# COSC428 Computer Vision



# Local Features

- Interest operators
- Correspondence
- Invariances
- Descriptors

# Local Features

Matching points across images important for recognition and pose estimation

Tracking vs. Indexing

# Today

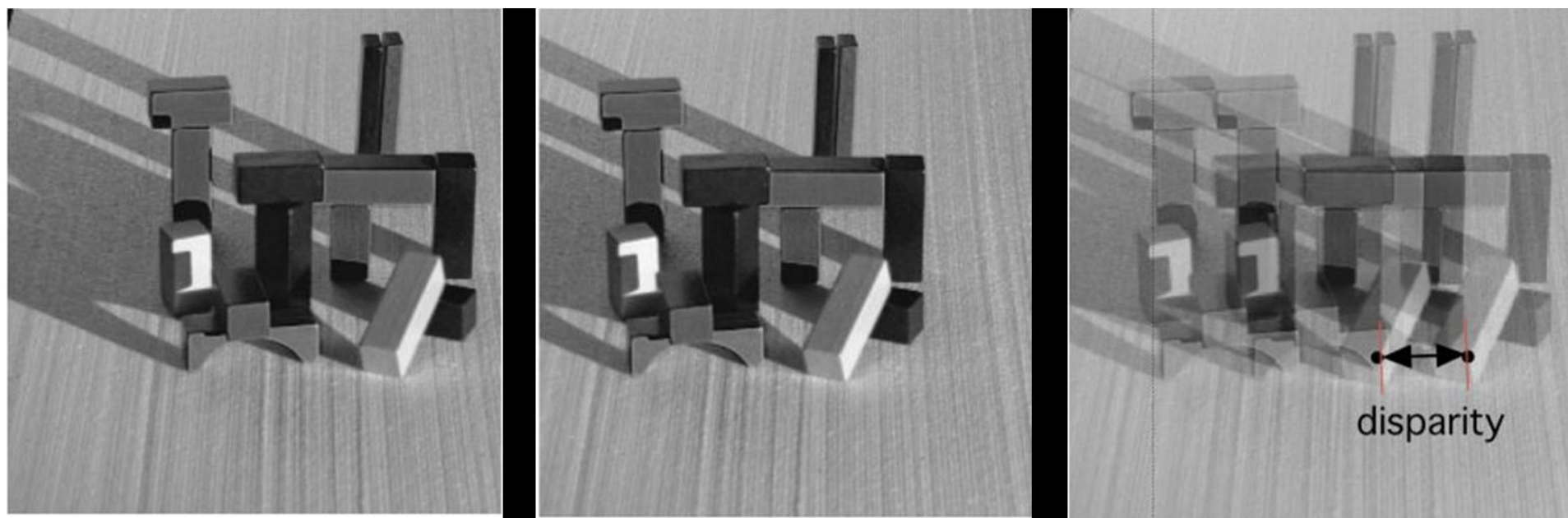Interesting points, correspondence, affine patch
  tracking

Scale and rotation invariant descriptors

# Correspondence using window matching

Points are highly individually ambiguous…
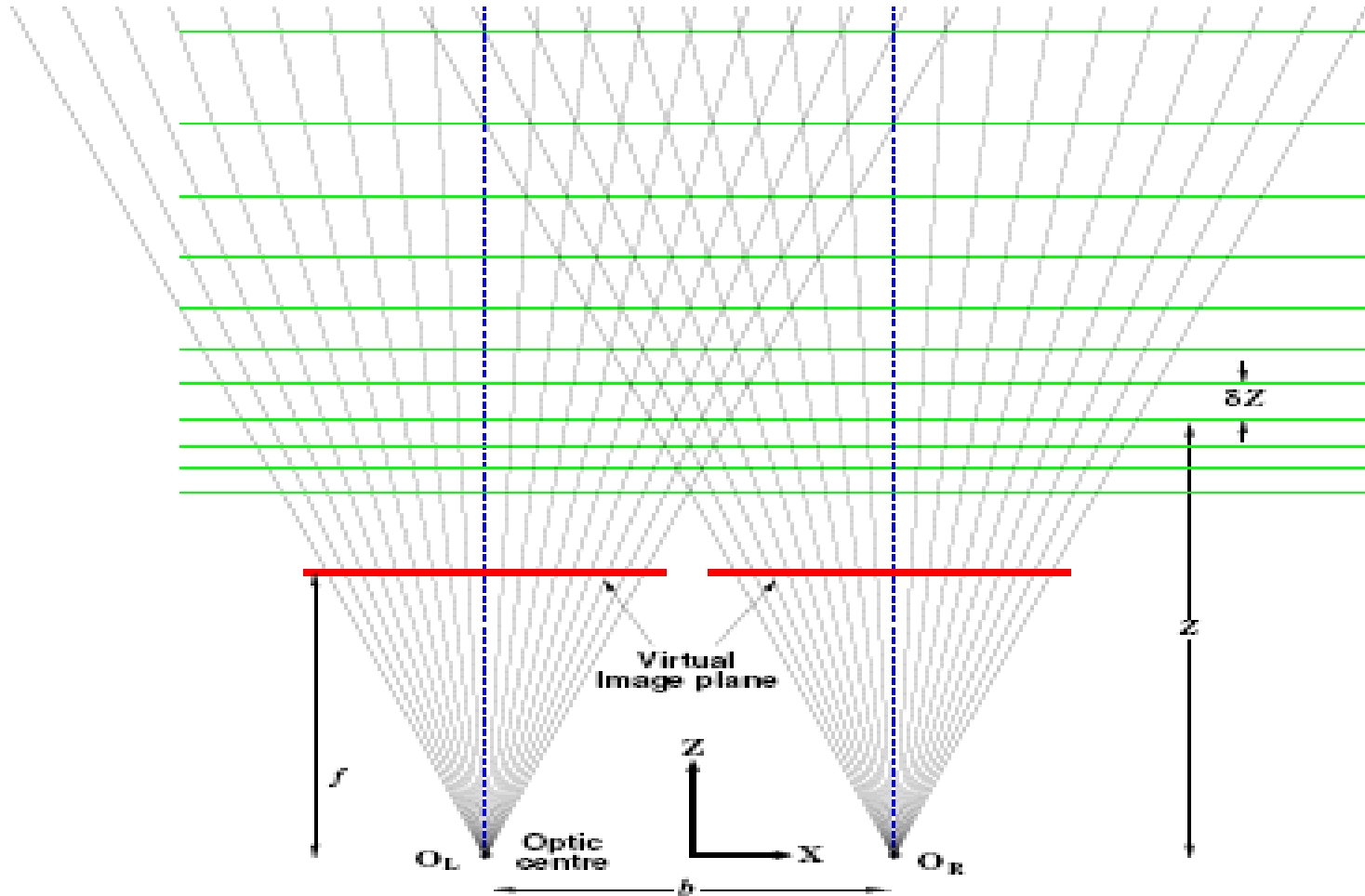
More unique matches are possible with small regions of image.

# Typical Stereo Camera Setup



disparity

- Left:       image from left camera (reference image)
- Middle:    image from right camera (match image)
- Right:      overlapping reference and match images

- Disparity: difference in pixel locations in reference and match images
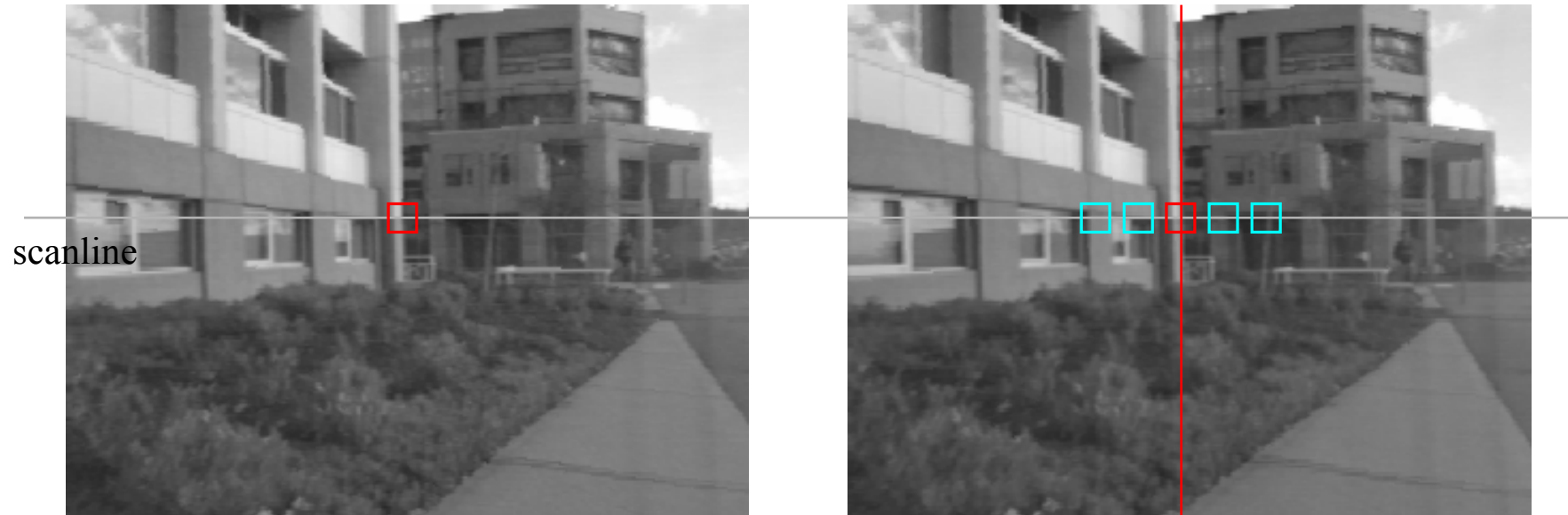
# Typical Stereo Camera Setup



- **Cameras are aligned so images are co-planer with (epipolar) lines passing through pixels on the same row in both left and right camera images (*rectified images*).**

- **Decreasing depth resolution/accuracy (larger $\delta Z$) further from camera.**

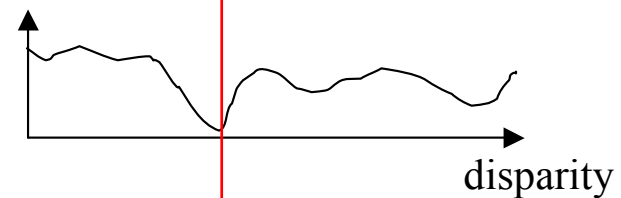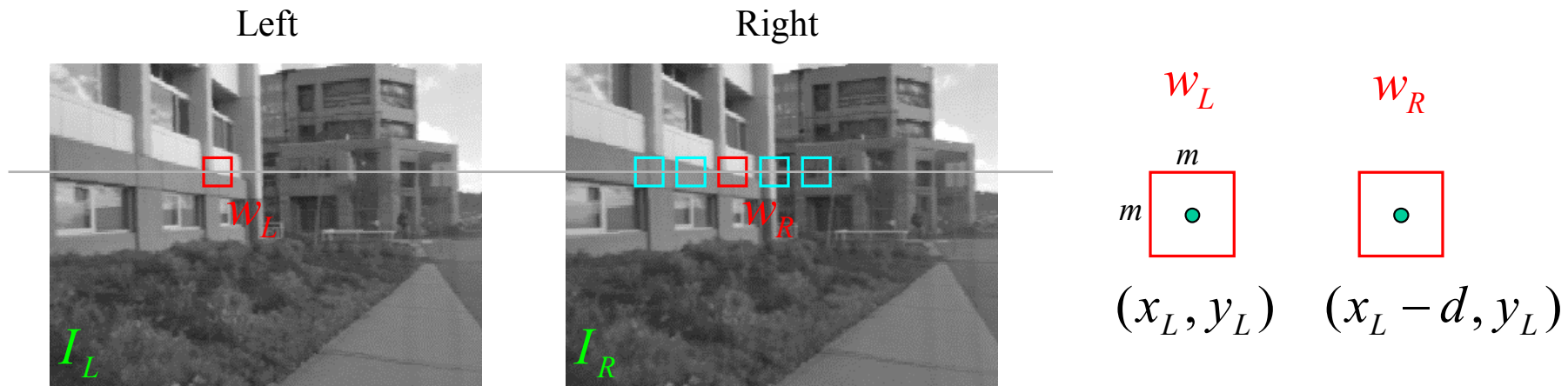# Correspondence using window matching

Left

Right



scanline

Criterion function:

error

disparity

# Sum of Squared (Pixel) Differences

Left

Right

$w_L$

$w_R$

$m$

$m$

$(x_L, y_L)$   $(x_L - d, y_L)$

$w_L$ and $w_R$ are corresponding $m$ by $m$ windows of pixels.

We define the window function :

$$W_m(x,y) = \{u,v \mid x - \tfrac{m}{2} \leq u \leq x + \tfrac{m}{2}, y - \tfrac{m}{2} \leq v \leq y + \tfrac{m}{2}\}$$

The SSD cost measures the intensity difference as a function of disparity :

$$C_r(x,y,d) = \sum_{(u,v) \in W_m(x,y)} [I_L(u,v) - I_R(u-d,v)]^2$$

6

# Image Normalization

- Even when the cameras are identical models, there can be differences in gain and sensitivity.
- The cameras do not see exactly the same surfaces, so their overall light levels can differ.
- For these reasons and more, it is a good idea to normalize the pixels in each window:

$$\bar{I} = \frac{1}{|W_m(x,y)|} \sum_{(u,v) \in W_m(x,y)} I(u,v) \qquad \text{Average pixel}$$

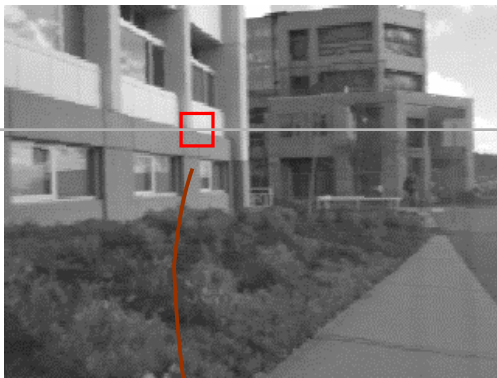$$\|I\|_{W_m(x,y)} = \sqrt{\sum_{(u,v) \in W_m(x,y)} [I(u,v)]^2} \qquad \text{Window magnitude}$$

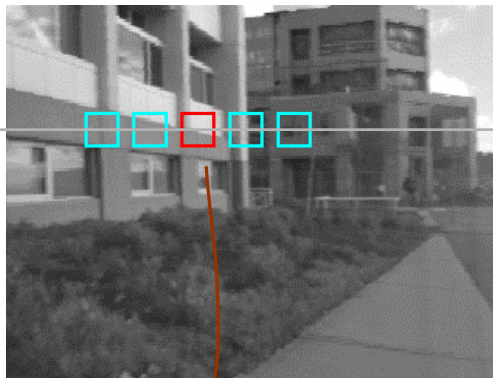$$\hat{I}(x,y) = \frac{I(x,y) - \bar{I}}{\|I - \bar{I}\|_{W_m(x,y)}} \qquad \text{Normalized pixel}$$
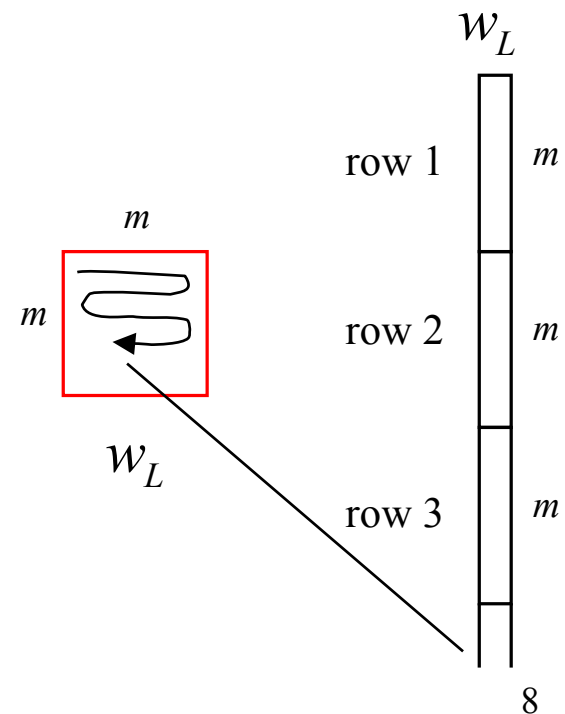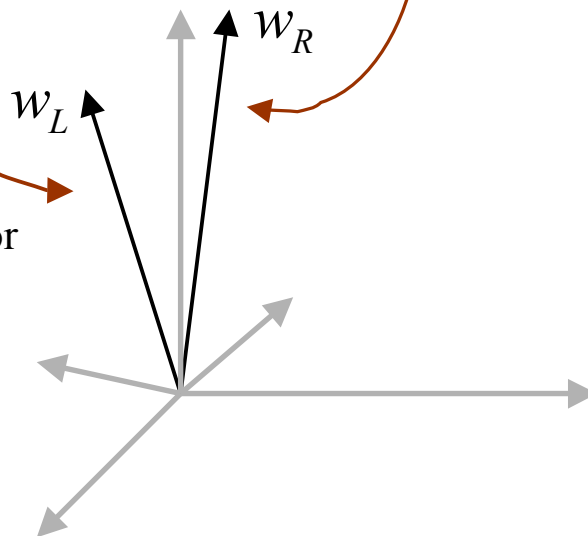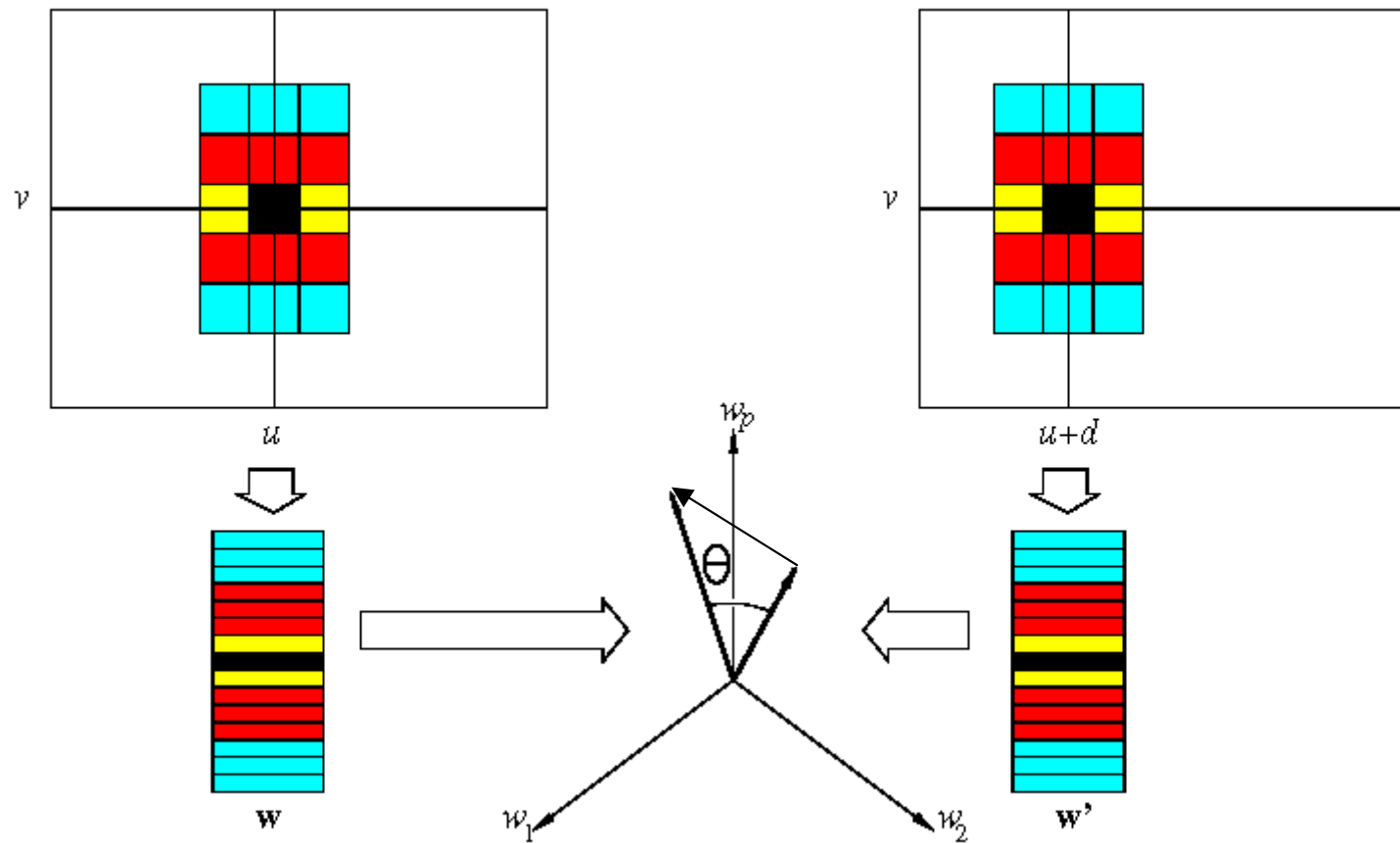
# Images as Vectors

Left

Right

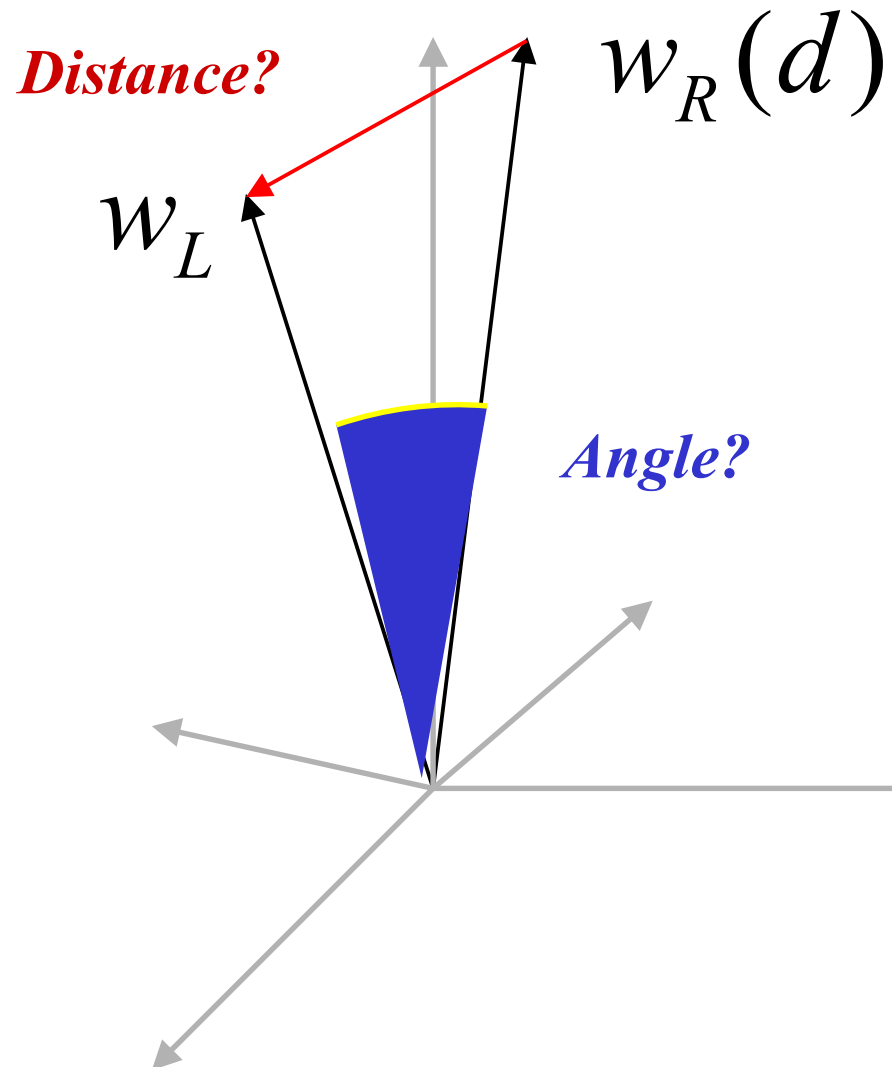"Unwrap" image to form vector, using raster scan order

$w_R$

$w_L$

$w_L$

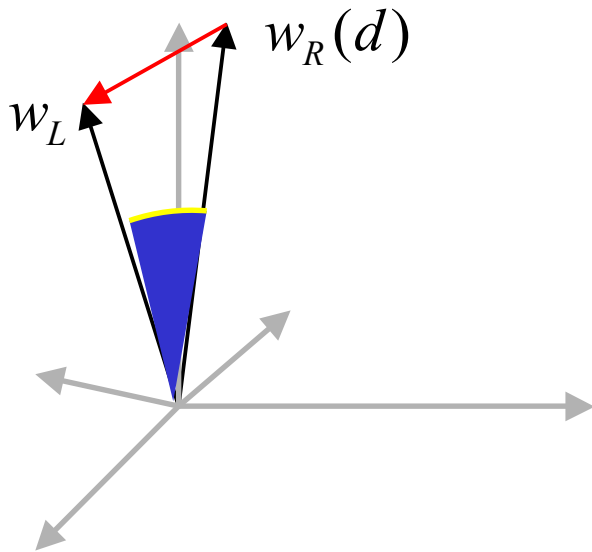Each window is a vector in an $m^2$ dimensional vector space. Normalization makes them unit length.

$m$

$m$

$w_L$

$w_L$

row 1 | $m$

row 2 | $m$

row 3 | $m$

8

# Image windows as vectors

# Possible metrics

**Distance?**

$w_R(d)$

$w_L$

**Angle?**

# Image Metrics

<span style="color:red">(Normalized) Sum of Squared Differences</span>

$$C_{\mathrm{SSD}}(d) = \sum_{(u,v)\in W_m(x,y)} [\hat{I}_L(u,v) - \hat{I}_R(u-d,v)]^2$$

$$= \left\| w_L - w_R(d) \right\|^2$$

<span style="color:blue">Normalized Correlation</span>

$$C_{\mathrm{NC}}(d) = \sum_{(u,v)\in W_m(x,y)} \hat{I}_L(u,v)\hat{I}_R(u-d,v)$$

$$= w_L \cdot w_R(d) = \cos\theta$$

$$d^* = \arg\min_d \left\| w_L - w_R(d) \right\|^2 = \arg\max_d \; w_L \cdot w_R(d)$$
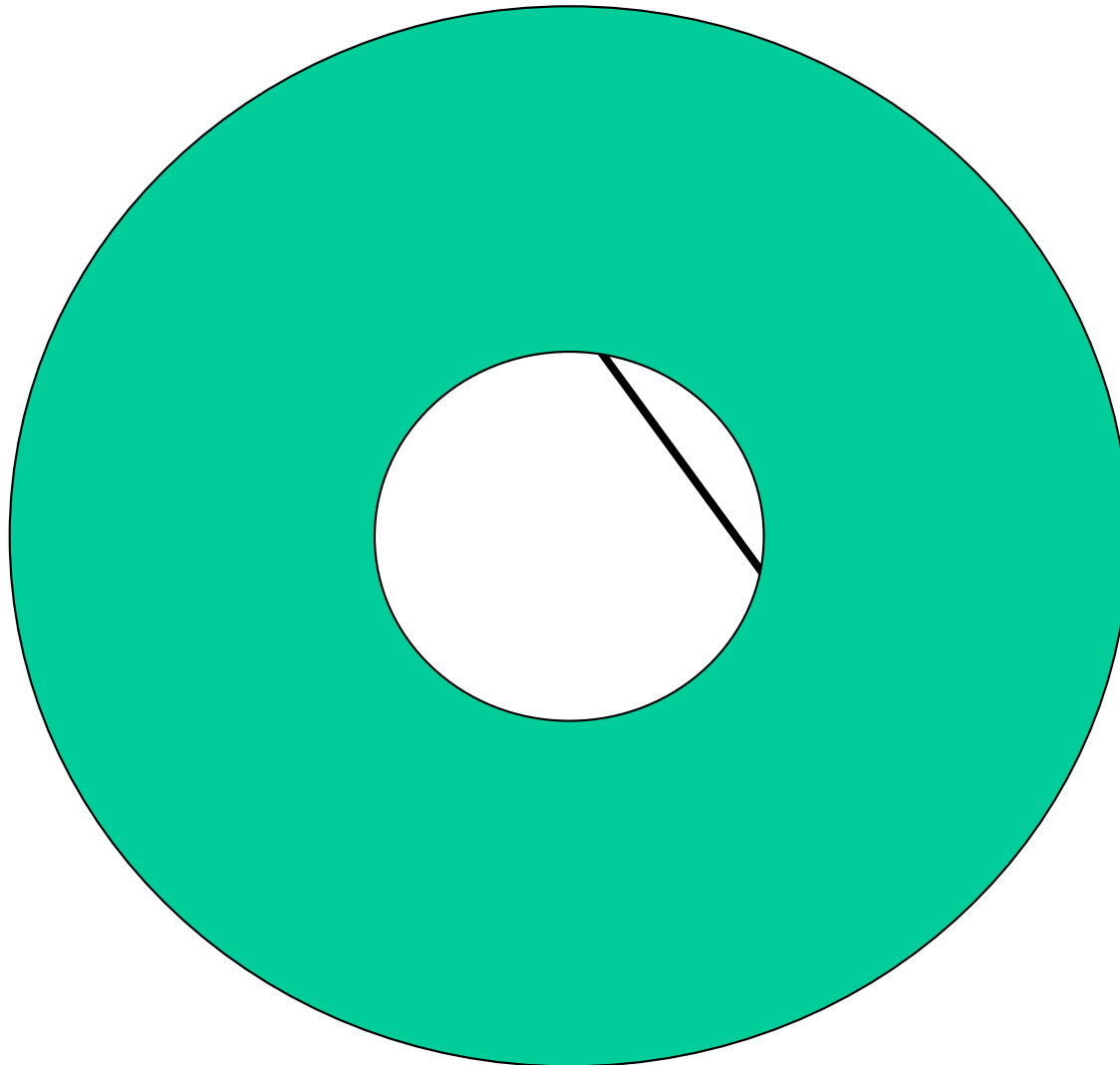
$w_R(d)$

$w_L$

11

# Local Features

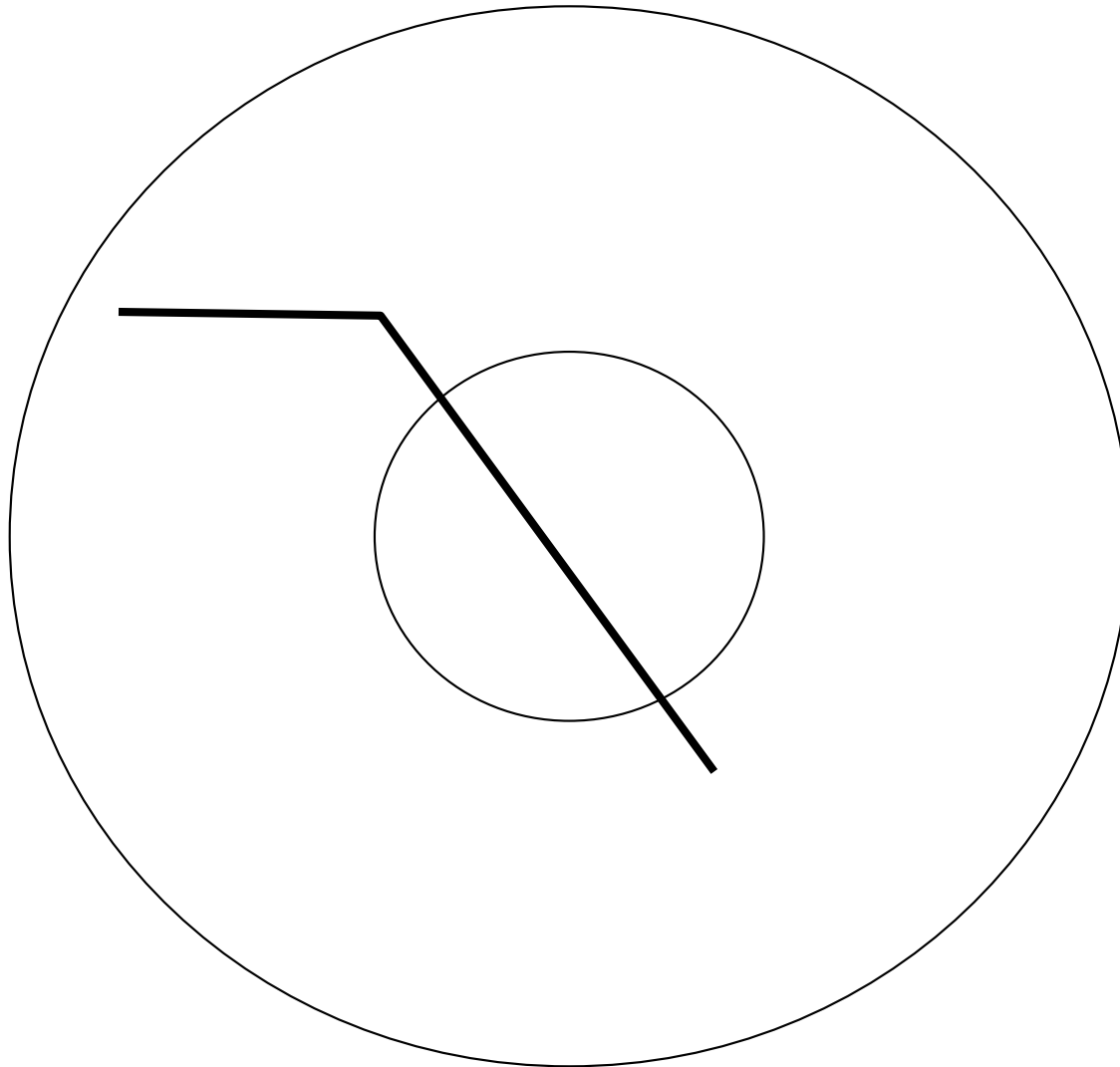Not all points are equally good for matching…

# Aperture Problem and Normal Flow
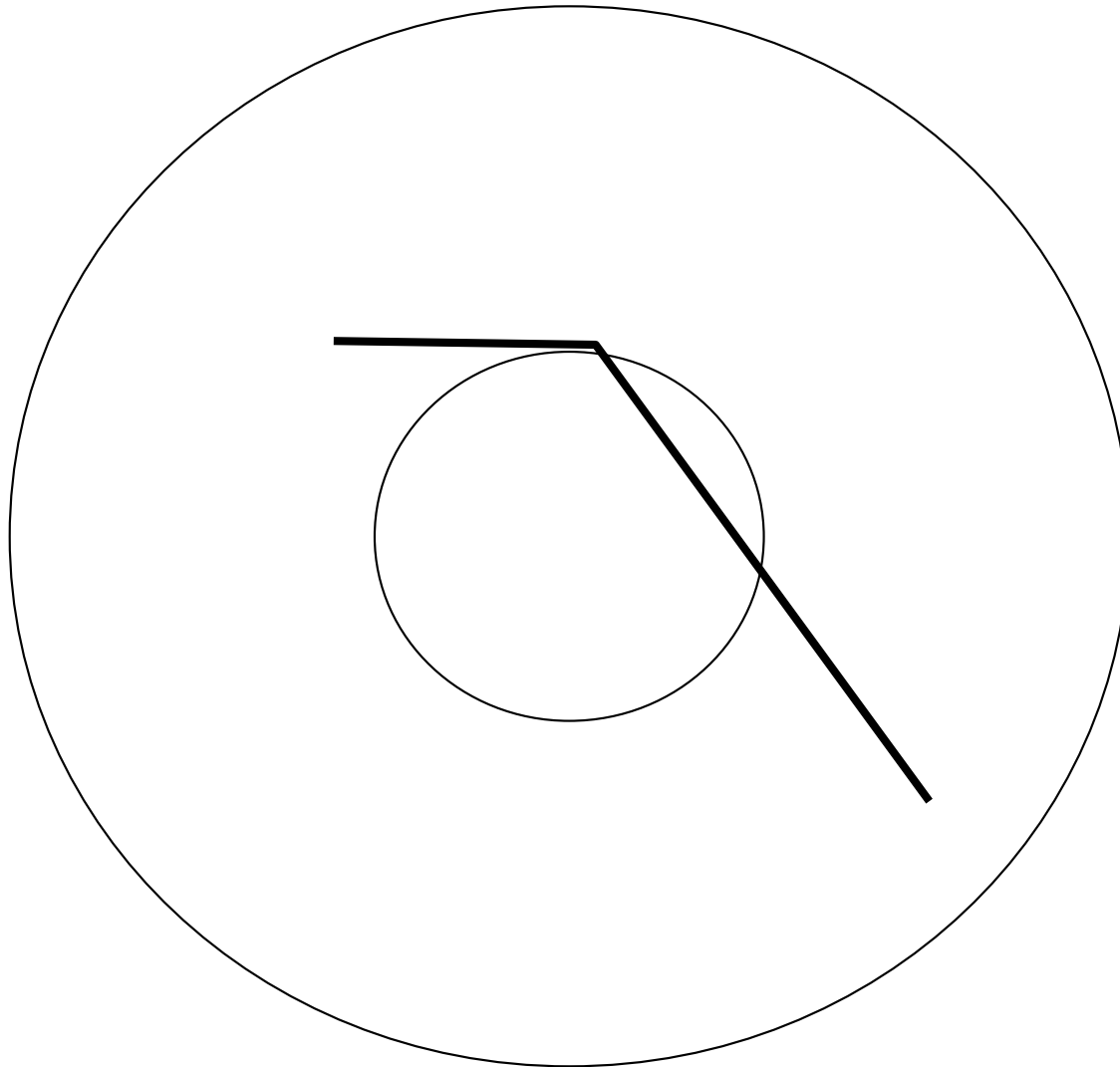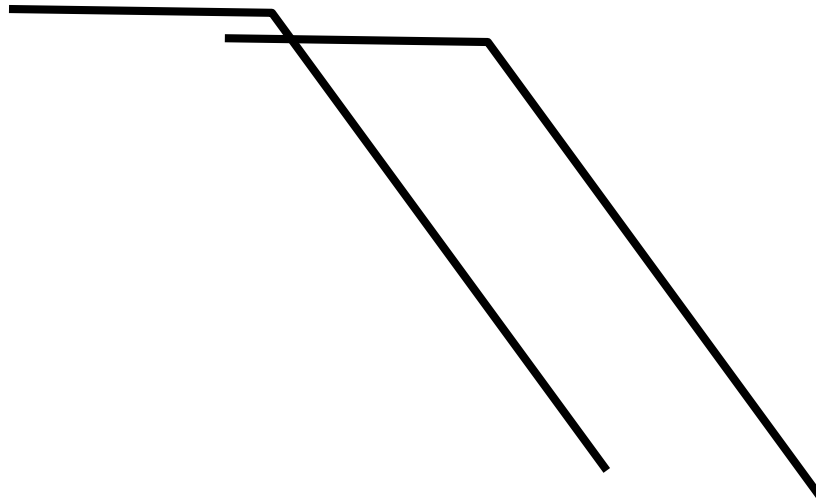
# Aperture Problem and Normal Flow

# Aperture Problem and Normal Flow

# Aperture Problem and Normal Flow

# Aperture Problem and Normal Flow

# Aperture Problem and Normal Flow

# (Review) Differential approach: Optical flow constraint equation

Brightness should stay constant as you track motion

$$I(x + u\delta t, y + v\delta t, t + \delta t) = I(x, y, t)$$

1st order Taylor series, valid for small $\delta t$

$$I(x, y, t) + u\delta t I_x + v\delta t I_y + \delta t I_t = I(x, y, t)$$

Constraint equation

$$\boxed{uI_x + vI_y + I_t = 0}$$

"BCCE" - Brightness Change Constraint Equation

# Aperture Problem and Normal Flow

**The gradient constraint:**

$$I_x u + I_y v + I_t = 0$$

$$\nabla I \bullet \vec{U} = 0$$

**Defines a line in the *(u,v)* space**

**Normal Flow:**

$$u_\perp = -\frac{I_t}{|\nabla I|} \frac{\nabla I}{|\nabla I|}$$

*v*

*u*

20

# Combining Local Constraints



$$\nabla I^1 \bullet U = -I_t^1$$

$$\nabla I^2 \bullet U = -I_t^2$$

$$\nabla I^3 \bullet U = -I_t^3$$

etc.

# Lucas-Kanade: Integrate gradients over a Patch

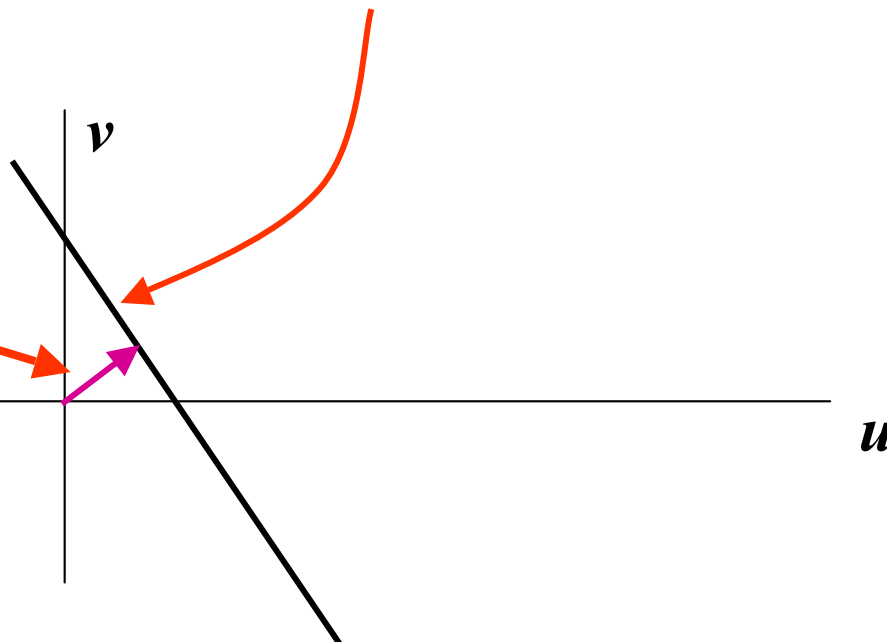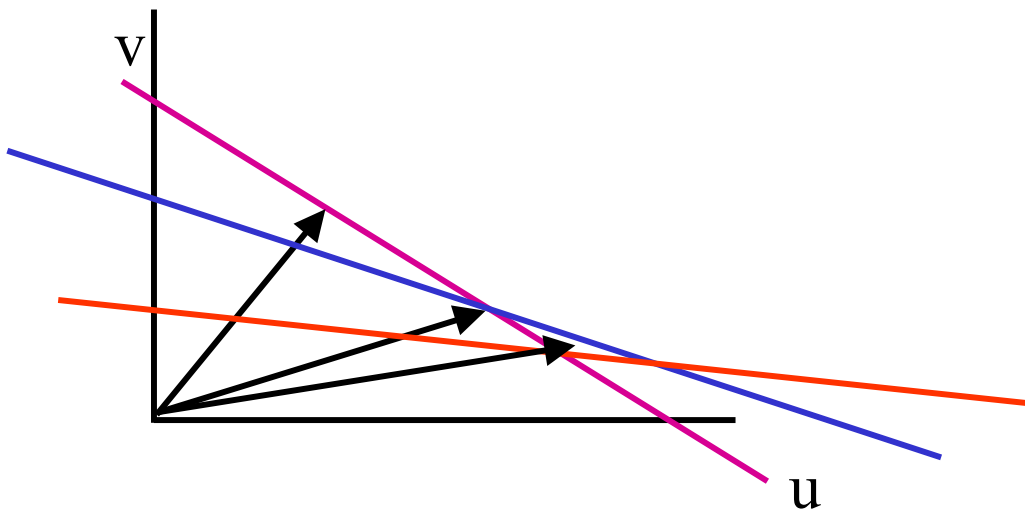Assume a single velocity for all pixels within an image patch

$$E(u,v) = \sum_{x,y \in \Omega} \left( I_x(x,y)u + I_y(x,y)v + I_t \right)^2$$

Solve with:

$$\begin{bmatrix} \sum I_x^2 & \sum I_x I_y \\ \sum I_x I_y & \sum I_y^2 \end{bmatrix} \begin{pmatrix} u \\ v \end{pmatrix} = -\begin{pmatrix} \sum I_x I_t \\ \sum I_y I_t \end{pmatrix}$$

On the LHS: sum of the 2x2 outer product tensor of the gradient vector

$$\left( \sum \nabla I \nabla I^T \right) \vec{U} = -\sum \nabla I I_t$$

# Local Patch Analysis

# Selecting Good Features

- What's a "good feature"?
  - Satisfies brightness constancy
  - Has sufficient texture variation
  - Does not have too much texture variation
  - Corresponds to a "real" surface patch
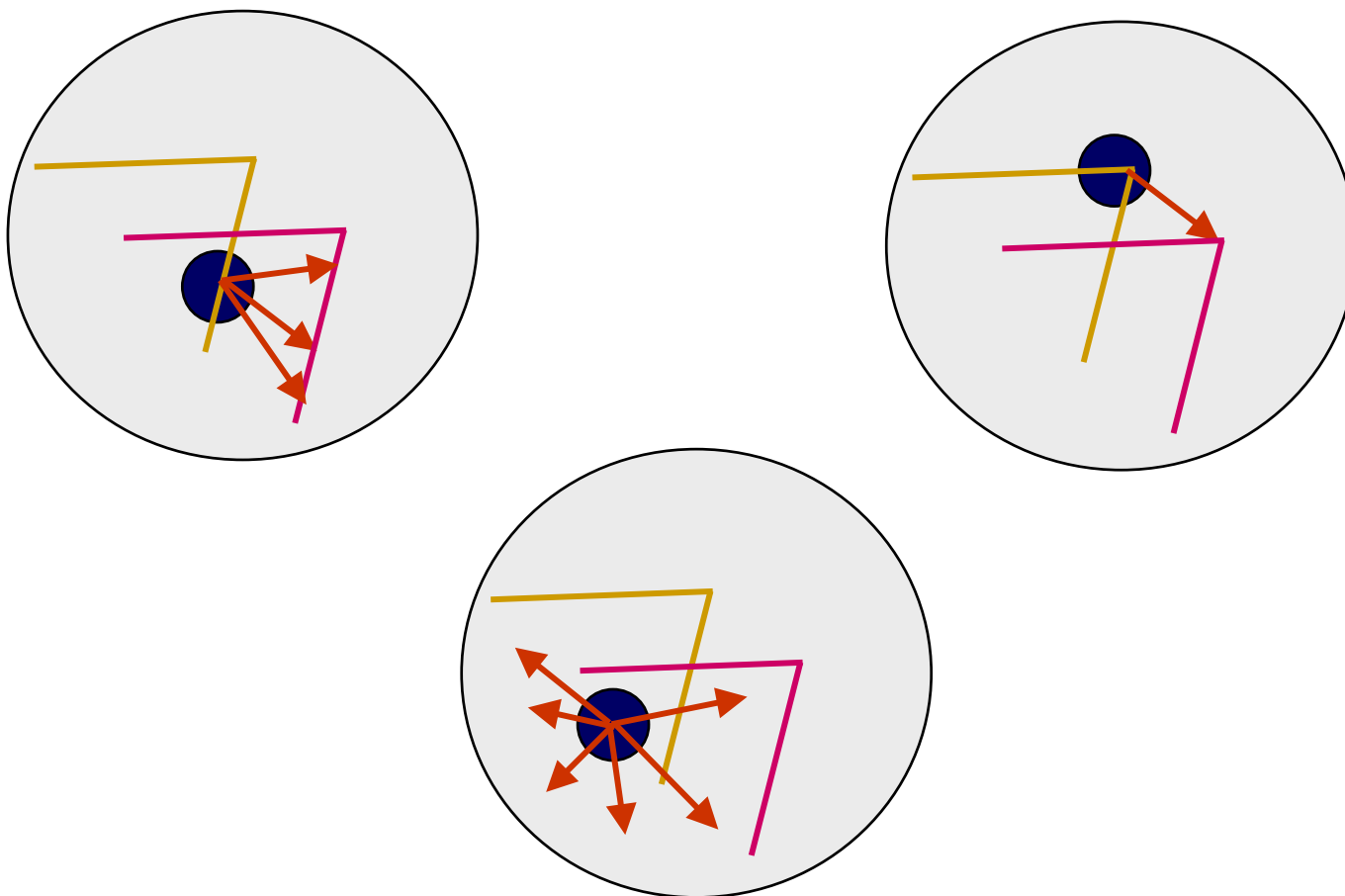  - Does not deform too much over time

# Good Features to Track

$$\begin{bmatrix} \sum I_x^2 & \sum I_x I_y \\ \sum I_x I_y & \sum I_y^2 \end{bmatrix} \begin{pmatrix} u \\ v \end{pmatrix} = -\begin{pmatrix} \sum I_x I_t \\ \sum I_y I_t \end{pmatrix}$$

$$\mathbf{A} \qquad \mathbf{u} \quad = \quad \mathbf{b}$$

## When is This Solvable?

- **A** should be invertible
- **A** should not be too small due to noise
  - eigenvalues $\lambda_1$ and $\lambda_2$ of **A** should not be too small
- **A** should be well-conditioned
  - $\lambda_1/\lambda_2$ should not be too large ($\lambda_1$ = larger eigenvalue)

Both conditions satisfied when $min(\lambda_1, \lambda_2) > c$

# Harris detector

Same idea, based on the idea of auto-correlation



Important difference in all directions => interest point

# Harris detector

Auto-correlation function for a point $(x, y)$ and a shift $(\Delta x, \Delta y)$

$$f(x, y) = \sum_{(x_k, y_k) \in W} (I(x_k, y_k) - I(x_k + \Delta x, y_k + \Delta y))^2$$

Discret shifts can be avoided with the auto-correlation matrix

with $I(x_k + \Delta x, y_k + \Delta y) = I(x_k, y_k) + (I_x(x_k, y_k) \quad I_y(x_k, y_k)) \begin{pmatrix} \Delta x \\ \Delta y \end{pmatrix}$

$$f(x, y) = \sum_{(x_k, y_k) \in W} \left( \left( I_x(x_k, y_k) \quad I_y(x_k, y_k) \right) \begin{pmatrix} \Delta x \\ \Delta y \end{pmatrix} \right)^2$$
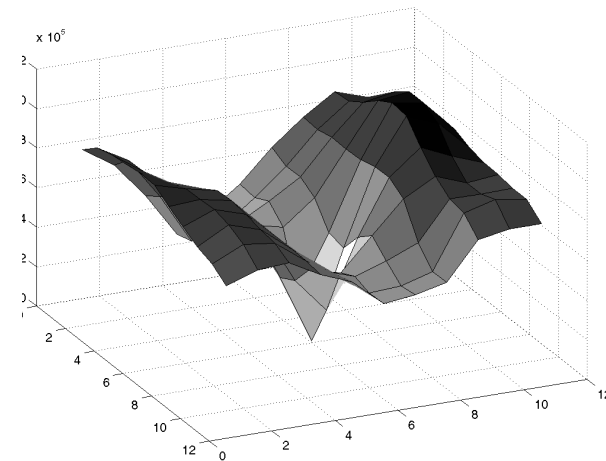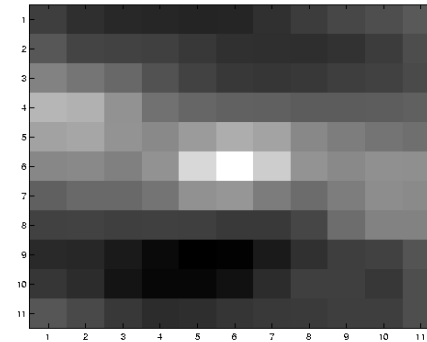
# Harris detector

Auto-correlation matrix

$$= \begin{pmatrix} \Delta x & \Delta y \end{pmatrix} \begin{bmatrix} \sum_{(x_k,y_k)\in W} (I_x(x_k,y_k))^2 & \sum_{(x_k,y_k)\in W} I_x(x_k,y_k)I_y(x_k,y_k) \\ \sum_{(x_k,y_k)\in W} I_x(x_k,y_k)I_y(x_k,y_k) & \sum_{(x_k,y_k)\in W} (I_y(x_k,y_k))^2 \end{bmatrix} \begin{pmatrix} \Delta x \\ \Delta y \end{pmatrix}$$
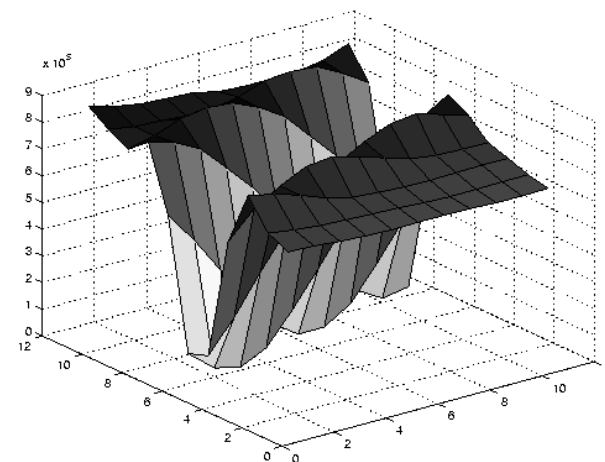
- Auto-correlation matrix
  - captures the structure of the local neighborhood
  - measure based on eigenvalues of this matrix
    - 2 strong eigenvalues  => interest point
    - 1 strong eigenvalue   => contour
    - 0 eigenvalue          => uniform region

- Interest point detection
  - threshold on the eigenvalues
  - local maximum for localization

# Selecting Good Features



$\lambda_1$ and $\lambda_2$ are large

# Selecting Good Features



large $\lambda_1$, small $\lambda_2$

30

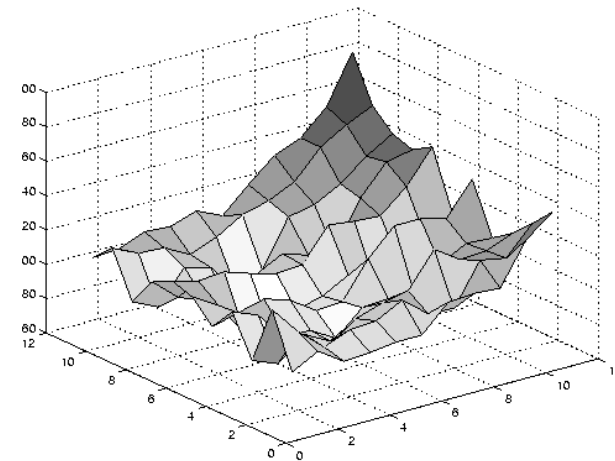# Selecting Good Features



small $\lambda_1$, small $\lambda_2$

# Feature Distortion

- Feature may change shape over time
    - Need a distortion model to really make this work



Find displacement (u,v) that minimizes SSD error over feature region

$$\sum_{(x,y)\in F\subset J}[I(W_x(x,y), W_y(x,y)) - J(x,y)]^2$$

(minimize with respect to $W_x$ and $W_y$)

*Shi and Tomasi: use affine model for verification*

$$W_x(x,y) = ax + by + c$$
$$W_y(x,y) = ex + fy + g$$

32

# Affine Motion



$$u(x, y) = a_0 + a_1 x + a_2 y$$

$$v(x, y) = a_3 + a_4 x + a_5 y$$

$\mathbf{u}(\mathbf{x}; \mathbf{a}) = (u(x,y), v(x,y))$

$I(\mathbf{x}, t-1)$

$\mathbf{x} + \mathbf{u}(\mathbf{x}; \mathbf{a})$

Warp

$\mathbf{x}$

$I(\mathbf{x}+\mathbf{u}(\mathbf{x}; \mathbf{a}), t-1) = I(\mathbf{x}, t)$

*(Brightness Constancy Assumption)*

33

# Affine Motion

$$u(x, y) = a_1 + a_2 x + a_3 y$$
$$v(x, y) = a_4 + a_5 x + a_6 y$$

Substituting into the B.C.C.E.:

$$I_x \cdot u + I_y \cdot v + I_t \approx 0$$

$$I_x (a_1 + a_2 x + a_3 y) + I_y (a_4 + a_5 x + a_6 y) + I_t \approx 0$$

**Each pixel provides 1 linear constraint in 6 *global* unknowns**
***(minimum 6 pixels necessary)***

Least Square Minimization  (over all pixels):

$$Err(\vec{a}) = \sum \left[ I_x (a_1 + a_2 x + a_3 y) + I_y (a_4 + a_5 x + a_6 y) + I_t \right]^2$$

# Tracking vs. Indexing

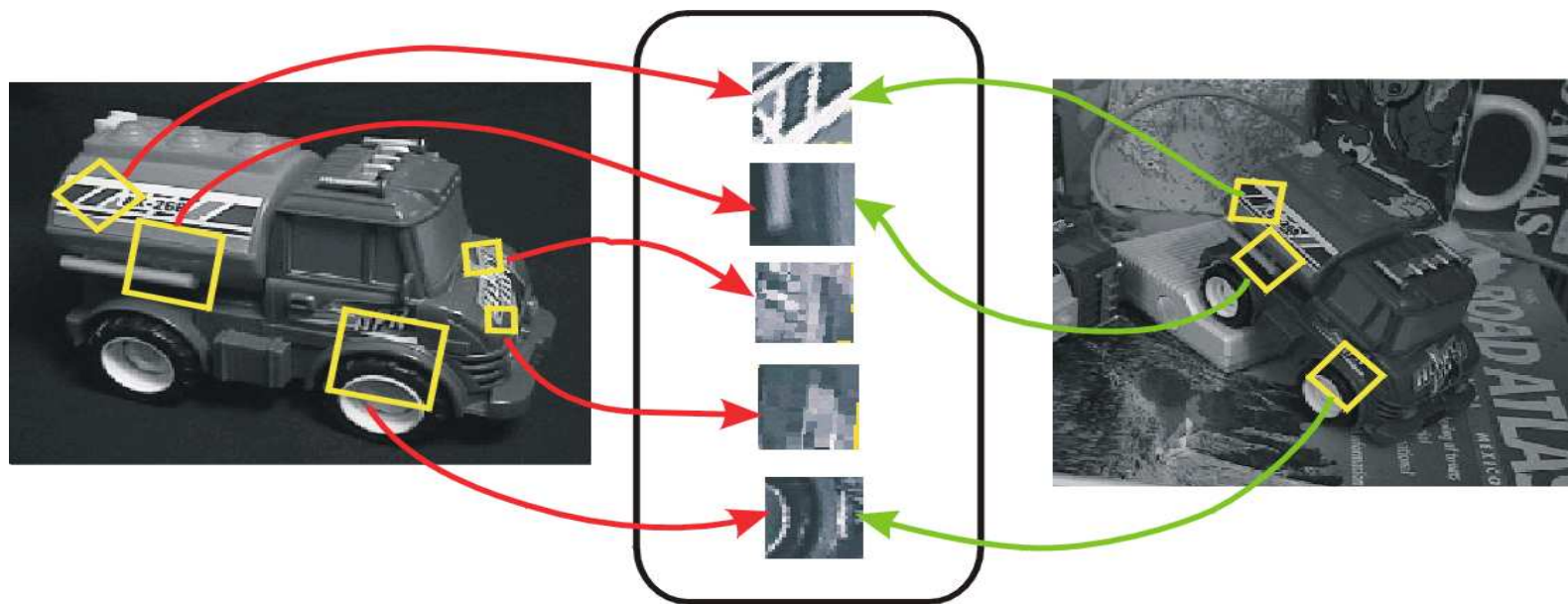But….

What if you can't track over time?

# Today

Interesting points, correspondence, affine patch tracking

**Scale and rotation invariant descriptors**

# Recognition and Matching Based on Local Invariant Features

# Invariant Local Features

- Image content is transformed into local feature coordinates that are invariant to translation, rotation, scale, and other imaging parameters



**SIFT Features**
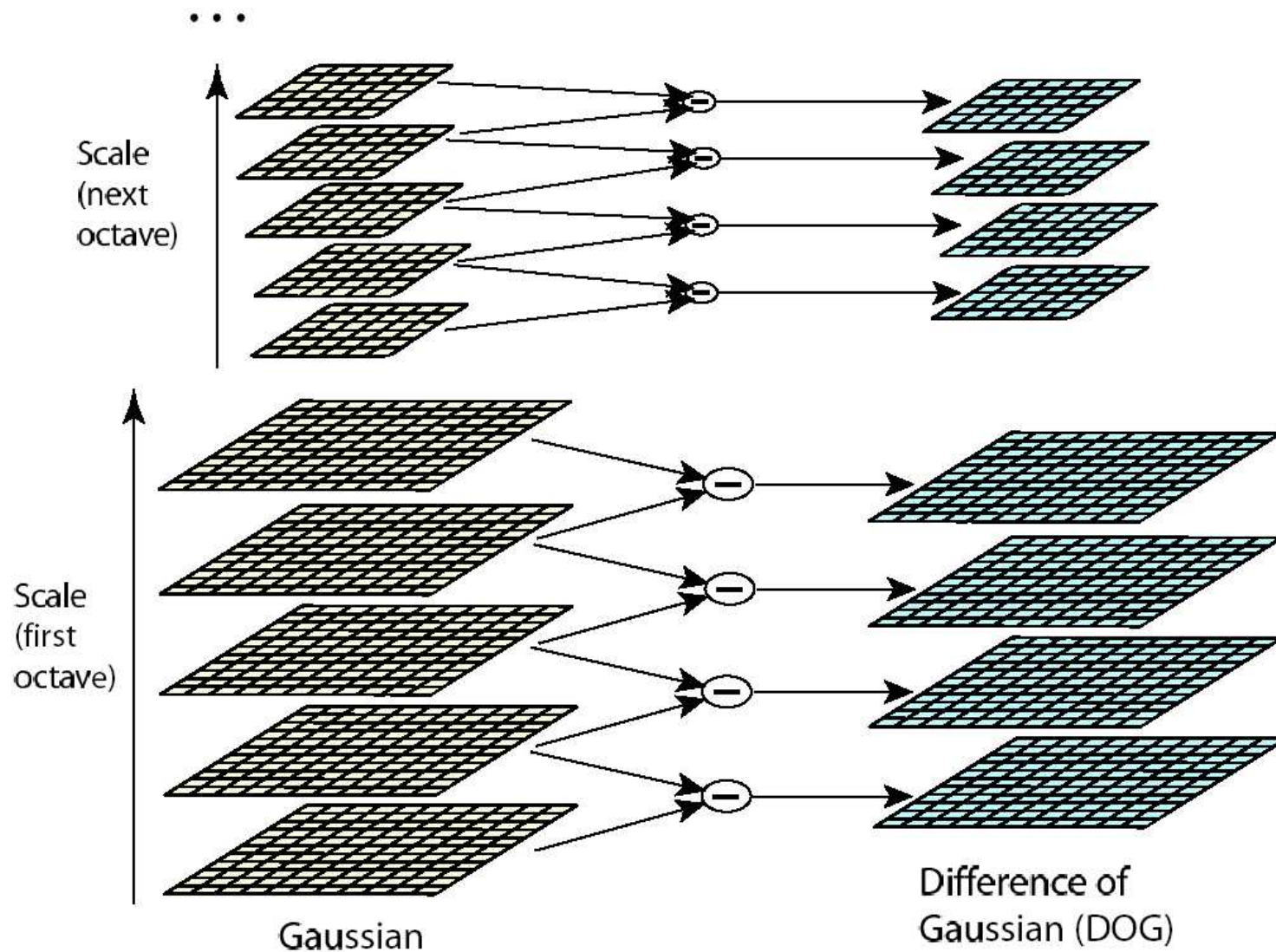
# Advantages of invariant local features

- **Locality:** features are local, so robust to occlusion and clutter (no prior segmentation)

- **Distinctiveness:** individual features can be matched to a large database of objects

- **Quantity:** many features can be generated for even small objects

- **Efficiency:** close to real-time performance

- **Extensibility:** can easily be extended to wide range of differing feature types, with each adding robustness

# Scale invariance

**Requires a method to repeatably select points in location and scale:**

- The only reasonable scale-space kernel is a Gaussian (Koenderink, 1984; Lindeberg, 1994)

- An efficient choice is to detect peaks in the difference of Gaussian pyramid (Burt & Adelson, 1983; Crowley & Parker, 1984 – but examining more scales)

- Difference-of-Gaussian with constant ratio of scales is a close approximation to Lindeberg's scale-normalized Laplacian (can be shown from the heat diffusion equation)

# Scale space processed one octave at a time



. . .

Scale (next octave)

Scale (first octave)

Gaussian

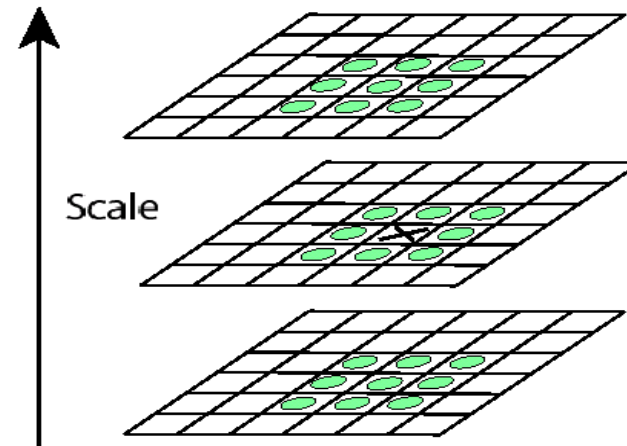Difference of Gaussian (DOG)

46

# Key point localization

- Detect maxima and minima of difference-of-Gaussian in scale space
- Fit a quadratic to surrounding values for sub-pixel and sub-scale interpolation (Brown & Lowe, 2002)
- Taylor expansion around point:

$$D(\mathbf{x}) = D + \frac{\partial D^T}{\partial \mathbf{x}} \mathbf{x} + \frac{1}{2}\mathbf{x^T}\frac{\partial^2 D}{\partial \mathbf{x}^2}\mathbf{x}$$
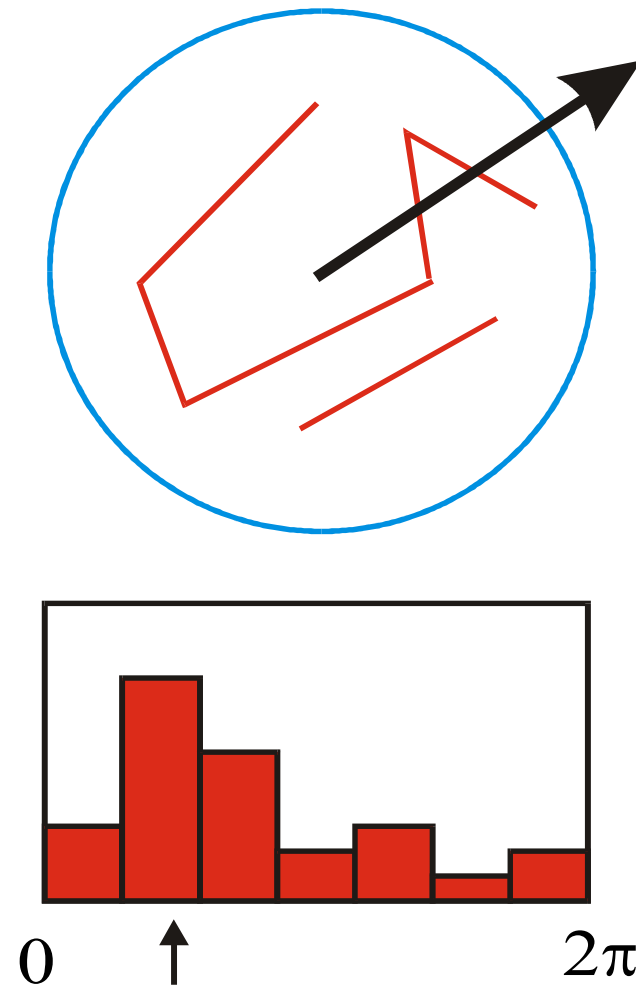
- Offset of extremum (use finite differences for derivatives):

$$\hat{\mathbf{x}} = -\frac{\partial^2 D}{\partial \mathbf{x}^2}^{-1} \frac{\partial D}{\partial \mathbf{x}}$$
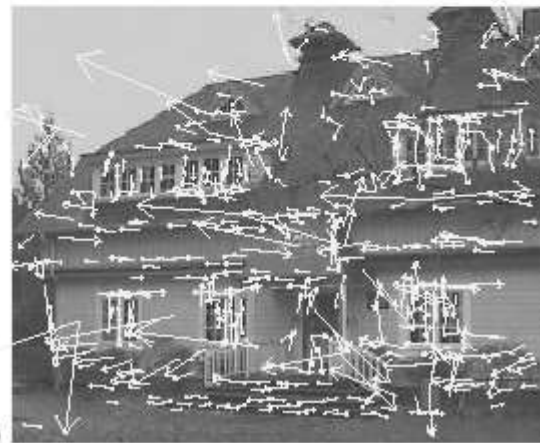


Scale

# Select canonical orientation

- Create histogram of local gradient directions computed at selected scale

- Assign canonical orientation at peak of smoothed histogram

- Each key specifies stable 2D coordinates (x, y, scale, orientation)

$0$     ↑      $2\pi$
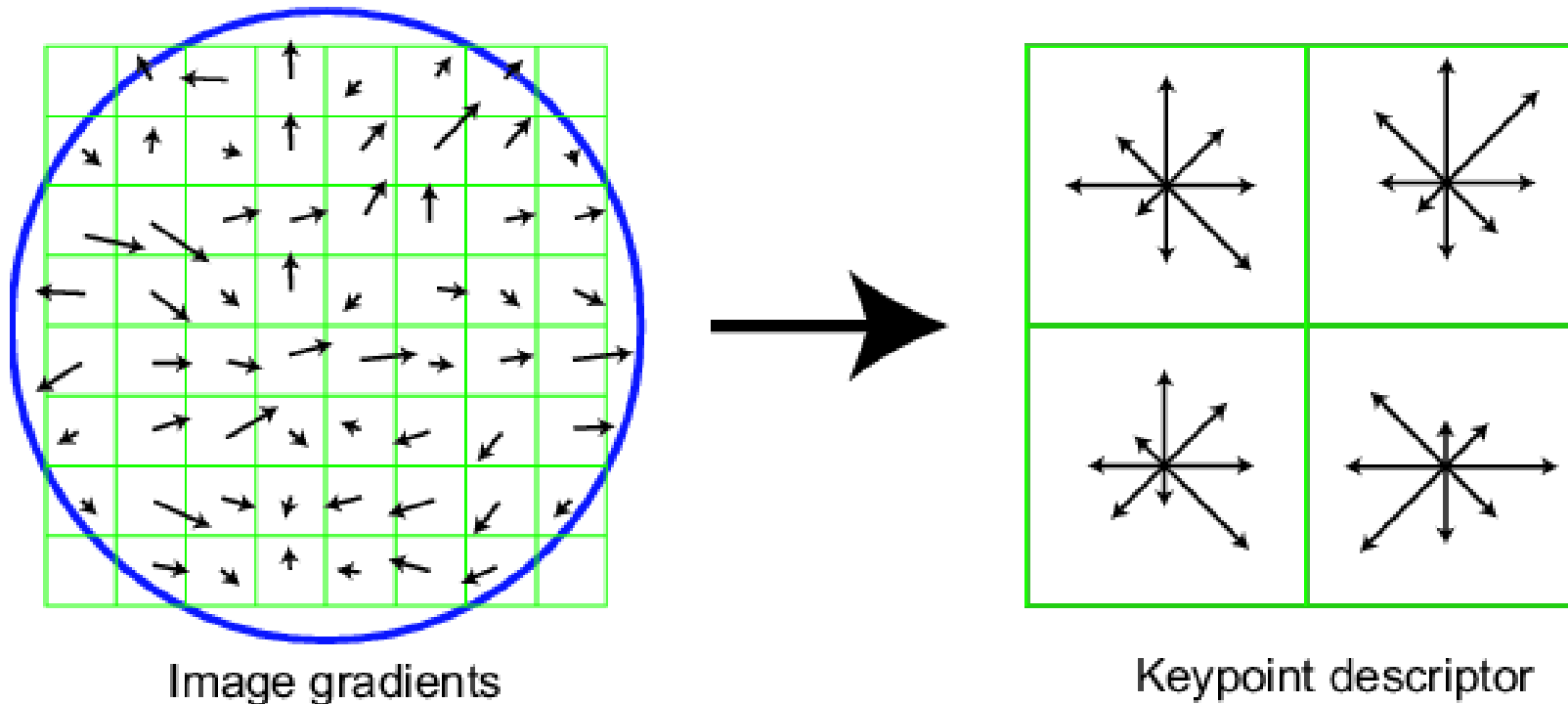
# Example of keypoint detection

Threshold on value at DOG peak and on ratio of principle curvatures (Harris approach)



**(a)** 233x189 image
**(b)** 832 DOG extrema
**(c)** 729 left after peak value threshold
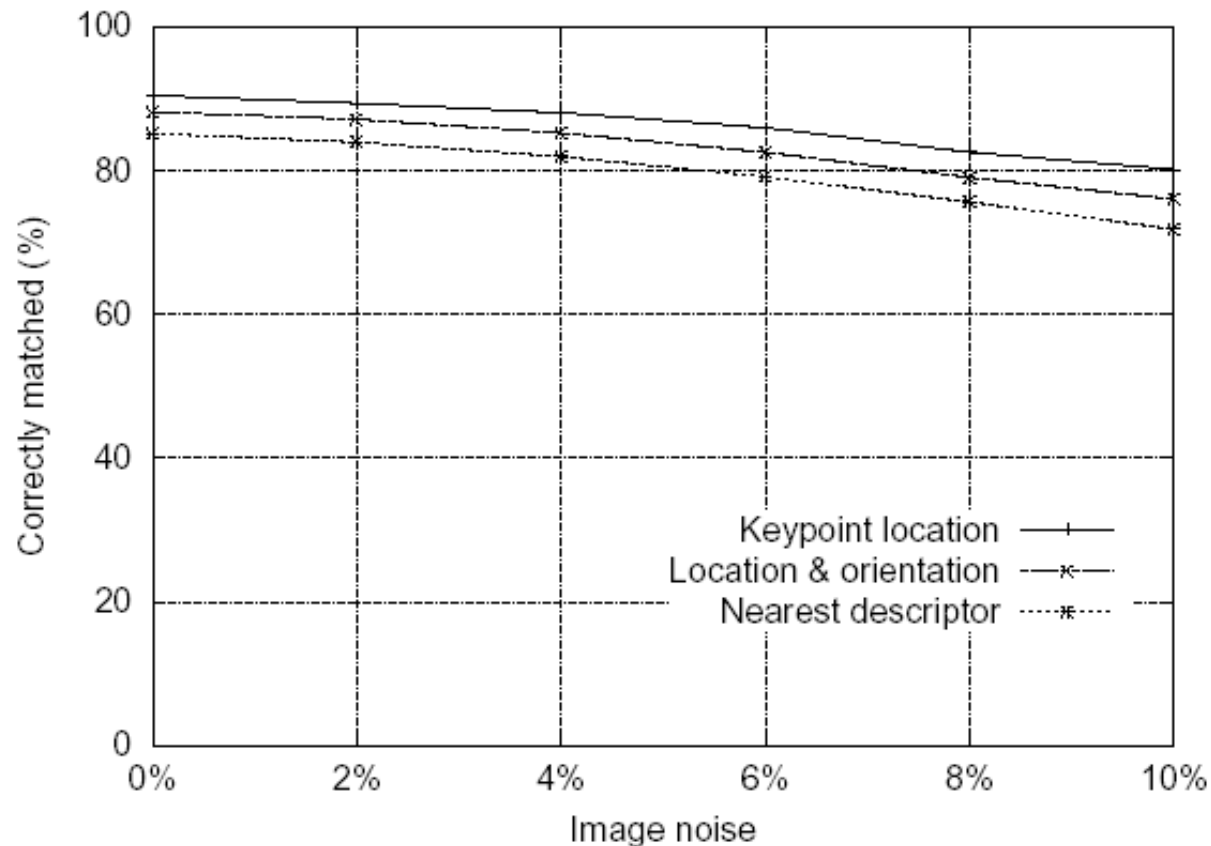**(d)** 536 left after testing ratio of principle curvatures

49

# SIFT vector formation

- Thresholded image gradients are sampled over 16x16 array of locations in scale space

- Create array of orientation histograms

- 8 orientations x 4x4 histogram array = 128 dimensions



Image gradients                    Keypoint descriptor
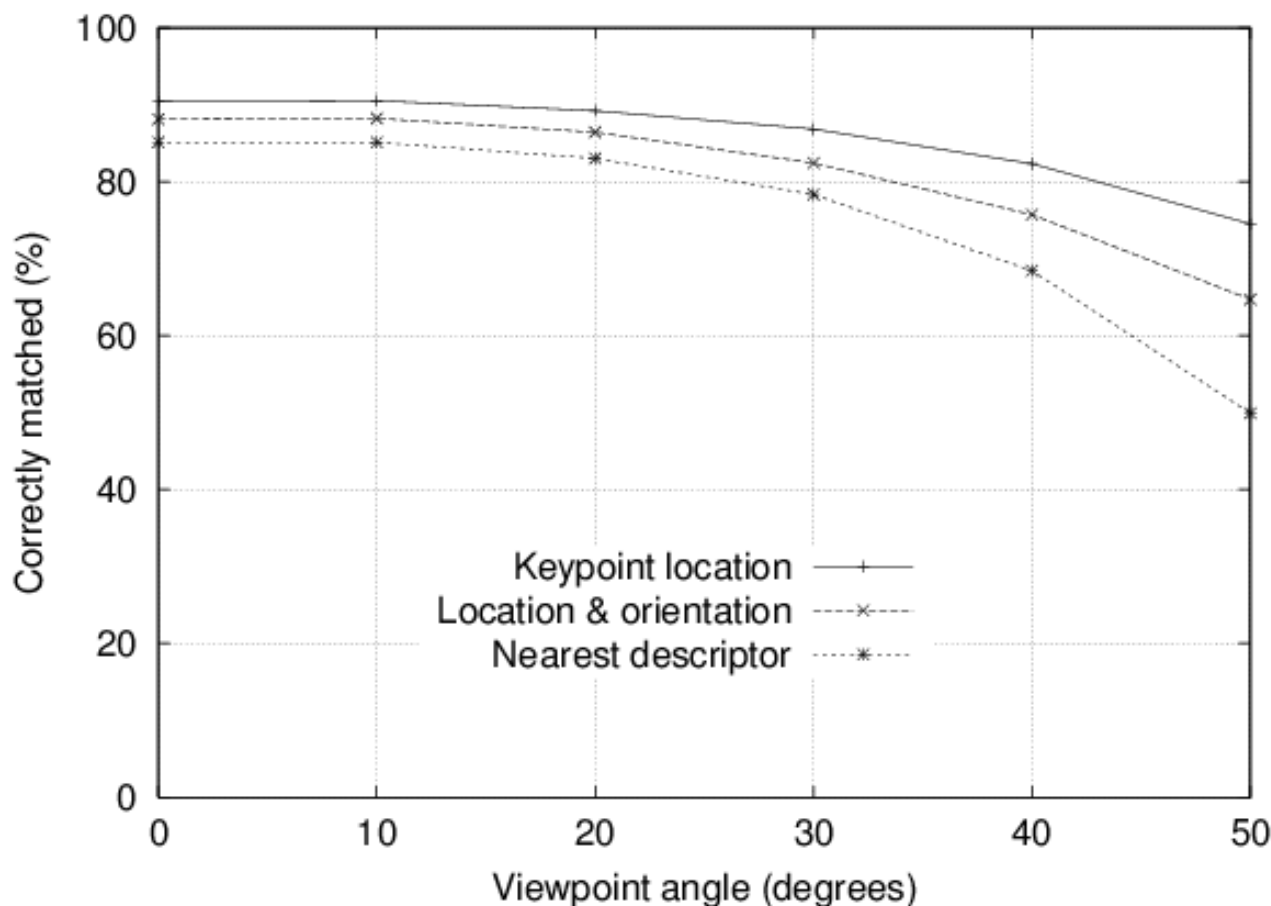
# Feature stability to noise

- Match features after random change in image scale & orientation, with differing levels of image noise

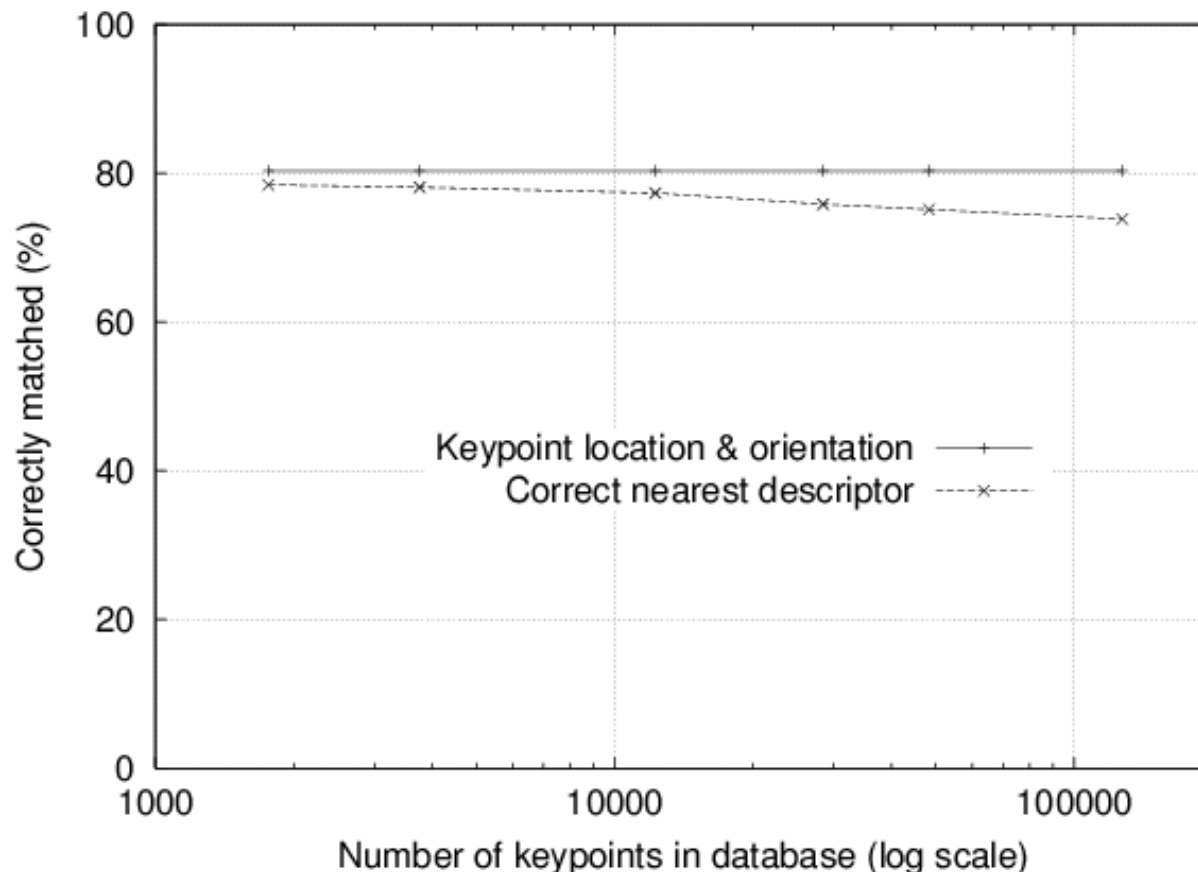- Find nearest neighbor in database of 30,000 features

# Feature stability to affine change

- Match features after random change in image scale & orientation, with 2% image noise, and affine distortion
- Find nearest neighbor in database of 30,000 features

# Distinctiveness of features

- Vary size of database of features, with 30 degree affine change, 2% image noise
- Measure % correct for single nearest neighbor match

# Today

Interesting points, correspondence, affine patch
tracking

Scale and rotation invariant descriptors