

# VISUAL OBJECT RECOGNITION

STATE-OF-THE-ART  
TECHNIQUES AND  
PERFORMANCE EVALUATION

# LECTURE 4: DESCRIPTOR CONSTRUCTION

- ✱ An opportunity for machine learning
- ✱ Scale Invariant Feature Transform (SIFT)
- ✱ Speeded Up Robust Features (SURF)
- ✱ Geometric Blur and Log-polar grids
- ✱ ContourSIFT
- ✱ Results from learning descriptors

# OPPORTUNITY FOR MACHINE LEARNING

- ✻ When we have a system that automatically detects canonical frames, we could potentially use machine learning to find a good way to describe the invariant image patches.



# OPPORTUNITY FOR MACHINE LEARNING

- ✱ Any method used to learn global appearance (LE1) could potentially be used.
- ✱ But, high-dimensional learning requires a large amount of training data.
- ✱ Solution: Parametrise space of solutions in intelligent ways.
- ✱ See: Winder&Brown, *Learning Local Image Descriptors*, CVPR'07 (more at end of lecture)

# THE SIFT DESCRIPTOR

- ✱ Scale Invariant Feature Transform
- ✱ Converts an image patch sampled in canonical frame to a 128-byte *descriptor vector*.
- ✱ Inherits geometric invariances from c-frame.

# THE SIFT DESCRIPTOR

- ✱ Compute x- and y-gradients through convolution:

$$\nabla \mathbf{f}(\mathbf{x}) = \begin{bmatrix} (d_x * f)(\mathbf{x}) \\ (d_y * f)(\mathbf{x}) \end{bmatrix}$$

- ✱ Rotate gradient map to direction from orhist:

$$\nabla \hat{\mathbf{f}}(\mathbf{x}) = \mathbf{R} \nabla \mathbf{f}(\mathbf{R}^T \mathbf{x})$$

- ✱ Compute gradient orientation histograms in 4x4 spatial regions:

$$h_{kl} = \sum_{\mathbf{x} \in \text{patch}_l} |\nabla \hat{\mathbf{f}}(\mathbf{x})| w(\mathbf{x} + \mathbf{d}_l) B_k(\tan^{-1} \nabla \hat{\mathbf{f}}(\mathbf{x}))$$

# THE SIFT DESCRIPTOR

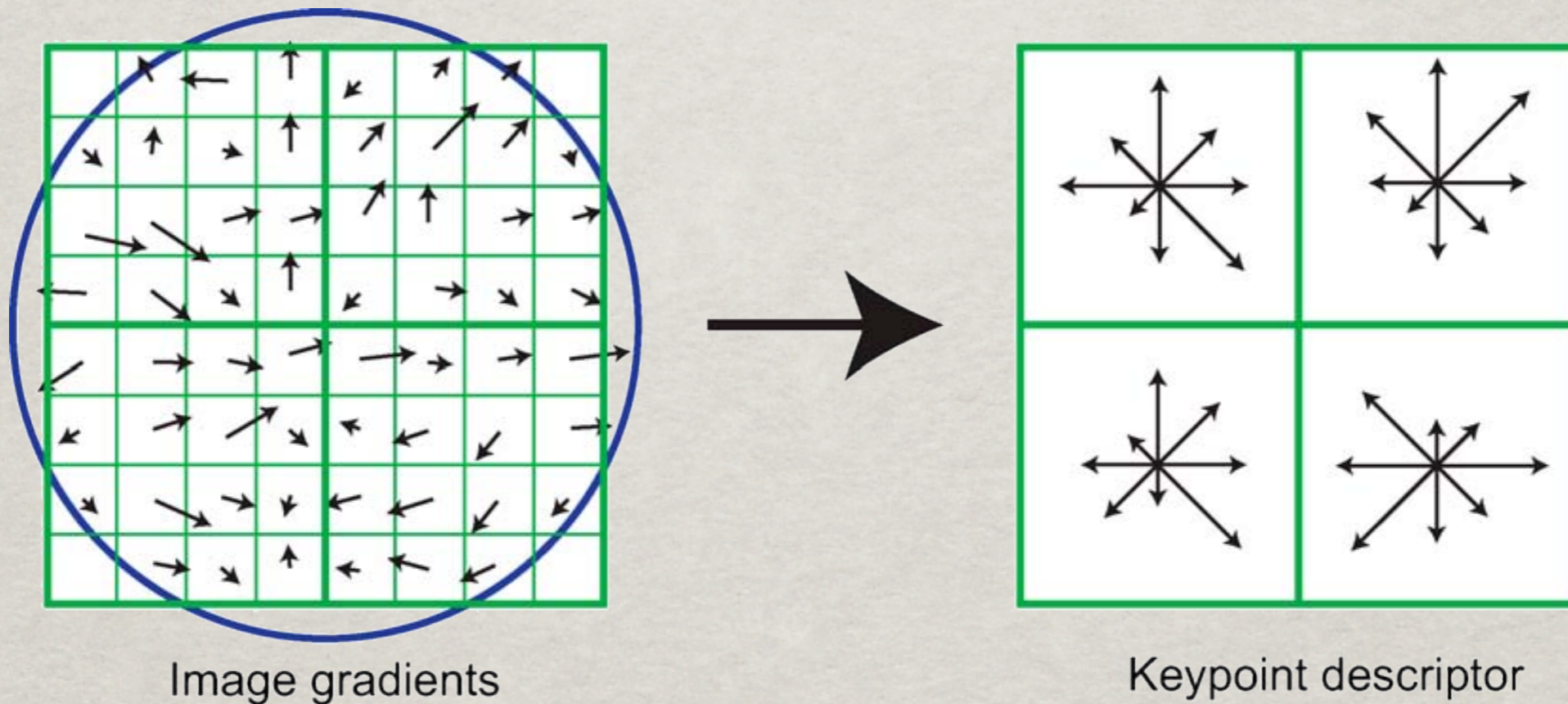
- ✱ Compute gradient orientation histograms in  $4 \times 4$  spatial regions :

$$h_{kl} = \sum_{\mathbf{x} \in \text{patch}_l} |\nabla \hat{\mathbf{f}}(\mathbf{x})| w(\mathbf{x} + \mathbf{d}_l) B_k(\tan^{-1} \nabla \hat{\mathbf{f}}(\mathbf{x}))$$

- ✱  $B_k(\mathbf{x})$  linear interpolation kernel  
Quadratic is better (Jonsson&Felsberg)
- ✱ Subwindows  $l \in [1 \dots 16]$  directions  $k \in [1 \dots 8]$
- ✱ Spatial weight  $w(\mathbf{x} + \mathbf{d}_l)$  (Gaussian decay)

# THE SIFT DESCRIPTOR

- ✱ Implementation with source code (A. Vedaldi):  
<http://vision.ucla.edu/~vedaldi/code/sift/sift.html>



Note that 4x4 regions are actually used, with 8 orientations, 128 elements

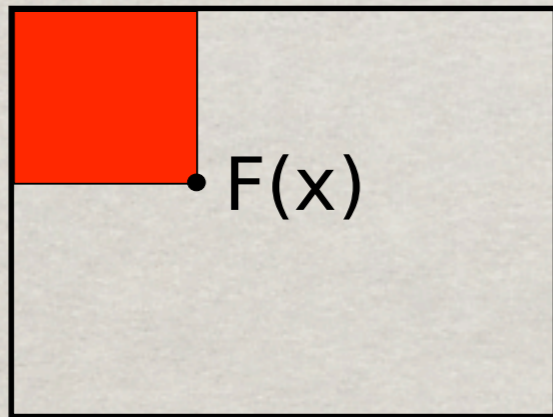


# THE SIFT DESCRIPTOR

- ✱ Affine illumination invariance by using gradients and normalising descriptor  $\hat{\mathbf{h}} = \mathbf{h}/\|\mathbf{h}\|$
- ✱ Some robustness by truncating and normalising again  $\hat{\mathbf{h}} = \min(t, \hat{\mathbf{h}})/\|\hat{\mathbf{h}}\|$
- ✱ The spatial histogramming gives robustness to scale/rotation/translation errors.
- ✱ Bears some similarity to the “Standard model” (see LE1)

# SURF

- ✱ Bay & Tuytelaars & van Gool, *SURF: Speeded Up Robust Features*, ECCV06
- ✱ Both detector and descriptor optimised for speed using *Integral images*. (Viola & Jones)



$$F(\mathbf{x}) = \sum_{k=0}^{x_1} \sum_{l=0}^{x_2} f(k, l)$$

- ✱ Very fast to compute:

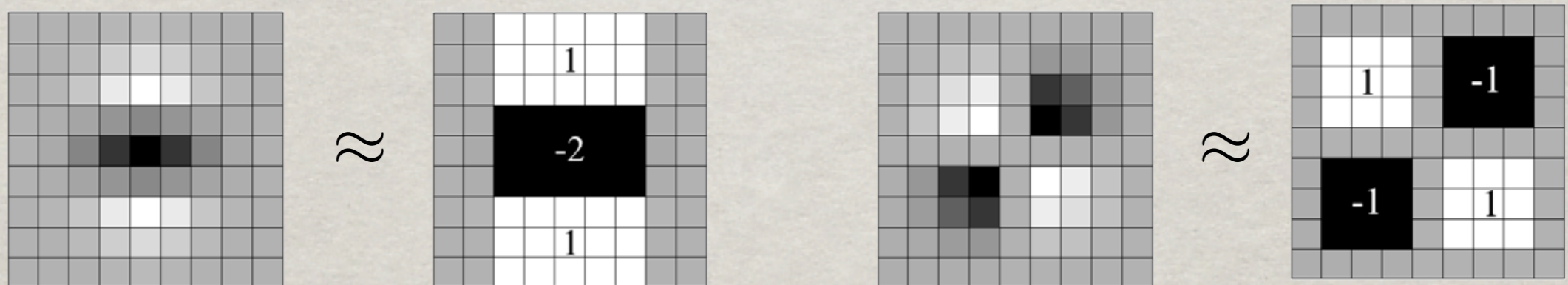
$$F(\mathbf{x}) = F(\mathbf{x} - \begin{pmatrix} 1 \\ 0 \end{pmatrix}) + F(\mathbf{x} - \begin{pmatrix} 0 \\ 1 \end{pmatrix}) - F(\mathbf{x} - \begin{pmatrix} 1 \\ 1 \end{pmatrix}) + f(\mathbf{x})$$

# THE SURF DETECTOR

- Find interest points using Hessian matrix in scale space:

$$\mathbf{H}(\mathbf{x}, \sigma) = \begin{bmatrix} f_{xx}(\mathbf{x}, \sigma) & f_{xy}(\mathbf{x}, \sigma) \\ f_{xy}(\mathbf{x}, \sigma) & f_{yy}(\mathbf{x}, \sigma) \end{bmatrix}$$

- Approximate filters using I-images:



- Scale space by increasing filter size.

# THE SURF DETECTOR

- ✱ Scale space maxima of  $\det(\mathbf{H})$

$$\det(\mathbf{H}) \approx H_{11}H_{22} - 0.9H_{12}^2$$

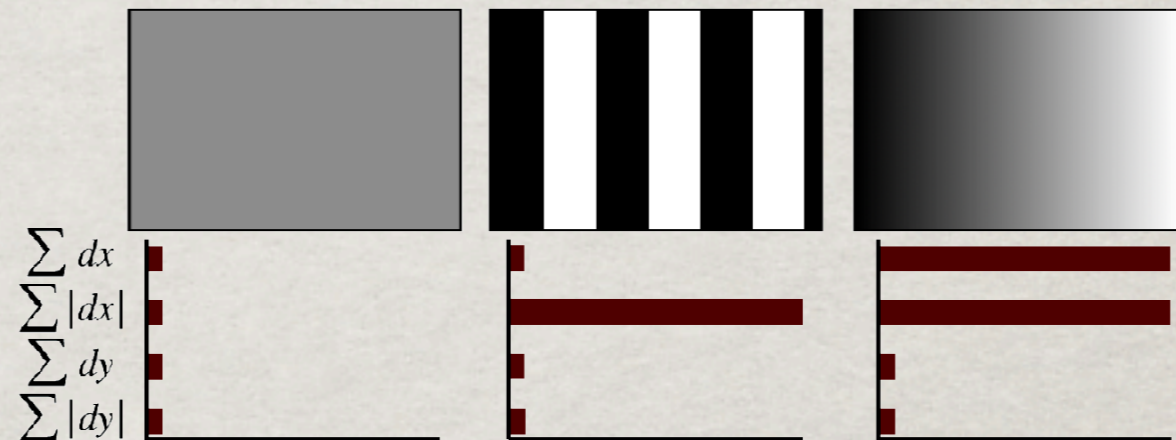
- ✱ Detection like in SIFT:

A.  $3 \times 3 \times 3$  non-max-suppression

B. Quadratic polynomial interpolation in scale space (Brown&Lowe'02)

# THE SURF DESCRIPTOR

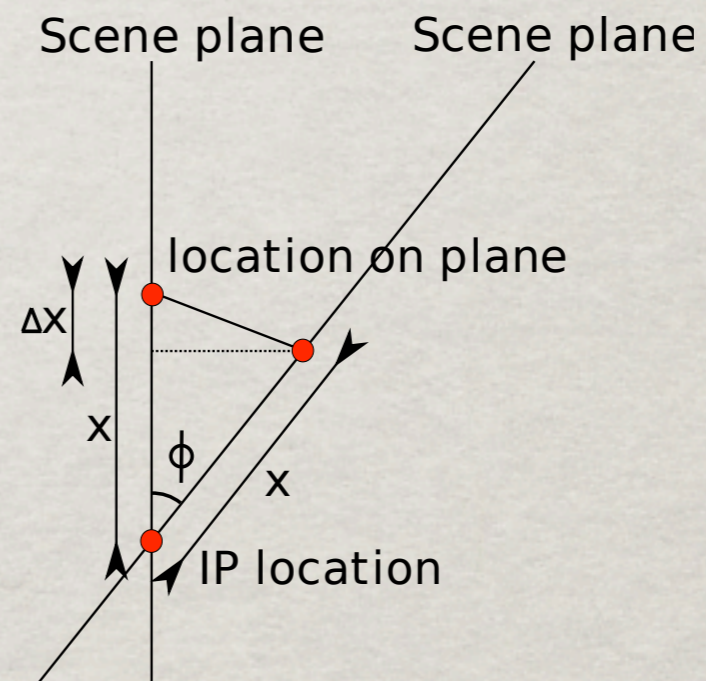
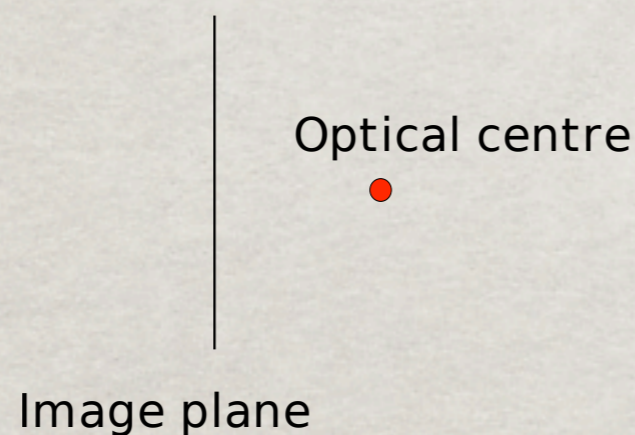
- ✱ Again using I-images



- ✱ 4 responses in each 4x4 block  $\Rightarrow$  64 elements in descriptor.
- ✱ Descriptor is normalised.
- ✱ Sign of Laplacian tells dark blobs from bright cf. MSER<sub>+</sub> and MSER<sub>-</sub>

# GEOMETRIC BLUR

- ✱ Introduced by Berg&Malik CVPR'01. As a descriptor at CVPR'05.

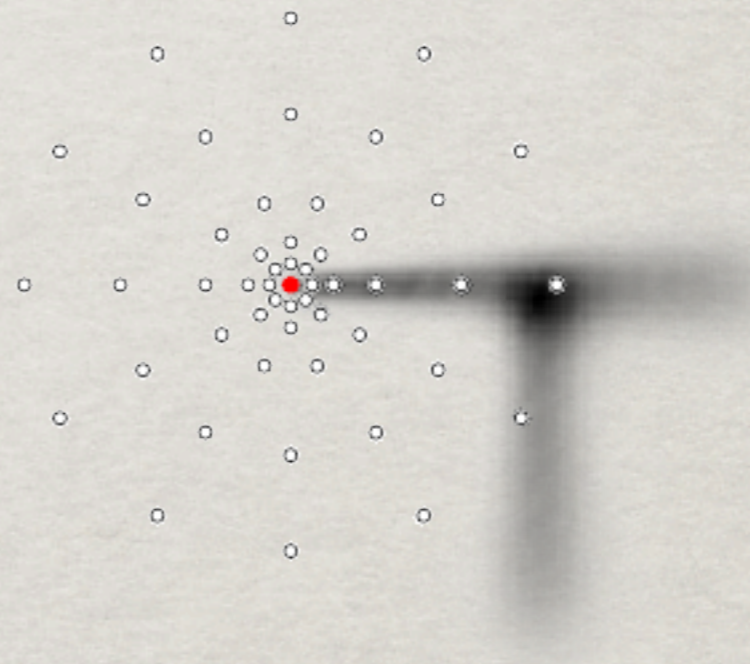
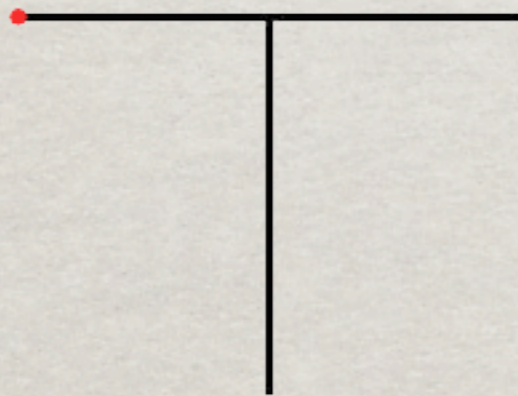


- ✱ If orthogonal projection: Point on a planar scene is displaced proportionally to the distance from the reference point:

$$\Delta x = (1 - \cos \phi)x$$

# GEOMETRIC BLUR

- ☼ Increase blur linearly with distance from interest point.



- ☼ Sample in 38 sparse points

# GEOMETRIC BLUR DESCRIPTOR

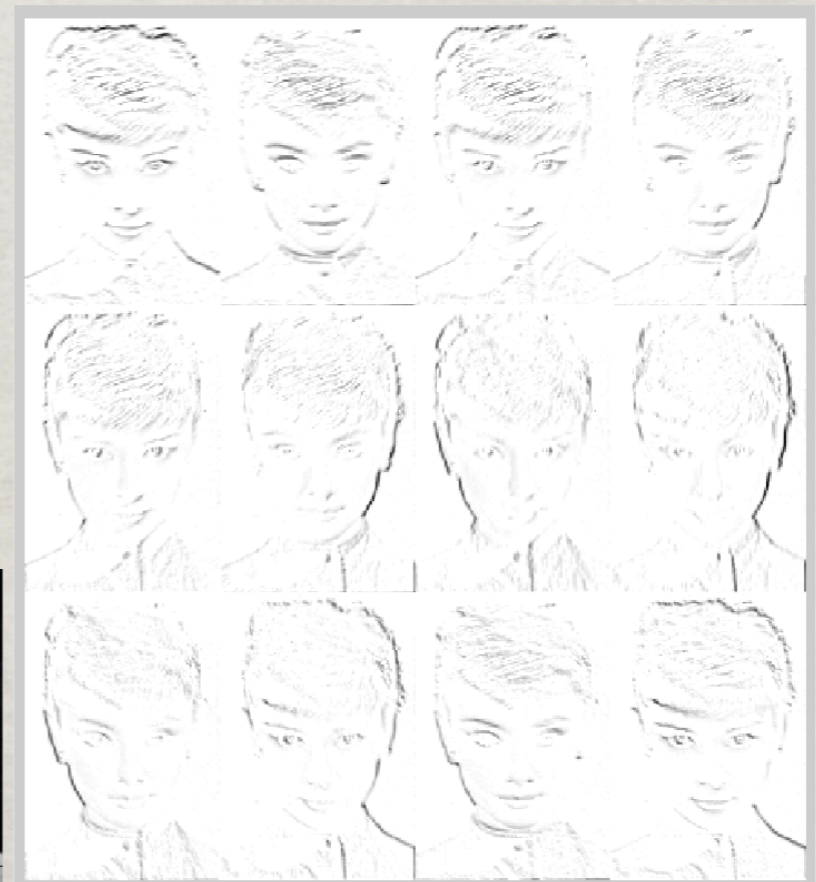
- Use *half-wave rectified* gradient responses in six directions:

$$h_k^+(\mathbf{x}) = \max(0, (f * g_k)(\mathbf{x}))$$

$$h_k^-(\mathbf{x}) = \max(0, -(f * g_k)(\mathbf{x}))$$

- Each is blurred with geometric blur before sampling.

- 228 elements

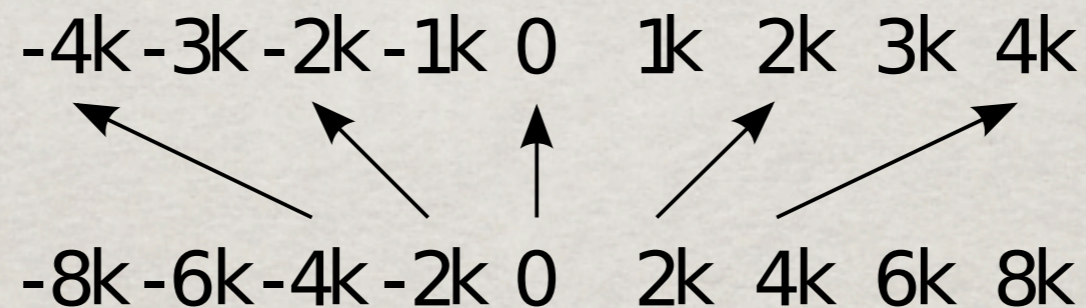




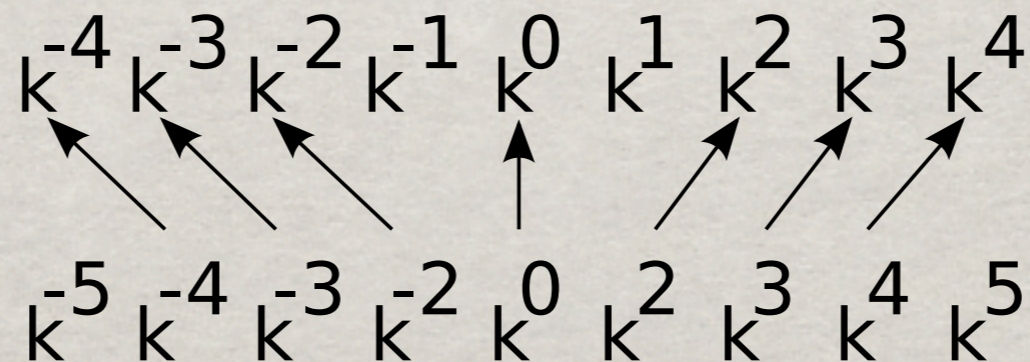
# LOG-POLAR PATCHES

☀ A similar sampling pattern can also be motivated from scale invariance.

☀ With linear distances, matching across scale is difficult:



☀ With exponential distances, scaling becomes a shift:



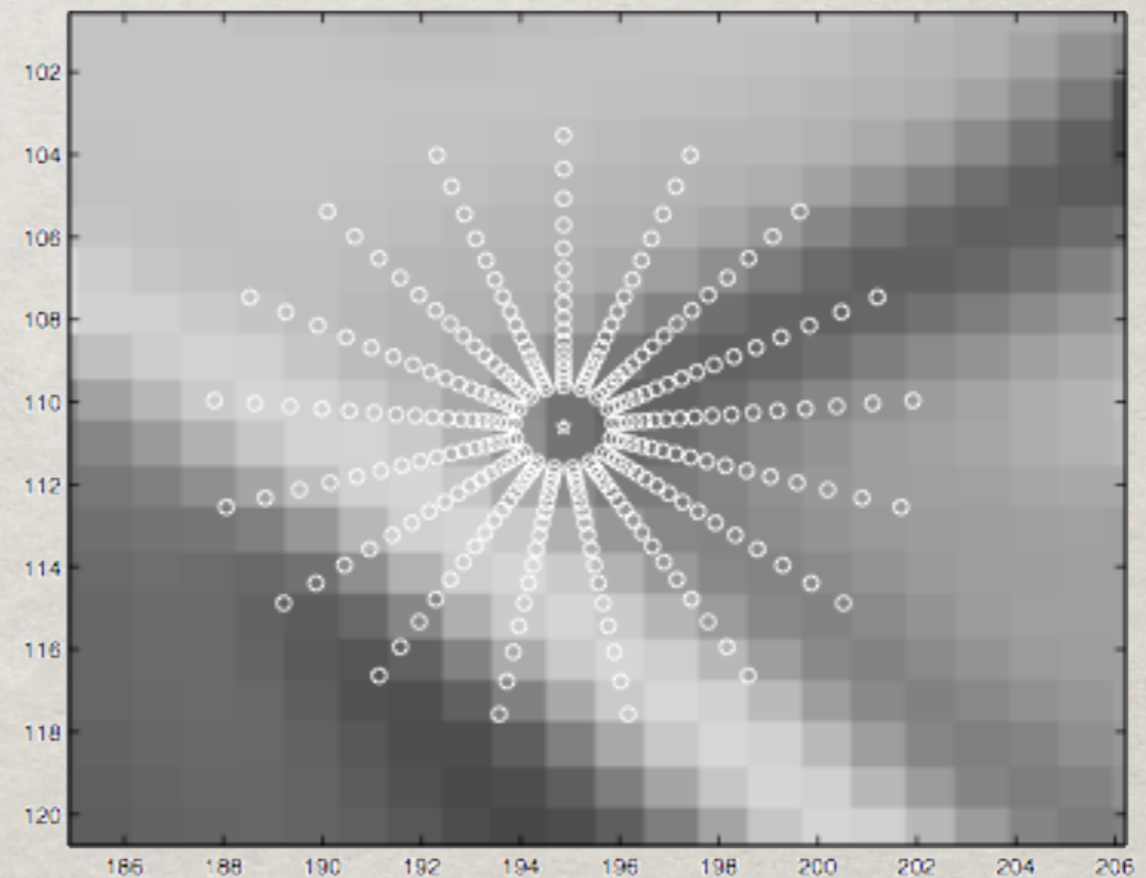
# LOG-POLAR PATCHES

✻ Viksten & Moe, *Local single-patch features for pose estimation, using the log-polar transform.*

IbPRIA'05



(a) Harris points.

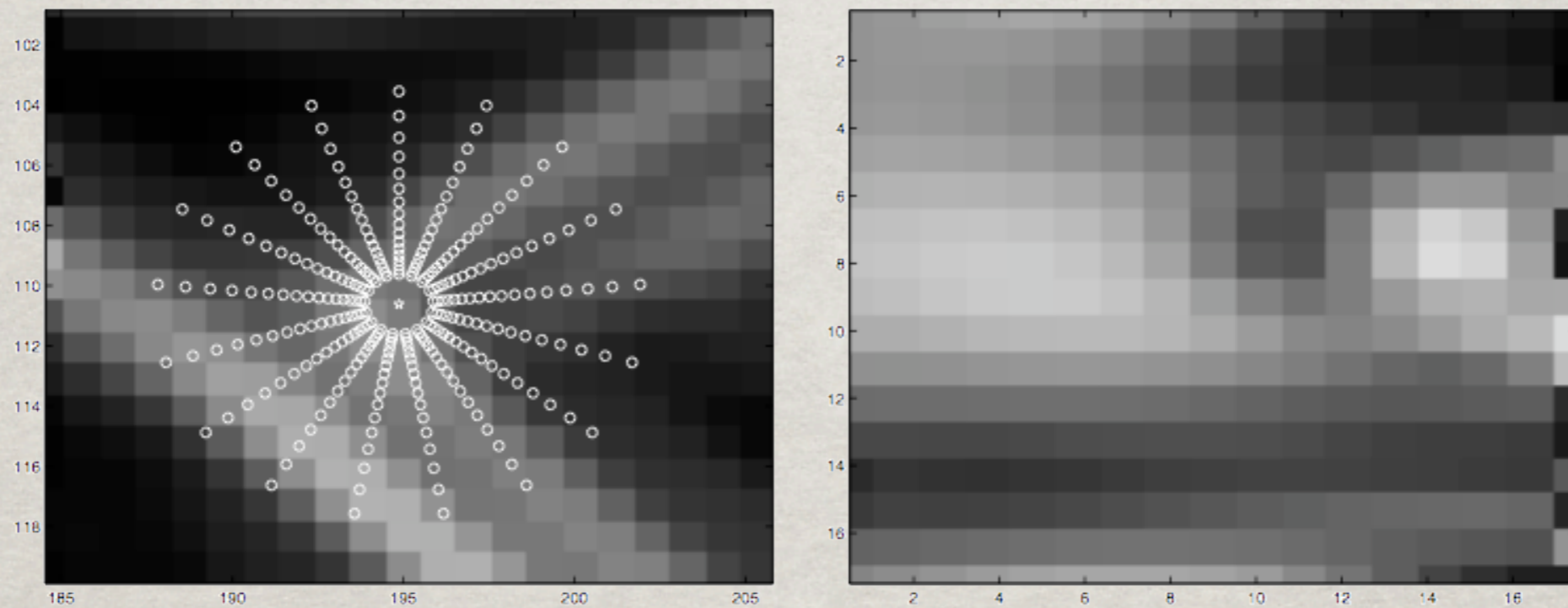


(b) Sample grid around Harris point.

# LOG-POLAR PATCHES

- ☼ Sample gradient image around Harris points in a log-polar pattern.

Figure 6.1: Harris points and log-polar sampling grid.



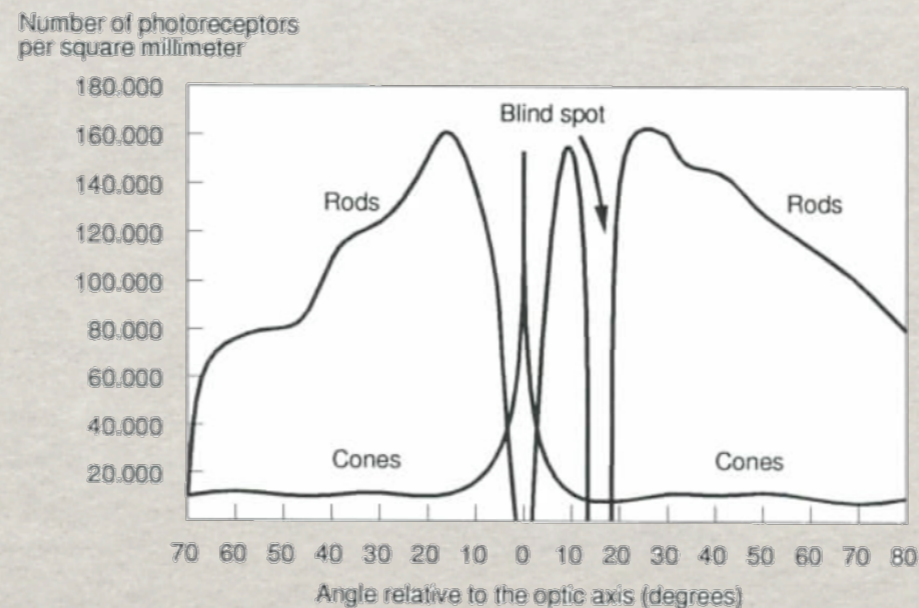
(a) Edge image with sample grid.

(b) Log-polar sampled edge image.

- ☼ Scale change can now be found using correlation.

# LOG-POLAR PATCHES

- ☼ Same sampling pattern is used in GLOH descriptor. Mikolajczyk&Schmid “A performance evaluation of local descriptors”, TPAMI 2005
- ☼ Fun fact: The eye also has a (partially) log polar sampling pattern.



# DIFFICULT CASES FOR DESCRIPTORS

- ✻ Background clutter in 3D scenes



- ✻ Patches cut out around features will have varying background.

# DIFFICULT CASES FOR DESCRIPTORS

- ☼ Large illumination changes



- ☼ Gradient strength changes non-uniformly.
- ☼ Contrast may be inverted.

# CONTOUR SIFT

- ✱ Idea: Use a detector that produces contours, e.g. MSER or MSCR



Input image



64 random MSER- regions

- ✱ Region shape is robust to changes outside the region

# CONTOUR SIFT

- ✱ Compute a descriptor from the binary mask of the region instead of the grey-scale patch.



- ✱ Less descriptive patches, but more robust to illumination and background clutter

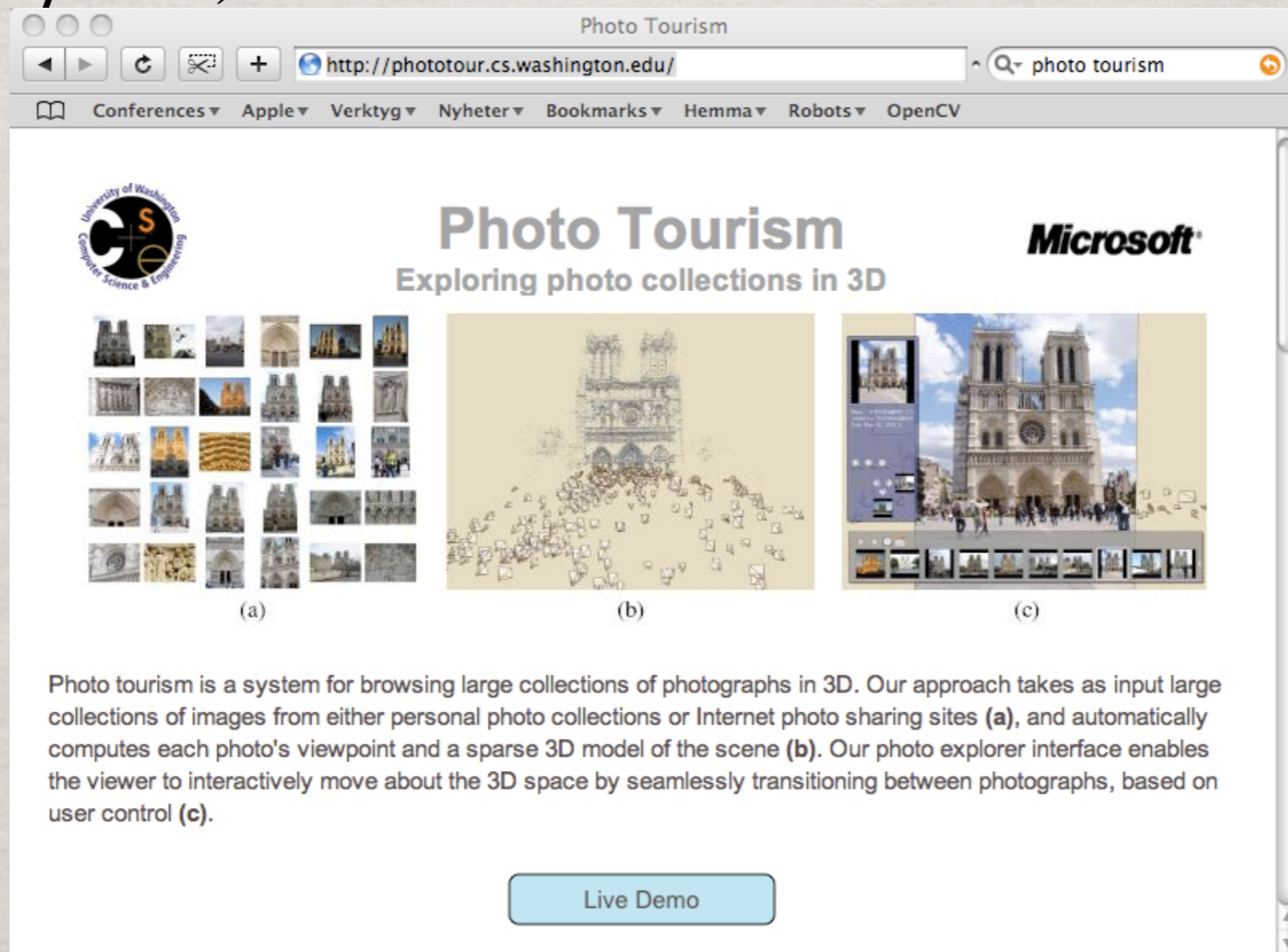


# CONTOUR SIFT

- ✱ Shape Descriptors for Maximally Stable Extremal Regions, Forssén&Lowe, ICCV'07
- ✱ Use the “standard SIFT pipeline”
- ✱ Re-tune all parameters to maximise performance on binary patches.
- ✱ Use Mikolajczyk's data set for training.

# LEARNING LOCAL DESCRIPTORS

✻ Winder Brown, *Learning Local Image Descriptors*, CVPR07



✻ Training data from Photo-tourism dataset.

# LEARNING LOCAL DESCRIPTORS

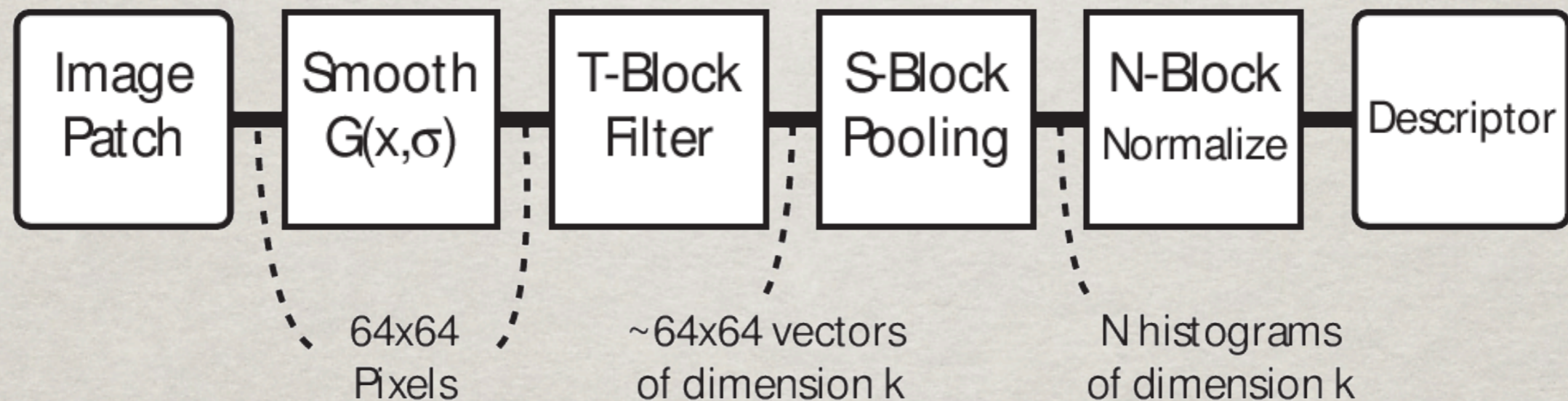
- ✻ Training data from Photo-tourism dataset.



- ✻ A research predecessor of MS Photosynth  
<http://livelabs.com/photosynth/>
- ✻ A neat way to arrange photo collections in 3D

# LEARNING LOCAL DESCRIPTORS

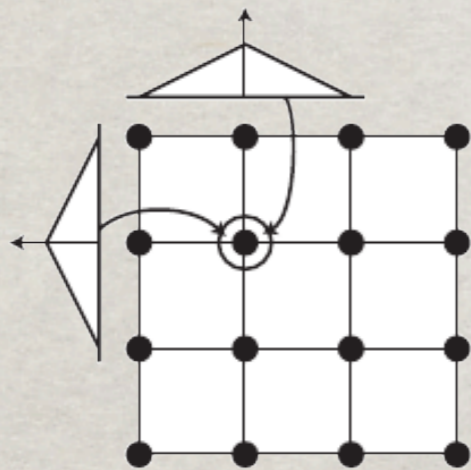
- ✱ Most descriptor algorithms can be described by the pipeline:



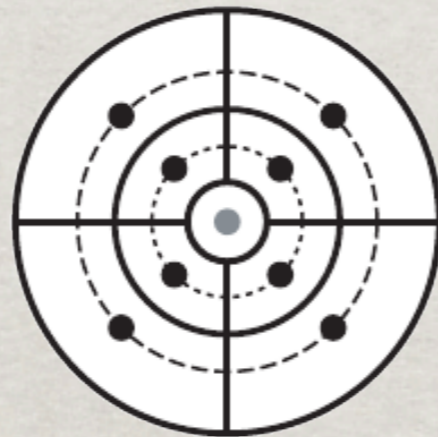
- ✱ Tune all steps of the algorithm with a robust Newton method (Powell's direction set method).

# LEARNING LOCAL DESCRIPTORS

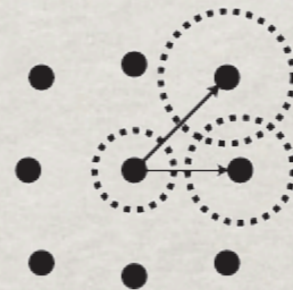
- Each block has many discrete variants, as well as continuous parameters.



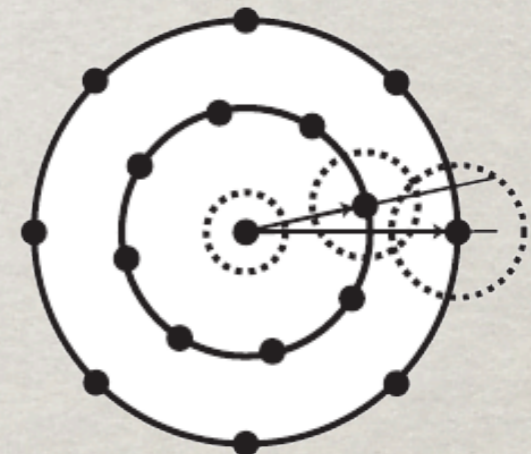
S1: SIFT grid with bilinear weights



S2: GLOH polar grid with bilinear radial and angular weights



S3: 3x3 grid with Gaussian weights

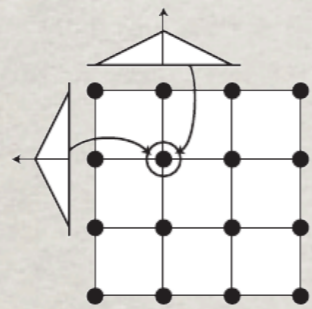


S4: 17 polar samples with Gaussian weights

- E.g. S-block (Spatial pooling).

# LEARNING LOCAL DESCRIPTORS

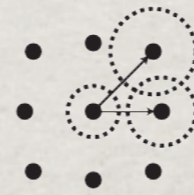
☼ Results for S-block:



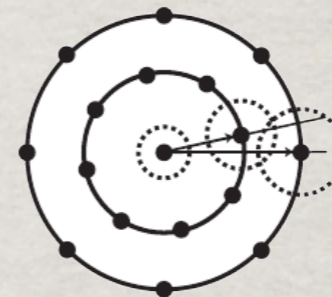
S1: SIFT grid with bilinear weights



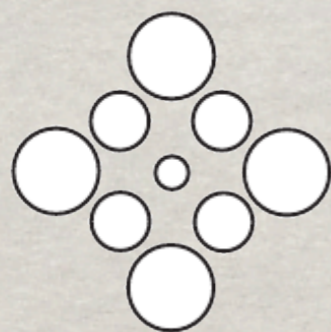
S2: GLOH polar grid with bilinear radial and angular weights



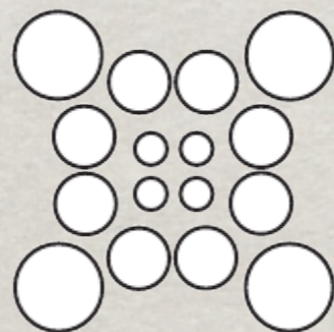
S3: 3x3 grid with Gaussian weights



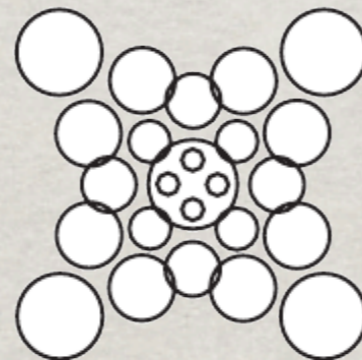
S4: 17 polar samples with Gaussian weights



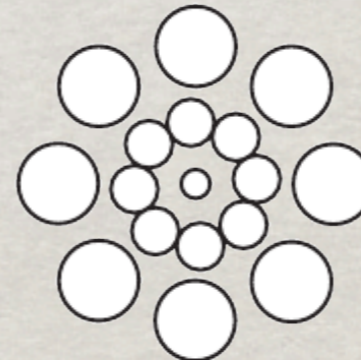
S3-9 (3x3)



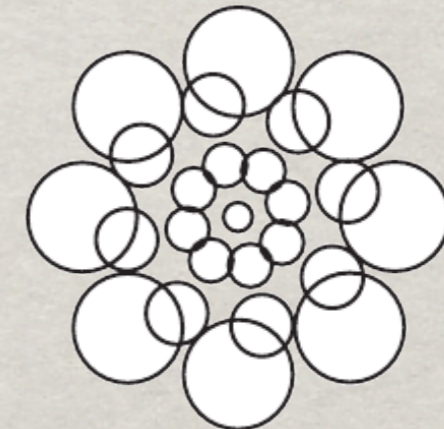
S3-16 (4x4)



S3-25 (5x5)



S4-17



S4-25

☼ Radially increasing blur seems useful.

# LEARNING LOCAL DESCRIPTORS

- ✱ Experiments used DoG detector.
- ✱ Winner was a 4th order quadrature descriptor with 400 elements
- ✱ Winner had an error rate of 2%, where the standard SIFT had 6%

# DISCUSSION

✻ Questions/comments on paper and lecture.