I}(i,j) \right\| \cdot \left\| T \right\| \right)$$

$$= \frac{1}{K(i,j)} \left( T(i,j) \circ I(i,j) - \left\| T(i,j) - \left\| T(i,j) - \left\| T(i,j) \circ I(i,j) - \left\| T(i,j) -$$

$$I(i,j) * T(-i,-j) \Leftrightarrow \Im \Big( I(i,j) \Big) \cdot \underbrace{\Im \Big( T(-i,-j) \Big)}_{= \Im \Big( I(i,j) \Big)} = \underbrace{\Im \Big( I(i,j) \Big)}_{= \Im \Big( I(i,j) \Big)} = \underbrace{I(u,v)}_{= \Im \Big( I(i,j) \Big)}_{= \Im \Big( I(i,j) \Big)} = \underbrace{I(u,v)}_{= \Im \Big( I(i,j) \Big)}_{= \Im \Big( I(i,j) \Big)} = \underbrace{I(u,v)}_{= \Im \Big( I(i,j) \Big)}_{= \Im \Big( I(i,j) \Big)} = \underbrace{I(u,v)}_{= \Im \Big( I(i,j) \Big)}_{= \Im \Big( I(i,j) \Big)} = \underbrace{I(u,v)}_{= \Im \Big( I(i,j) \Big)}_{= \Im \Big( I(i,j) \Big)}_{= \Im \Big( I(i,j) \Big)} = \underbrace{I(u,v)}_{= \Im \Big( I(i,j) \Big)}_{= \Im \Big( I(i,j) \Big)}_{= \Im \Big( I(i,j) \Big)} = \underbrace{I(u,v)}_{= \Im \Big( I(i,j) \Big)}_{= \Im \Big( I(i,j) \Big)}_{=$$

IFFT

$$NCC(i, j) = \frac{1}{K(i, j)} \mathfrak{I}^{-1} \left( I(u, v) \cdot \overline{T}(u, v) \right)$$

FFT

## Template Matching via FFT (2)



As far as the SSD is concerned:

$$SSD(i, j) = \sum_{m=0}^{M-1} \sum_{n=0}^{N-1} \left( I(i+m, j+n) - T(m, n) \right)^{2}$$

$$\sum_{m=0}^{M-1} \sum_{n=0}^{N-1} I(i+m, j+n)^{2} + \sum_{m=0}^{M-1} \sum_{n=0}^{N-1} T(m, n)^{2} - 2 \sum_{m=0}^{M-1} \sum_{n=0}^{N-1} I(i+m, j+n) \cdot T(m, n)$$

$$\|\tilde{I}(i, j)\|^{2} + \|T\|^{2} - 2 \left( T(i, j) \circ I(i, j) \right)$$
IFFT

FFT

FFT

#### Template Matching via FFT (3)



Computational complexity of Template Matching in the signal domain:

$$O(M \cdot N \cdot W \cdot H)$$

Computational complexity of Template Matching in the Fourier domain:

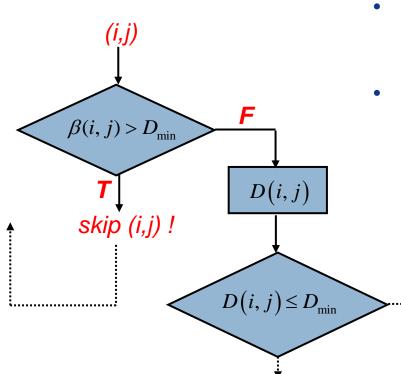
$$O(W \cdot H \cdot \log_2(W \cdot H))$$
 Independent of template size!

 Usually, Template Matching turns out faster in the Fourier domain when template size is large. Nonetheless, it is worth pointing out that template matching in the Fourier domain requires *floating-point* arithmetic operations which are often exceedingly slow in embedded architectures without FPU (*Floating Point Unit*).

#### **Bound-based methods**



• Let D(i,j) be the *dissimilarity* function to be minimized to detect T and let us assume that there exists a *lower bound* of this function:  $\beta(i,j) \leq D(i,j)$ .



- The final outcome is the same as it would have been provided by a full-search!
- The approach is faster than a full-search provided that  $\beta$  (i,j) can be calculated much more efficiently than D(i,j) (efficient bound) and it turns out enough accurate to skip the computation of D(i,j) at a large number of image positions (effective bound).
  - In case of template matching based on similarity functions,
     β (i,j) must be an upperbound: β(i,j) ≥ S(i,j).

Efficient and effective bounds have been proposed for SAD, SSD, NCC and ZNCC.

#### Successive Elimination Algorithm

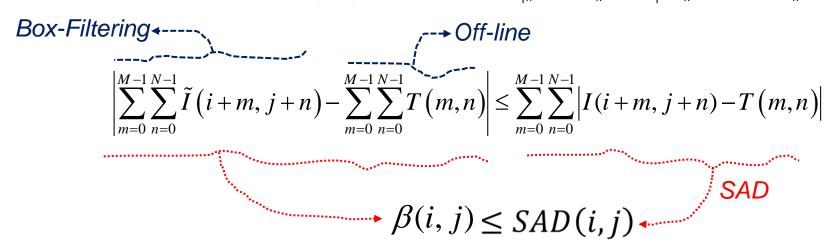


$$\mathbf{X} = \begin{bmatrix} x_1 \\ \vdots \\ x_{N} \end{bmatrix} \Rightarrow \|\mathbf{X}\|_p = \left(\sum_{i=1}^N |x_i|^p\right)^{1/p}$$

Triangle Inequality:

$$\left\| \left\| \mathbf{X} \right\|_{p} - \left\| \mathbf{Y} \right\|_{p} \right| \leq \left\| \mathbf{X} - \mathbf{Y} \right\|_{p}$$

$$p = 1$$
,  $\mathbf{X} = \tilde{I}(i, j)$ ,  $\mathbf{Y} = T$   $\longrightarrow$   $\|\tilde{\mathbf{I}}(i, j)\| - \|\mathbf{T}\| \le \|\tilde{\mathbf{I}}(i, j) - \mathbf{T}\|$ 

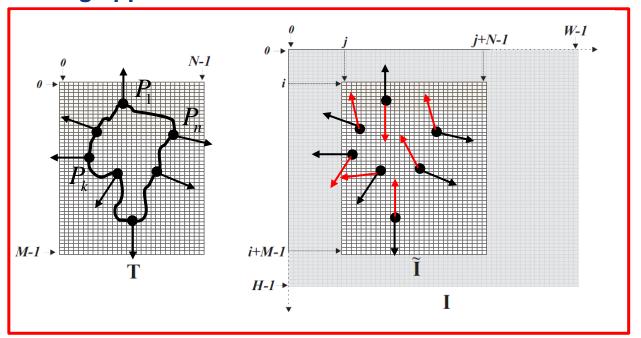


SEA was originally proposed (Lee & Salari, TIP, 1995) for the *block matching* step carried out to estimate motion in video compression, then it has been used widely also for template matching. An advanced and more effective approach was been proposed later (Tombari, Mattoccia & Di Stefano, PAMI 2009).

#### **Shape-based Matching**



• Shape-based Matching (Steger, 2002) may be thought of as an <u>edge-based</u> template matching approach.



- First, a set of control points,  $P_k$ , is extracted from the model image by an Edge Detection operation (e.g. Sobel) and the gradient direction at each  $P_k$  is recorded.
- Then, at each position (i,j) of the target image, the recorded gradient directions associated with control points are compared to those at their <u>corresponding</u> image points,  $P_k(i,j)$ , in order to compute a similarity function.

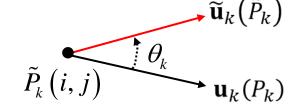
#### Similarity Function



$$\mathbf{G}_k(P_k) = \begin{bmatrix} I_{\chi}(P_k) \\ I_{y}(P_k) \end{bmatrix}, \mathbf{u}_k(P_k) = \frac{1}{\|\mathbf{G}_k(P_k)\|} \begin{bmatrix} I_{\chi}(P_k) \\ I_{y}(P_k) \end{bmatrix}, k = 1..n$$

$$\widetilde{\mathbf{G}}_{k}(\widetilde{P}_{k}) = \begin{bmatrix} I_{x}(\widetilde{P}_{k}) \\ I_{y}(\widetilde{P}_{k}) \end{bmatrix}, \widetilde{\mathbf{u}}_{k}(\widetilde{P}_{k}) = \frac{1}{\|\widetilde{\mathbf{G}}_{k}(\widetilde{P}_{k})\|} \begin{bmatrix} I_{x}(\widetilde{P}_{k}) \\ I_{y}(\widetilde{P}_{k}) \end{bmatrix}, k = 1..n$$

$$S(i,j) = \frac{1}{n} \sum_{k=1}^{n} \mathbf{u}_{k}(P_{k}) \cdot \widetilde{\mathbf{u}}_{k}(\widetilde{P}_{k}) = \frac{1}{n} \sum_{k=1}^{n} \cos \theta_{k}$$



- The above defined similarity function spans the interval [-1; 1]. It takes its maximum value when all the gradients at the control points in the current window of the target image are perfectly aligned to those at the control points of the model image.
- Choosing a detection threshold,  $S_{min}$ , can be thought of as specifying the fraction of model points which must be seen in the image to trigger a detection.

#### More robust similarity functions



 Certain application settings call for invariance to global inversion of contrast polarity along object's contours, as the object may appear either darker or brighter than the background in the target image. This kind of invariance can be achieved by a slight modification to the similarity function defined previously:

$$S(i,j) = \frac{1}{n} \left| \sum_{k=1}^{n} \mathbf{u}_{k}(P_{k}) \cdot \widetilde{\mathbf{u}}_{k}(\widetilde{P}_{k}) \right| = \frac{1}{n} \left| \sum_{k=1}^{n} \cos \theta_{k} \right|$$

• The following function is even more robust due to the ability to withstand local contrast polarity inversions:

$$S(i,j) = \frac{1}{n} \sum_{k=1}^{n} \left| \mathbf{u}_{k}(P_{k}) \cdot \widetilde{\mathbf{u}}_{k}(\widetilde{P}_{k}) \right| = \frac{1}{n} \sum_{k=1}^{n} \left| \cos \theta_{k} \right|$$

#### Main Properties



- Steger's similarity function is based on gradient directions only: as such it is invariant to affine intensity changes.
- A peculiar trait of Shape-based Matching consists in edges (i.e. control points) being detected in the model image (7) only, never in the target image (1). This renders the method very robust wrt the potential fragility of carrying out edge detection on the target image (e.g. how to set thresholds under varying lighting conditions?)
- In principle, the similarity function is potentially robust to occlusion. Let us assume that a fraction of T is occluded in I and that gradient orientations are uniformly distributed across the occluding surface (as it would be, e.g., should the occluder be texture-less). Thus, the cosines of the angles between the gradients at the control points and their corresponding image points are uniformly distributed alike, such that the expected value of the sum of the contributions to the similarity function due to occluded control points is null. Accordingly, as already mentioned, the detection threshold determines the degree of occlusion that one wishes to withstand in the addressed application.
- In practise, a small threshold related to gradient magnitude at image points is enforced, so as to ensure nullifying the contributions due to uniformly textured occluders.

#### Speeding-up the search



• *Call-Out*: given the chosen detection threshold  $S_{min}$  and the already computed partial similarity,  $S_p(i,j)$ , it is possible to check whether the computation of the similarity function S(i,j) needs to continue or not:

$$S_p(i,j) = \sum_{k=1}^p \mathbf{u}_k(P_k) \cdot \widetilde{\mathbf{u}}_k(\widetilde{P}_k) \longrightarrow S(i,j) \le \frac{1}{n} \Big( S_p(i,j) + (n-p) \Big)$$

$$S(i,j) < S_{min} \qquad \frac{1}{n} \Big( S_p(i,j) + (n-p) \Big) < S_{min}$$

$$S_p(i,j) + (n-p) < n \cdot S_{min} \qquad \longrightarrow \qquad S_p(i,j) < n \cdot S_{min} + (p-n)$$

If the condition holds, the computation of S(i,j) is terminated

#### Overall search process (1)



- In the original formulation (Steger, 2002), object's pose is given by a similarity (roto-translation plus scale change). The method has then be extended so as to allow estimation of a homography (Hofhauser et.al., 2008).
- The search process follows a pyramidal approach, with the number of levels,  $I_{\text{max}}$ , chosen by the user. In the off-line training phase, the model image (T), is repeatedly scaled according to the chosen pyramidal structure.
- Given the rotation ( $\phi_{\min}...\phi_{\max}$ ) and scale variation ( $\sigma_{\min}...\sigma_{\max}$ ) ranges, all possible rotated and scaled version of the model image are pre-computed for each pyramid level. Steger suggests to rotate and scale the model image and then extract the edges rather than rotate and scale the edges extracted from the model image.
- The rotation and scale quantization steps can be either determined automatically (Borgefors, PAMI, 1988) or chosen by the user. At each next pyramid level, the rotation step is doubled.

#### Overall search process (2)



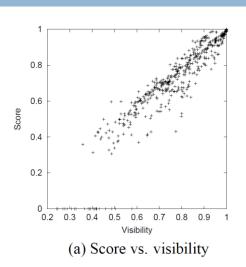
- In the on-line search phase, all the pre-computed scaled and rotated versions of *T* are sought for at the top level of the pyramid by the usual sliding window approach.
- Candidate matches are obtained by *non-maxima suppression* followed by thresholding  $(S_{min})$  of the similarity function, S(i,j), at the found extrema. The previously introduced *call-out* condition is deployed to speed-up the search.
- Candidate matches are then propagated through pyramid levels. At each, a small neighbourhood of the pose estimated at the previous level is explored, so as to either reject the match or validate and propagate it to the next one.
- Once a match is established by a local search in the original image, *I*, the pose estimation resolution is increased by an extrapolation process relying on fitting a second degree polynomial (Steger, PAMI, 1998). This allows translation to be estimated with sub-pixel resolution and scale and rotation with a finer resolution than the chosen quantization steps. Optionally, accuracy may be improved by a final least-squares pose refinement procedure (Steger, 2002).

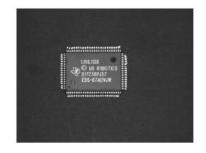
#### Some Results



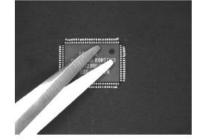












500 images of an IC characterized by different occlusion levels

Smin=0.3  $\rightarrow$ : RR (*Recogniton Rate*): 478/500 (95.6%). If the images featuring an occlusion level higher than > 0.7 are not considered (15 out of 22 *misdetections*), the RR increases to 98.6%.

#### The Hough Transform



- The Hough Transform (HT) enables to detect objects having a known shape (e.g. lines, circles, ellipses..) based on projection of the input data into a suitable space referred to as *parameter* or *Hough* space.
- The HT turns a *global* detection problem into a *local* one.
- The HT is usually applied after an edge detection process (i.e. the actual input data consist of the edge pixels extracted from the original image).
- The HT is robust to noise and allows for detecting the sought shape even though it is partially occluded into the image (up to a certain userselectable degree of occlusion).
- The HT was invented to detect lines and later extended to other analytical shapes (circle, ellipses) as well as to arbitrary shapes (Generalized Hough Transform).
- The GHT principle is widely deployed also within object detection pipelines relying on local invariant features such as, e.g., SIFT.

#### Basic Principle (1)



#### Let's consider first the HT formulation for lines:

$$y - mx - c = 0$$

In the usual image space interpretation of the line equation the parameters  $(\hat{m},\hat{c})$  are fixed

$$y - \hat{m}x - \hat{c} = 0$$

so that the equation represents the mapping  $(1 \to \infty)$  from point  $(\hat{m}, \hat{c})$  of the parameter space to the image points belonging to the line.

However, we may instead fix  $(\hat{x}, \hat{y})$ 

$$\hat{y} - m\hat{x} - c = 0$$

so as to interpret the equation as the mapping  $(1 \rightarrow \infty)$  from image point  $(\hat{x}, \hat{y})$  to the parameter space providing all the lines through the image point.

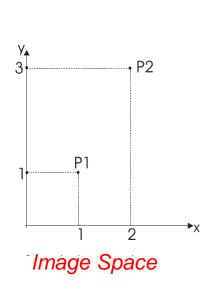
The equation still represents a line, so all the parameter values representing lines through image point  $(\hat{x}, \hat{y})$  lay on a line of the parameter space.

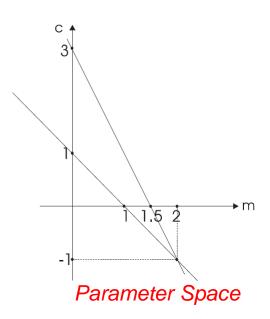
#### Basic Principle (2)



Accordingly, if we consider two image points P<sub>1</sub>,P<sub>2</sub> and map both into the parameter space, we get two lines intersecting at the parameter space point representing the image line through P<sub>1</sub>,P<sub>2</sub>:

$$\begin{cases} \hat{y}_1 - m\hat{x}_1 - c = 0 \\ \hat{y}_2 - m\hat{x}_2 - c = 0 \end{cases} \Longrightarrow \begin{cases} m = \frac{\hat{y}_2 - \hat{y}_1}{\hat{x}_2 - \hat{x}_1} \\ c = \frac{\hat{x}_2\hat{y}_1 - \hat{x}_1\hat{y}_2}{\hat{x}_2 - \hat{x}_1} \end{cases}$$





More generally, if we map n image points we get as many intersections as

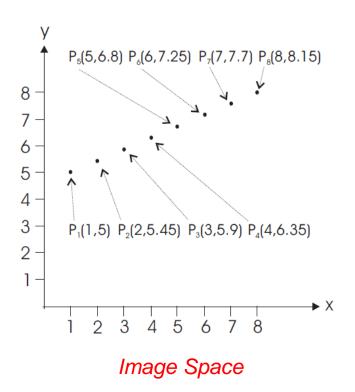
$$\frac{n(n-1)}{2}$$

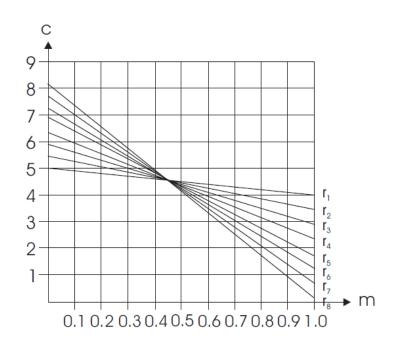
i.e. the number of lines through the n image points.

## Basic Principle (3)



 Considering then n collinear image points, we can notice that their corresponding transforms (i.e. parameter space lines) will intersect at a single parameter space point representing the image line along which such n points lay:





Parameter Space

#### Basic Principle (4)

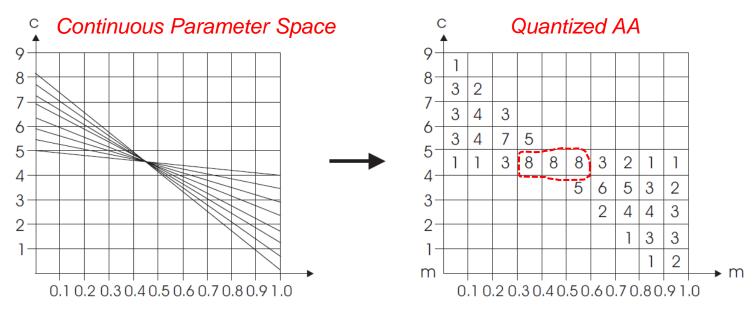


- Therefore, given a sought analytic shape represented by a set of parameters, the HT consists in mapping image points (i.e. usually edge points) so as to create curves into the parameter space of the shape.
- Intersections of parameter space curves indicate the presence of image points
  explained by <u>a certain instance of the shape</u>, the more the intersecting curves the
  more are such image points and thus the higher is the evidence of the presence of
  that instance in the image.
- Detecting objects through the HT consists in finding parameter space points through which many curves do intersect, indeed a *local* rather than *global* detection problem.
- To make it work in practice, the parameter space needs to be quantized and allocated as a memory array, which is often refereed to as *Accumulator Array (AA)*.
- Then, curves are "drawn" into the AA by a so called *voting* process: the transform equation is repeatedly computed to increment the bins satisfying the equation. Accordingly, a high number of intersecting curves at a point of the parameter space will provide a high number of votes into a bin of the AA. Finding parameter space points through which many curves do intersect is thus implemented in practise by finding *peaks* of the AA, i.e. local maxima showing a high number of votes.

#### Basic Principle (5)



Exemplar AA for line detection:



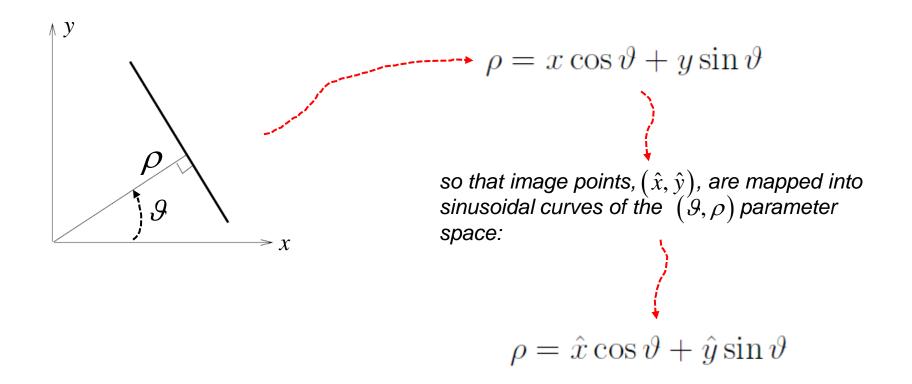
The AA highlights the presence of a line with  $m \in [0.3, 0.6]$ ,  $c \in [4, 5]$ To detect the line more accurately, the AA should be quantized more finely.

 The HT is robust to noise because spurious votes due to noise unlikely accumulate as coherently into a bin as to trigger a false detection. A partially occluded object can still be detected provided that the threshold on the minimum number of votes requited to declare a detection is lowered according to the degree of occlusion to be handled.

#### HT for Line Detection (1)



- The usual line parametrization considered so far (i.e. *y-mx-c=0*) is impractical due to *m* spanning an infinite range.
- The so called normal parametrization is therefore adopted in the HT for lines:



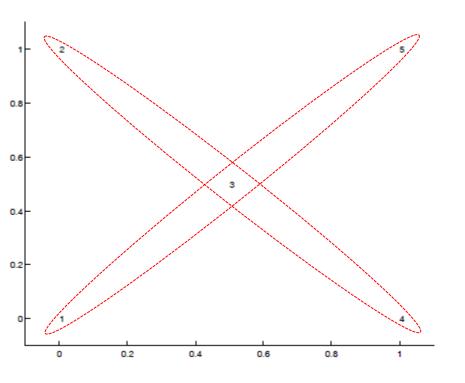
#### HT for Line Detection (2)

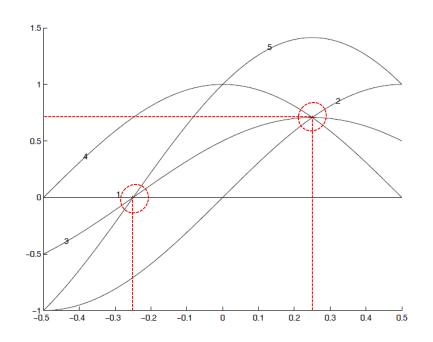


With the normal parametrization: 
$$\vartheta \in \left[ -\frac{\pi}{2}, \frac{\pi}{2} \right], \quad \rho \in \left[ -\rho_{\text{max}}, \rho_{\text{max}} \right]$$

while  $\rho_{max}$  is usually taken as large as the image diagonal:

$$N \times N$$
 pixels  $\rightarrow \rho_{\text{max}} = N \cdot \sqrt{2}$ 

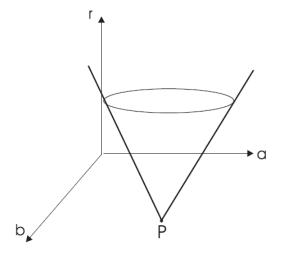




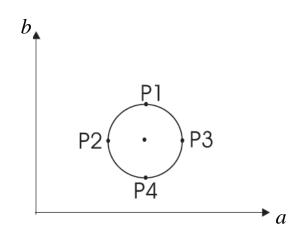
#### HT for Circle Detection (1)



• In the case of circular objects, the equation  $(x-a)^2+(y-b)^2=r^2$  is interpreted as a mapping from the image space (x,y) to the 3-dimensional parameter space (a,b,r). Accordingly, equation  $(\hat{x}-a)^2+(\hat{y}-b)^2=r^2$  provides all the parameter triplets representing circles through image point  $(\hat{x},\hat{y})$ 



 $HT(P(\hat{x}, \hat{y}))$ : cone with vertex at P



Intersection with plane  $r = \overline{r}$ 

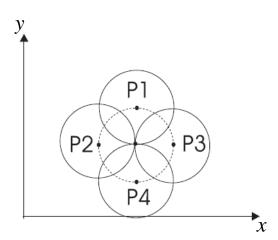


Image circles definied by  $(P_i, \overline{r})$ 

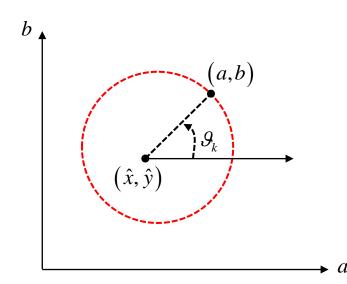
#### HT for Circle Detection (2)



In the practical implementation, a stack of images is adopted as 3D Accumulator Array, each layer of the stack is associated with a radius value. To cast votes, circles are successively "drawn" in each layer:

$$\forall \overline{r} \in [r_{\min}, r_{\max}]: \begin{cases} a = \hat{x} + \overline{r} \cdot \cos \theta_k \\ b = \hat{y} + \overline{r} \cdot \sin \theta_k \end{cases} \qquad \theta_k \in [0, 2\pi]$$

$$\mathcal{G}_k \in [0, 2\pi]$$

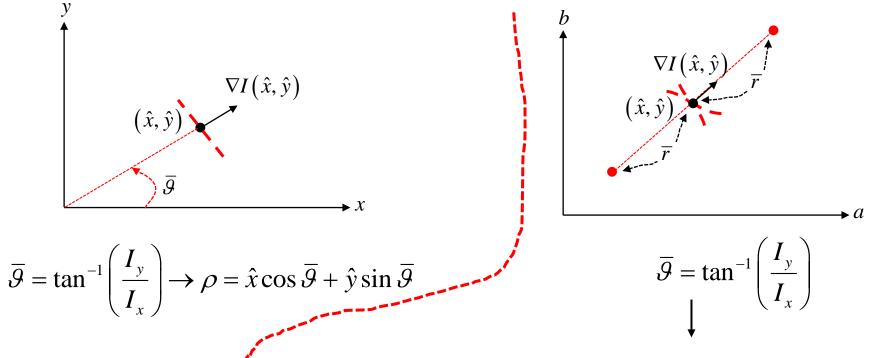


$$\begin{cases} \Delta \mathcal{P} = \frac{2\pi}{M} \\ \mathcal{P}_k = k\Delta \mathcal{P}, \ k = 0...M \end{cases}$$

#### Deployment of gradient information



• The voting phase of the HT can be made faster by deployment of gradient information. In the case of lines (left) and circles (right):

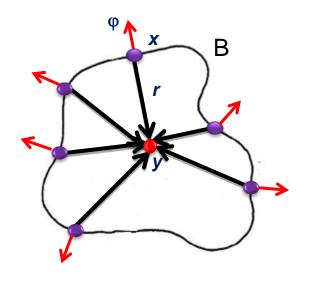


$$\begin{cases} a_1 = \hat{x} + \overline{r} \cdot \cos \overline{\mathcal{G}} \\ b_1 = \hat{y} + \overline{r} \cdot \sin \overline{\mathcal{G}} \end{cases}, \begin{cases} a_2 = \hat{x} + \overline{r} \cdot \cos \left(\overline{\mathcal{G}} + \pi\right) = \hat{x} - \overline{r} \cdot \cos \overline{\mathcal{G}} \\ b_2 = \hat{y} + \overline{r} \cdot \sin \left(\overline{\mathcal{G}} + \pi\right) = \hat{y} - \overline{r} \cdot \sin \overline{\mathcal{G}} \end{cases}$$

#### GHT: Generalized Hough Transform (1)



The HT has been extended to detect arbitrary (i.e. non analytical) shapes:



Off-line Phase (to build the object's model)

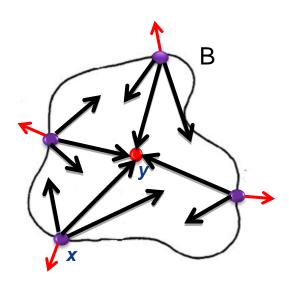
- 1. A reference point y is chosen (e.g. barycentre).
- 2. Gradient direction is quantized according to a chosen step  $\Delta \phi$
- 3. For each point x belonging to object's border B:
  - 1. Compute gradient direction  $\varphi(x)$
  - 2. Compute vector r from y to x (i.e. r = y-x).
- 4. Store r as a function of  $\Delta \varphi$  (R-Table)

i	$\phi_i$	$R_{\phi_i}$
0	0	$\{r y-r=x, x\in B, \phi(x)=0\}$
1	$\Delta\phi$	$\{r y-r=x, x\in B, \phi(x)=\Delta\phi\}$
2	$2\Delta\phi$	$\{r y-r=x, x\in B, \phi(x)=2\Delta\phi\}$

An entry in the R-Table can contain several r vectors.

#### GHT: Generalized Hough Transform (2)





#### On-line Phase (object detection)

1. An image A[y] is iniatialized as accumulator array.

For each edge pixel x of the input image :

- 2. Compute gradient direction  $\varphi$
- 3. Quantize  $\varphi$  to index the R-Table. For each  $r_i$  vector stored into the accessed row:
  - a) Compute the position the reference point y:  $y=x+r_i$
  - b) Cast a vote into the accumulator array:

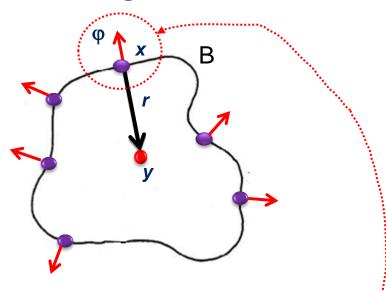
    A[y]++
- 4. As usual with the HT, instances of the sought object are detected by finding peaks of the accumulator array.

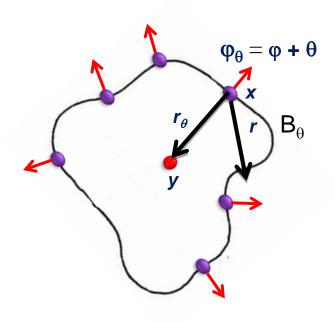
i	$\phi_i$	$R_{\phi_i}$
0	0	$\{r y-r=x, x\in B, \phi(x)=0\}$
1	$\Delta\phi$	r1,r2,r3
2	$2\Delta\phi$	$\{r y-r=x, x\in B, \phi(x)=2\Delta\phi\}$

## GHT: Generalized Hough Transform (3)



#### Handling rotation





$$T_{g}\left[R\left[\varphi\right]\right] = R\left[\left(\varphi_{\theta} - \theta\right) \bmod 2\Pi\right] \xrightarrow{r} \Rightarrow \forall r : y = x + r_{\theta} = x + ROT\left(r, \theta\right), A\left[y, \theta\right] + r + r_{\theta} = x + r_{\theta} = x + ROT\left(r, \theta\right), A\left[y, \theta\right] + r + r_{\theta} = x + r_{$$

#### Handling scale

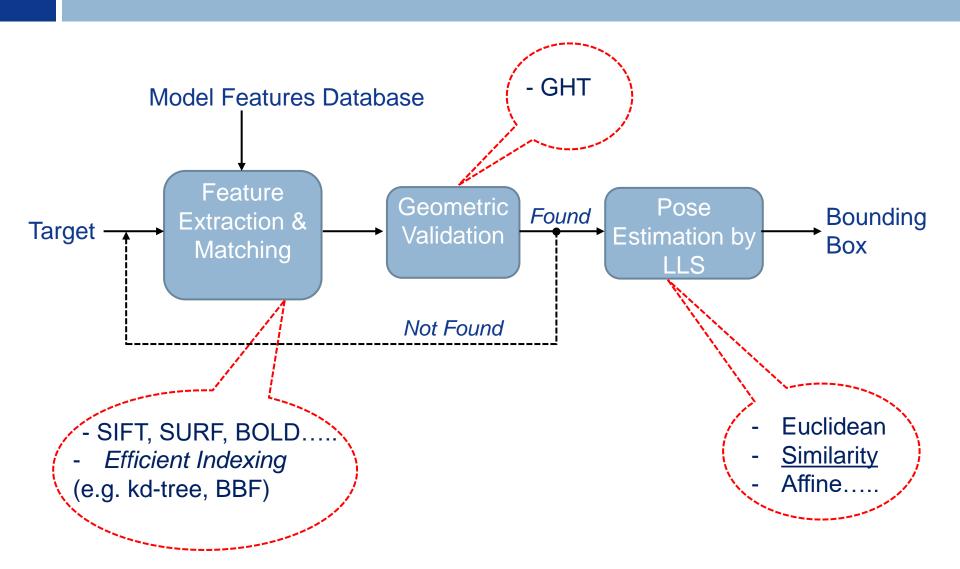
$$\forall r: y = x + r_s = x + s \cdot r, A[y, s] + +$$

#### Handling rotation and scale

$$\forall r: y = x + s \cdot r_{\theta}, A[y, \theta, s] + +$$

#### Detection by local invariant features





#### Geometric Validation: Star Model



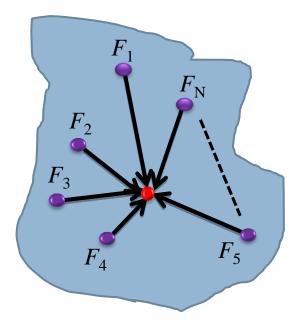
#### Off-line phase:

$$F = \{F_1, F_2 \dots F_N\}, F_i = (\mathbf{P}_i, \varphi_i, S_i, \mathbf{D}_i)$$

$$\mathbf{P}_{C} = \frac{1}{N} \sum_{i=1}^{N} \mathbf{P}_{i} \rightarrow \mathbf{V}_{i} = \mathbf{P}_{C} - \mathbf{P}_{i}$$







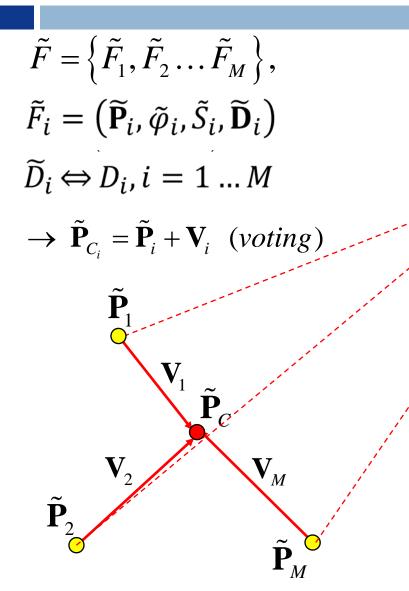
$$\forall F_i \in F$$

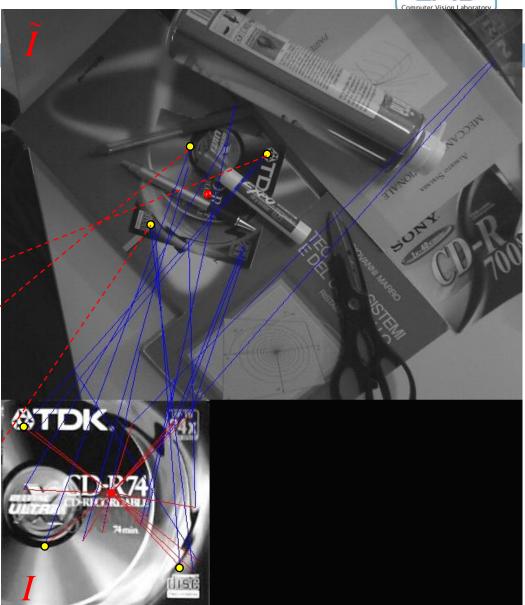
$$\downarrow$$

$$F_i = (\mathbf{P}_i, \varphi_i, S_i, \mathbf{D}_i, \mathbf{V}_i)$$

## Geometric Validation: On-line phase

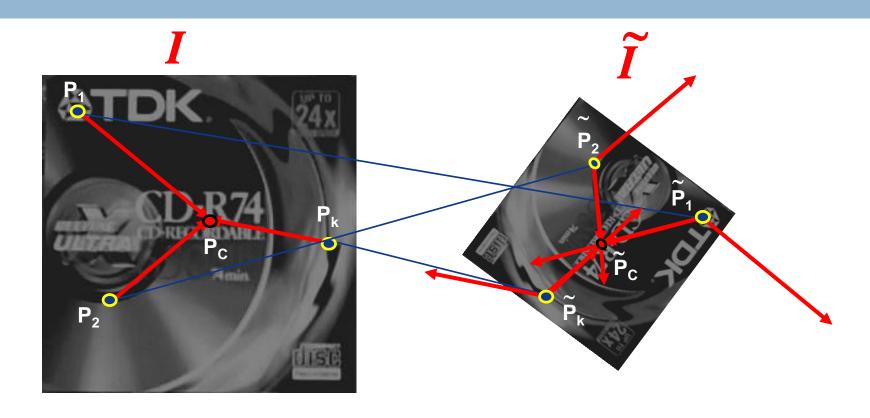






## Similarity Invariant Voting





$$\tilde{F}_{i} = (\tilde{\mathbf{P}}_{i}, \tilde{\varphi}_{i}, \tilde{S}_{i}, \tilde{\mathbf{D}}_{i}) \Leftrightarrow F_{i} = (\mathbf{P}_{i}, \varphi_{i}, S_{i}, \mathbf{D}_{i}, \mathbf{V}_{i})$$

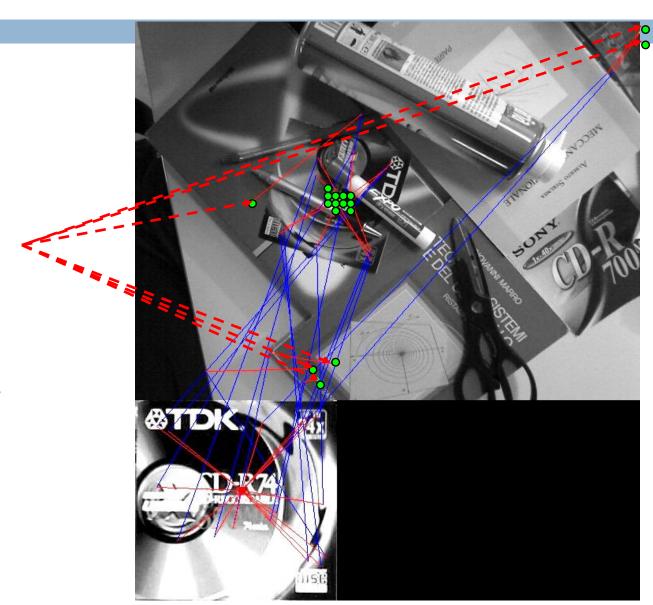
$$\Delta \varphi_{i} = \tilde{\varphi}_{i} - \varphi_{i}, \quad s_{i} = \frac{\tilde{S}_{i}}{S_{i}}$$

$$\tilde{\mathbf{P}}_{C_{i}} = \tilde{\mathbf{P}}_{i} + s_{i} \cdot \mathbf{R} (\Delta \varphi_{i}) \mathbf{V}_{i}$$

## Validating matches by the voting process

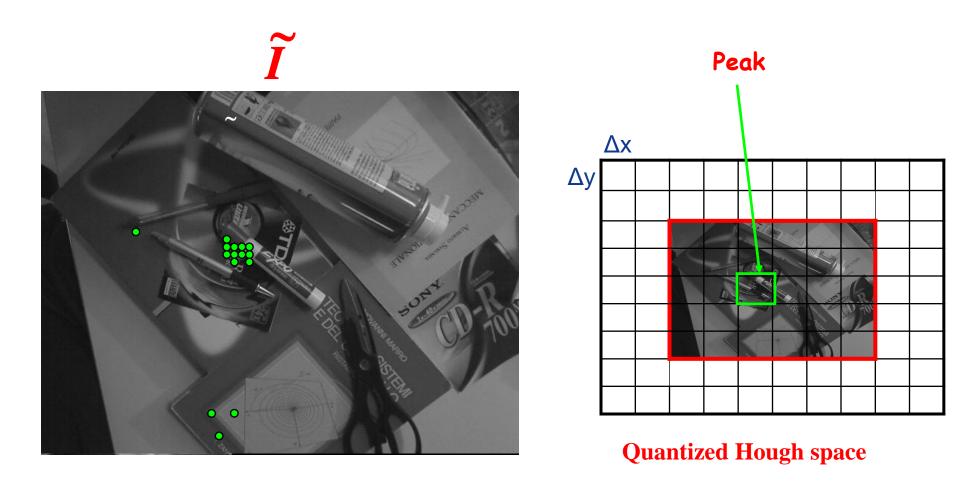


Similarly to the edge-based GHT, correct correspondences generate votes that accumulate coherently at the centre of the object in the target image. Conversely, wrong correspondences yield votes spread across the image, this providing no evidence for the presence and pose of the sought object.



# Quantization of the Hough (pose) space

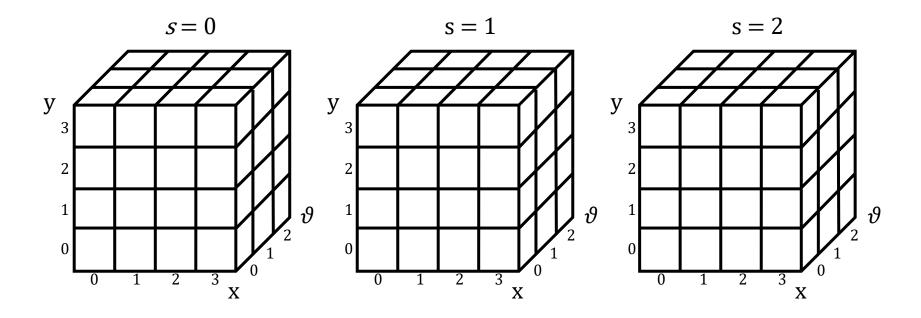




## Quantized Hough Space (1)



- 4D Hough Space
- Translation  $(x_c, y_c)$ , Rotation  $(\theta = \Delta \varphi)$ , Scale (s)

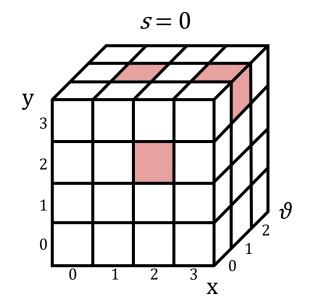


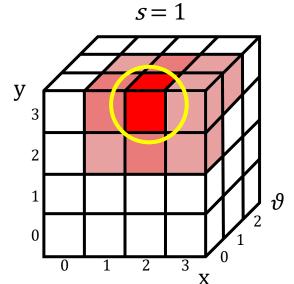
Exemplar 4D Hough Space (4 bins for  $x_C$  e  $y_C$ , 3 bins per  $\theta$  e s)

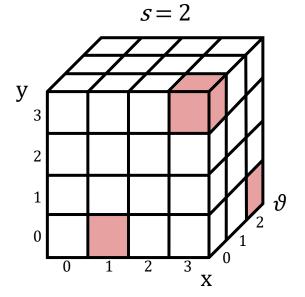
# Quantized Hough Space (2)



- Each correspondence provides a pose (translation, rotation, scale)
  hypothesis that is then accumulated into a certain bin within the
  quantized 4D Hough space.
- Coherent correspondences will cast votes into the very same bin.
- Thus, peaks in the 4D space highlight the most likely pose hypotheses concerning the sought object. Peaks are then thresholded to decide whether each such hypothesis should be accepted or rejected.

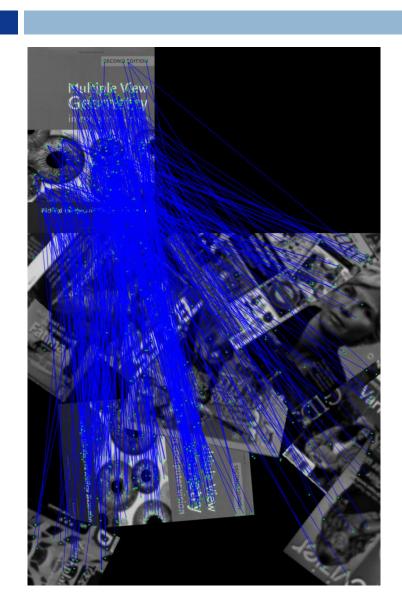


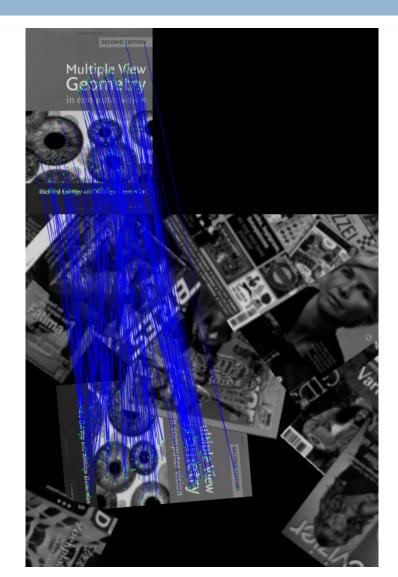




# **Exemplar Results**

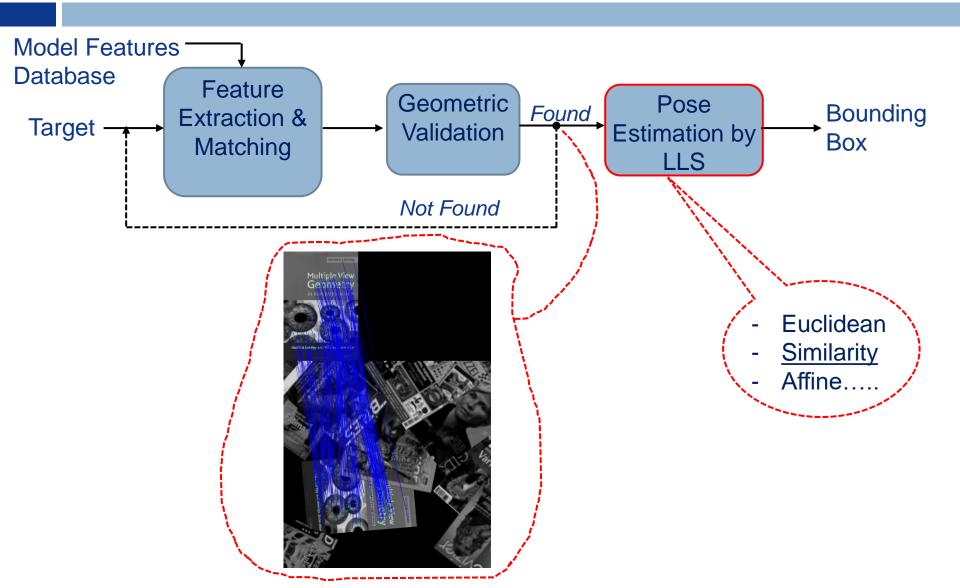






## Finally.....





# Similarity Estimation (1)



Given the positions of the feature pairs validated by the GHT:

$$\tilde{\mathbf{P}}_i \iff \mathbf{P}_i \ i = 1...n, \quad \tilde{\mathbf{P}}_i = \begin{pmatrix} \tilde{u}_i \\ \tilde{v}_i \end{pmatrix}, \quad \mathbf{P}_i = \begin{pmatrix} u_i \\ v_i \end{pmatrix}$$

We seek to estimate the similarity (pose) which brings the model into the target image:

$$\tilde{\mathbf{P}}_{i} = s\mathbf{R}\mathbf{P}_{i} + \mathbf{t} \quad \Rightarrow \quad \begin{pmatrix} \tilde{u}_{i} \\ \tilde{v}_{i} \end{pmatrix} = \begin{pmatrix} s\cos\theta & -s\sin\theta \\ s\sin\theta & s\cos\theta \end{pmatrix} \begin{pmatrix} u_{i} \\ v_{i} \end{pmatrix} + \begin{pmatrix} t_{u} \\ t_{v} \end{pmatrix}$$

$$m = s \cos \theta, \ n = s \sin \theta \implies \begin{pmatrix} \tilde{u}_i \\ \tilde{v}_i \end{pmatrix} = \begin{pmatrix} m & -n \\ n & m \end{pmatrix} \begin{pmatrix} u_i \\ v_i \end{pmatrix} + \begin{pmatrix} t_u \\ t_v \end{pmatrix}$$

# Similarity Estimation (2)



Let us rewrite the equation by grouping the unknown similarity parameters into a vector

$$\begin{pmatrix} u_i & -v_i & 1 & 0 \\ v_i & u_i & 0 & 1 \end{pmatrix} \begin{pmatrix} m \\ n \\ t_u \\ t_v \end{pmatrix} = \begin{pmatrix} \tilde{u}_i \\ \tilde{v}_i \end{pmatrix}$$

and then set-up an overdetermined (2*n* equations in 4 unknowns) linear system by considering all the available correspondences.

$$\begin{pmatrix}
u_1 & -v_1 & 1 & 0 \\
v_1 & u_1 & 0 & 1 \\
\vdots & \vdots & \vdots & \vdots \\
u_n & -v_n & 1 & 0 \\
v_n & u_n & 0 & 1
\end{pmatrix}
\begin{pmatrix}
m \\
n \\
t_u \\
t_v
\end{pmatrix}
=
\begin{pmatrix}
\tilde{u}_1 \\
\tilde{v}_1 \\
\vdots \\
\tilde{u}_n \\
\tilde{v}_n
\end{pmatrix}
\Rightarrow \mathbf{A}\mathbf{x} = \mathbf{b}$$

# Similarity Estimation (3)



The previous system admits a least-square solution, which can be computed either by the *normal equations* or the *pseudo-inverse*:

$$\mathbf{x} = \left(\mathbf{A}^T \mathbf{A}\right)^{-1} \mathbf{A}^T \mathbf{b} = \mathbf{A}^+ \mathbf{b}$$

• Once (m,n) are found, rotation (9) and scale (s) may be easily disentangled:

$$\widehat{\mathcal{G}} = \tan^{-1}\left(\frac{n}{m}\right), \ \widehat{\mathcal{G}} \in \left[-\frac{\pi}{2}, \frac{\pi}{2}\right]$$

$$\widehat{\mathcal{G}} \in \left[-\frac{\pi}{2}, \frac{\pi}{2}\right] \Rightarrow \mathcal{G} \in [0, 2\pi] :$$

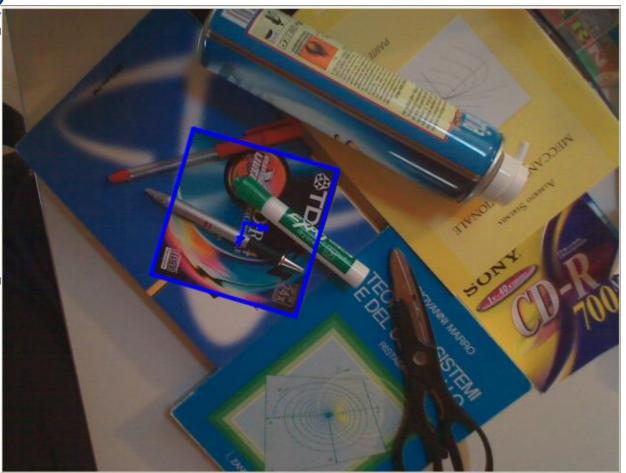
$$\widehat{\mathcal{G}} \in \left[0, \frac{\pi}{2}\right] \Rightarrow \mathcal{G} = \widehat{\mathcal{G}} + sign(m) \cdot \pi, \quad \widehat{\mathcal{G}} \in \left[-\frac{\pi}{2}, 0\right] \Rightarrow \mathcal{G} = \widehat{\mathcal{G}} + 2\pi - sign(m) \cdot \pi$$
with  $sign(m) = \begin{cases} 1, & \text{if } m < 0 \\ 0, & \text{otherwise} \end{cases}$ 

$$s = \frac{m}{\cos \theta} = \frac{n}{\sin \theta}$$

# Exemplar Results (1)

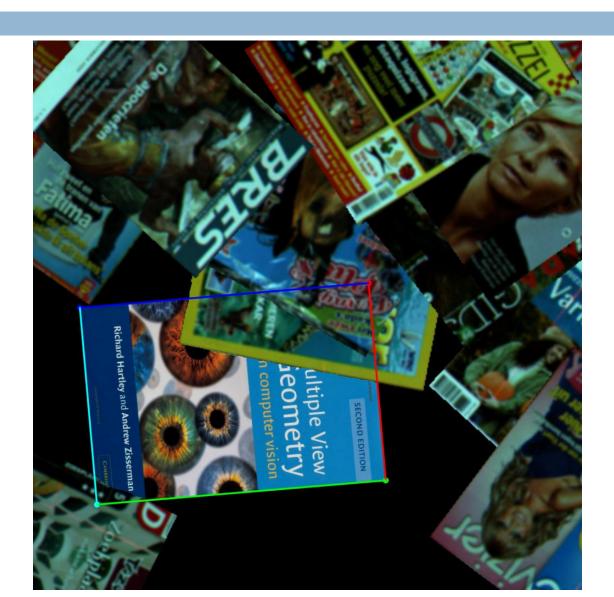






# Exemplar Results (2)





# Exemplar Results (3)















# Exemplar Results (BOLD)

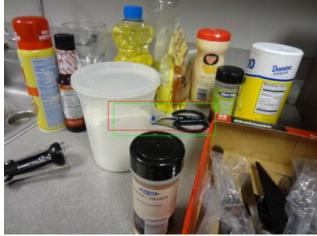














https://www.youtube.com/watch?v=lqVBsyOVMro







# That's all folks!



#### Main References



- 1) Martin, J. and Crowley, J. (1995). *Experimental comparison of correlation techniques*. In Proc. Int. Conf. On Intelligent Autonomous Systems, volume 4, pages 86–93.
- 2) W. Li and E. Salari, *Successive elimination algorithm for motion estimation*, IEEE Trans. on Image Processing, vol. 4, no. 1, pp. 105–107, 1995
- 3) F. Tombari, S. Mattoccia, L. Di Stefano, "Full search-equivalent pattern matching with Incremental Dissimilarity Approximations", IEEE Trans. on Pattern Analysis and Machine Intelligence, Vol. 31, N.1, pp. 129-141, Jan. 2009
- 4) Steger, C, Occlusion, clutter, and illumination invariant object recognition. In Kalliany, R., Leberl, F., eds.: International Archives of Photogrammetry, Remote Sensing, and Spatial Information Sciences. Volume XXXIV, part 3A., Graz (2002).
- 5) A. Hofhauser, C. Steger, N. Navab, *Edge-based Template Matching and Tracking for Perspectively Distorted Planar Objects*, Proc. 4th International Symposium on Advances in Visual Computing ISVC '08, 2008.
- 6) Borgefors, G., 1988. *Hierarchical chamfer matching: A parametric edge matching algorithm.* IEEE Trans on Pattern Analysis and Machine Intelligence 10(6), pp. 849–865.
- 7) Steger, C., 1998. *An unbiased detector of curvilinear structures*. IEEE Trans. on Pattern Analysis and Machine Intelligence 20(2), pp. 113–125.