

# Tracking surfaces / objects across video frames

Readings:

Paper, First dozen pages:  
Baker&Matthews, IJCV

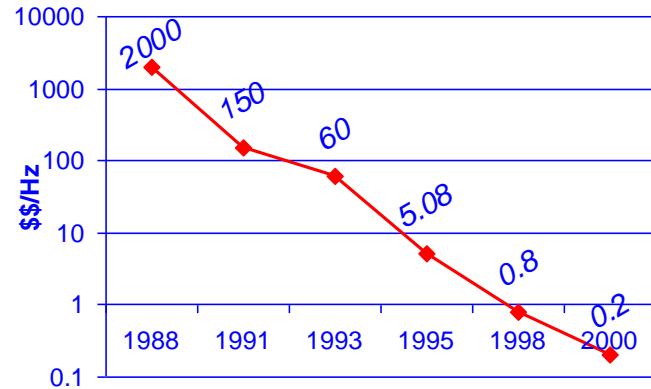
“Lukas-Kanade 20 years on”

Szeliski 8.1.3, 4.1

Ma et al Ch 4.

Forsythe and Ponce Ch 17.

# Why Visual Tracking?

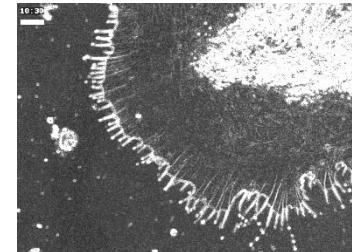
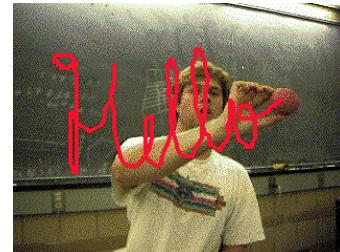
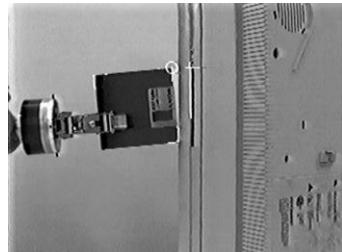


Computers are fast enough!



Related technology is cheap!

Applications abound



Motion Control

HCI

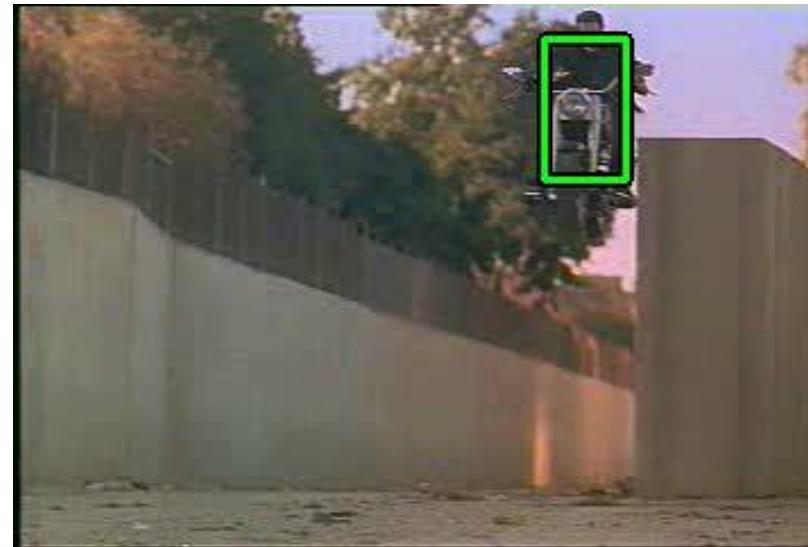
Measurement

# Applications: Watching a moving target

- Camera + computer can determine how things move in a scene over time.

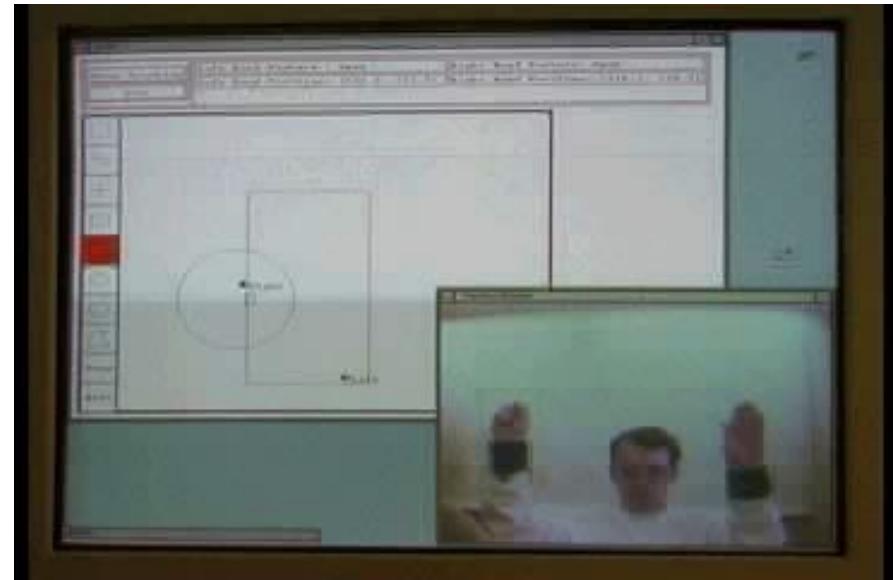
Uses:

- Security: e.g. monitoring people moving in a subway station or store
- Measurement: Speed, alert on colliding trajectories etc.



# Applications Human-Computer Interfaces

- Camera + computer tracks motions of human user and interprets this in an on-line interaction.
- Can interact with menus, buttons and e.g. drawing programs using hand movements as mouse movements, and gestures as clicking
- Furthermore, can interpret physical interactions

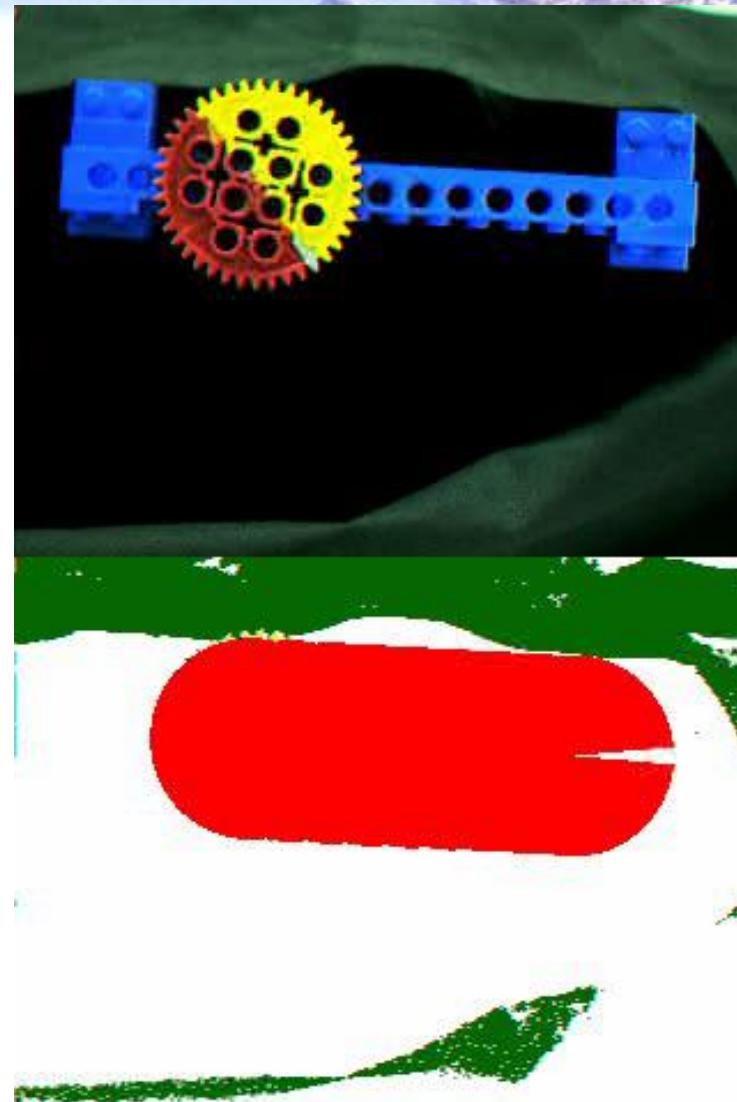


# Applications

## Interpreting tasks and functionality

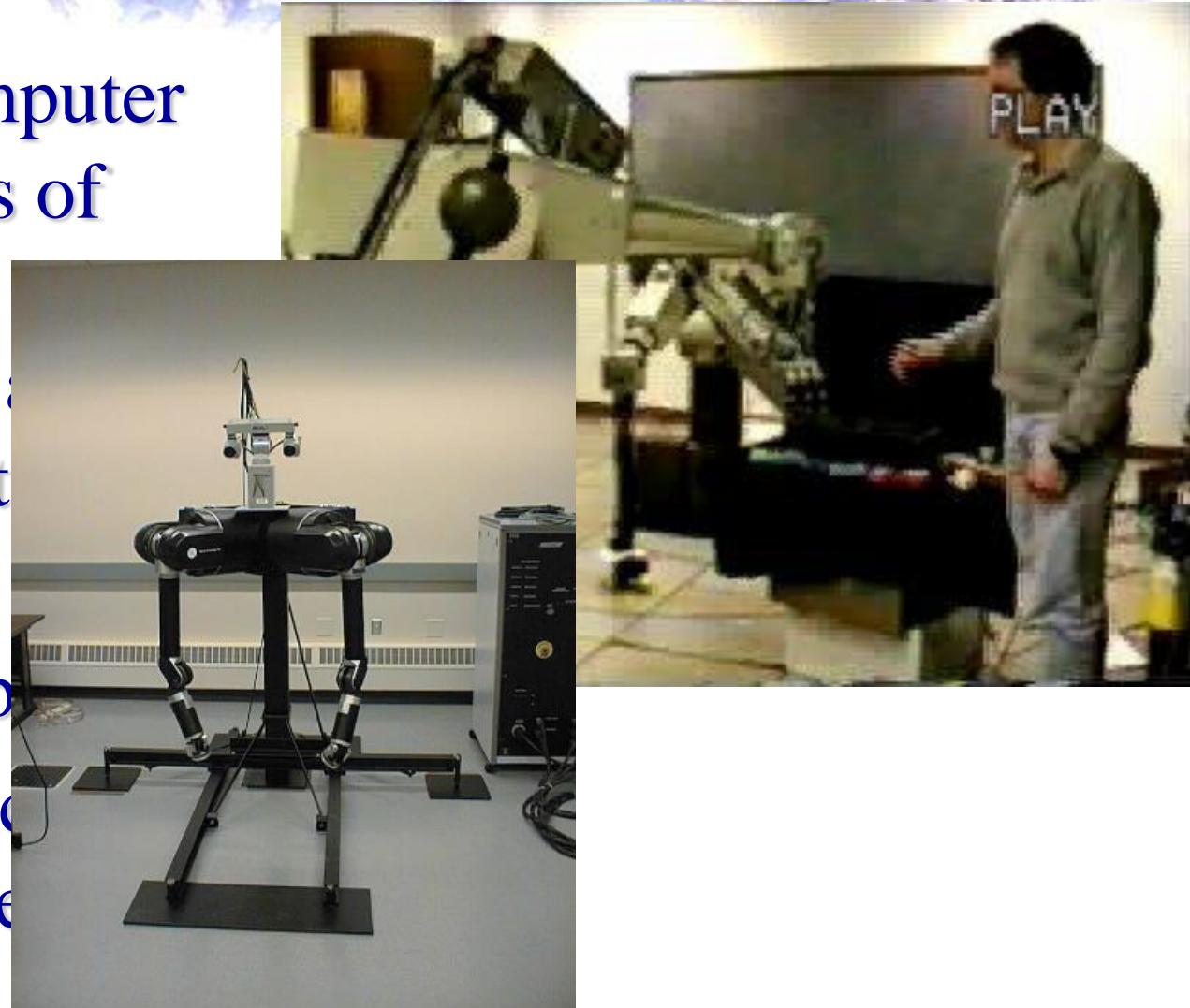
Use tracking to interpret:

1. Physics of a mechanism
2. How a car is repaired.
3. How a human cooks in a kitchen



# Applications Human-Machine Interfaces

- Camera + computer tracks motions of human user, interprets this as machine/robot to carry out task.
- Remote manipulation
- Service robotics for the handicapped and elderly



# Applications Human-Human Interfaces

- Camera + computer tracks motions and expressions of human user, interprets, codes and transmits to render at remote location
- Ultra low bandwidth video communication
- Handheld video cell phone



# BANDWIDTH and PROCESSING REQUIREMENTS: TV camera

## Binary

1 bit \*  $640 \times 480 \times 30 = 9.2 \text{ Mbits/second}$

## Grey

1 byte \*  $640 \times 480 \times 30 = 9.2 \text{ Mbytes/second}$

## Color

3 bytes \*  $640 \times 480 \times 30 = 27.6 \text{ Mbytes/second}$  (actually about 37 mbytes/sec)

Typical operation: 8x8 convolution on grey image  
64 multiplies + adds → 600 Mflops

Consider 800x600, 1200x1600 (2MP)...

Today's PC's are getting to the point they  
can process images at frame rate

For real time: process small window in image

# Characteristics of Video Processing

$\text{abs}(\text{Image 1} - \text{Image 2}) = ?$



Note: Almost all pixels change!

Constancy: The physical scene is the same

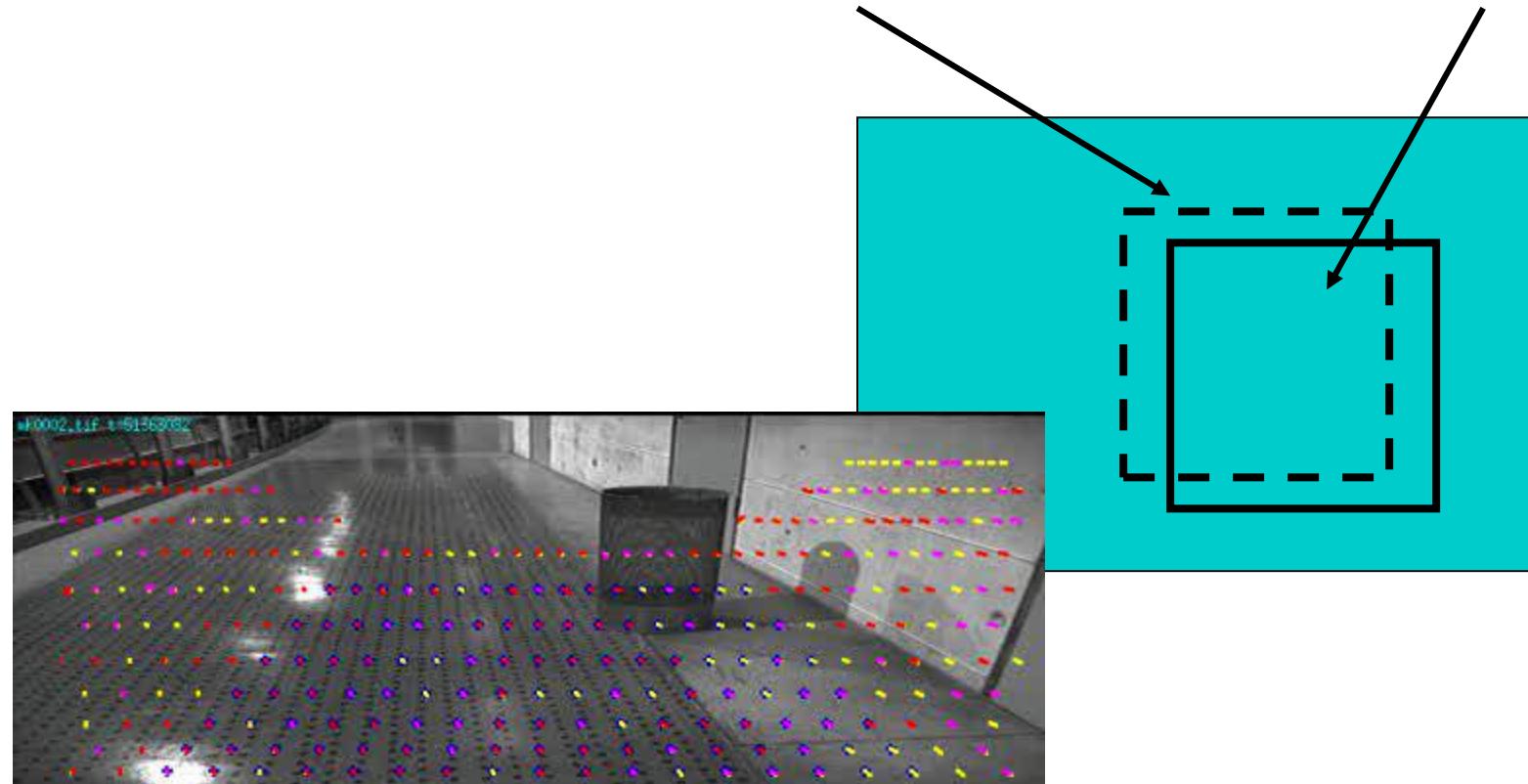
How do we make use of this?

# Fundamental types of video processing

## “Visual motion detection”

- Relating two adjacent frames: (small differences):

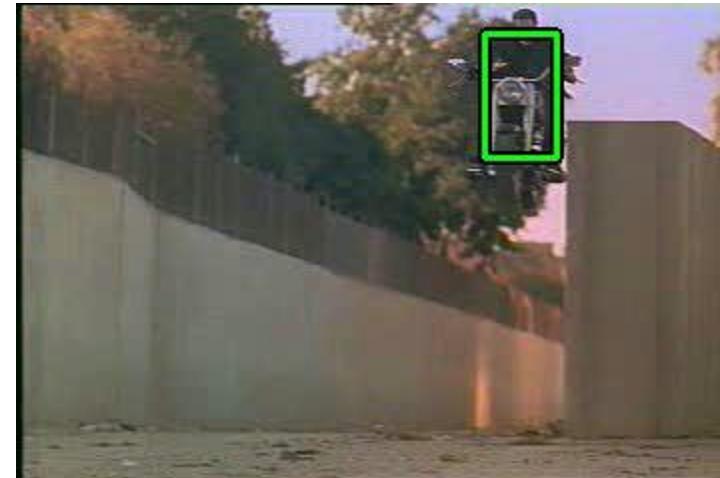
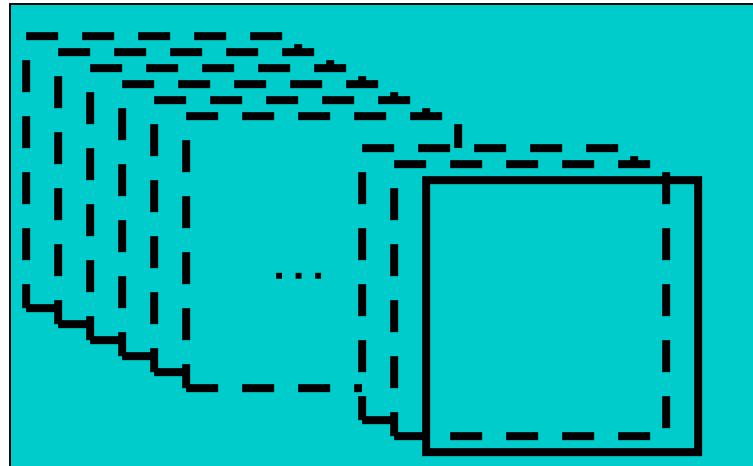
$$\text{Im}(x + \delta x, \ y + \delta y, \ t + \delta t) = \text{Im}(x, \ y, \ t)$$



# Fundamental types of video processing

## “Visual Tracking” / Stabilization

- Globally relating all frames: (large differences):



# Types of tracking

- **Point tracking**

*Extract the point (pixel location) of a particular image feature in each image.*

- **Segmentation based tracking**

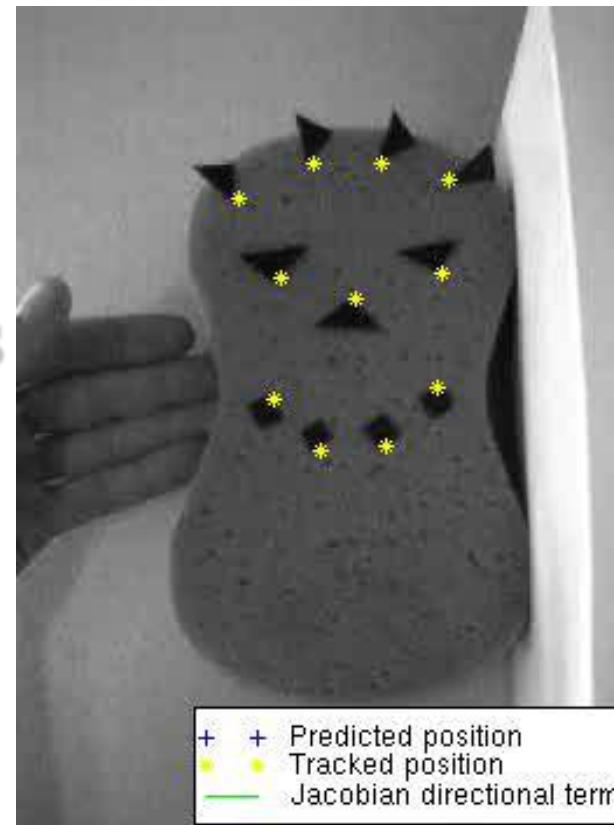
*Define an identifying region property (e.g. color, texture statistics). Repeatedly extract this region in each image.*

- **Stabilization based tracking**

*Formulate image geometry and appearance equations and use these acquire image transform parameters for each image.*

# Point tracking

- Simplest technique. Already commercialized in motion capture systems.
- Features commonly LED's (visible or IR), special markers (reflective, patterns)
- Detection: e.g. pick brightest pixel(s), cross correlate image to find best match for known pattern.



# Commercial point trackers: “Motion capture”



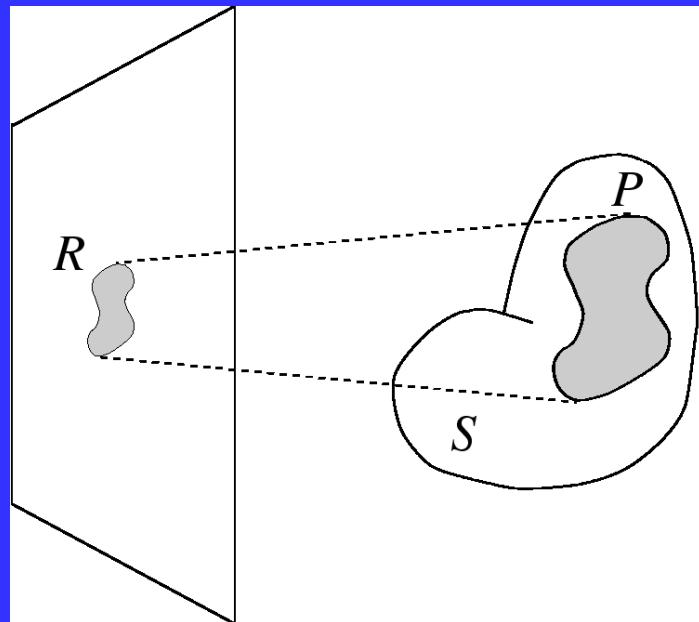
# Region Tracking (Segmentation-Based)

- Select a “cue:”  $\gamma(\mathbf{I}(x, y))$ 
  - foreground enhancement
  - background subtraction
- Segment
  - threshold
  - “clean up”
- Compute region geometry
  - centroid (first moment)
  - orientation (second moment)
  - scale (second moment)

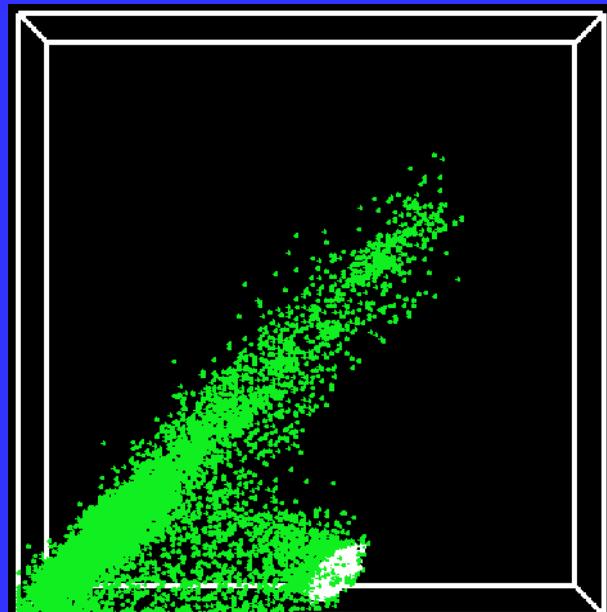
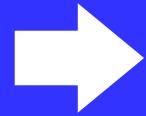
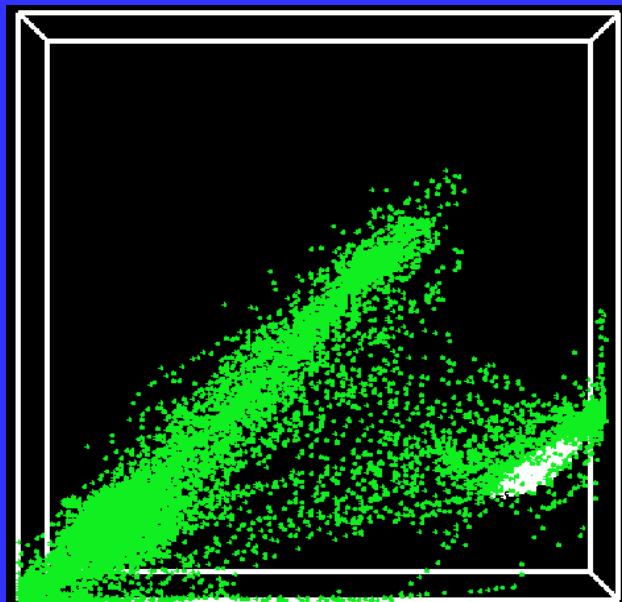
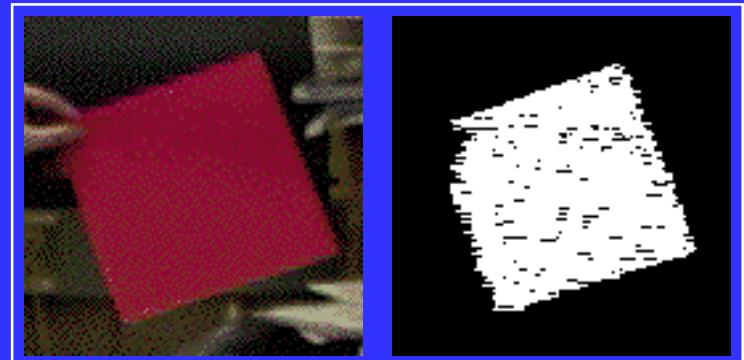
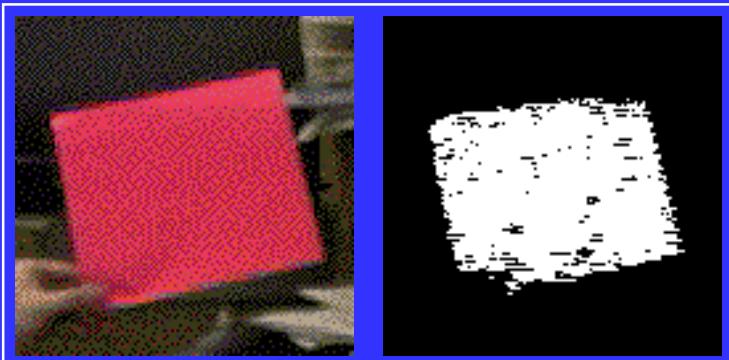
$$\begin{aligned} u &= (x, y)' \\ m_i &= \sum_i S(u) u^i \\ c &= m_1 / m_0 \\ \Lambda &= m_2 / m_0 - m_1^2 \end{aligned}$$

# Regions

- $c_P = P \rightarrow \Re^3$  (color space): intrinsic coloration
  - Homogeneous region:  $c_P(P)$  is roughly constant
  - Textured region:  $c_P(P)$  has significant intensity gradients horizontally and vertically
  - Contour: (local) rapid change in contrast



# Homogeneous Color Region: Photometry

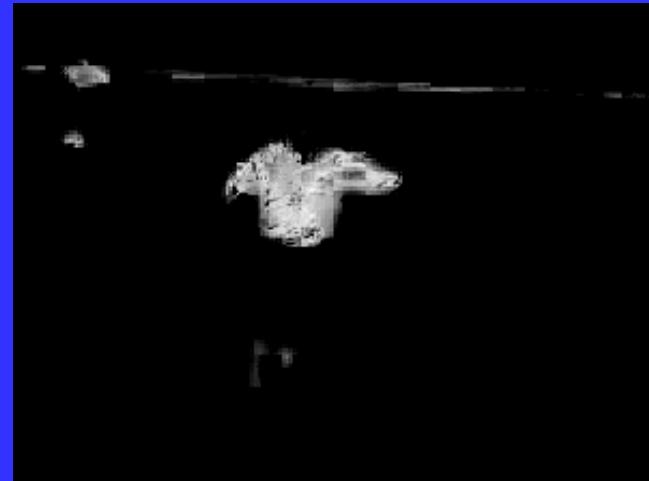


# Color Representation

- $c_R = R \rightarrow \Re^3$  is irradiance of region  $R$
- Color representation
  - DRM [Klinker et al., 1990]: if  $P$  is Lambertian, has *matte* line and *highlight* line
  - User selects *matte* pixels in  $R$
  - PCA fits ellipsoid  $(S, \mathbf{R}^T, T)$  to matte cluster
  - Color similarity  $\gamma(I(x, y))$  is defined by Mahalanobis distance

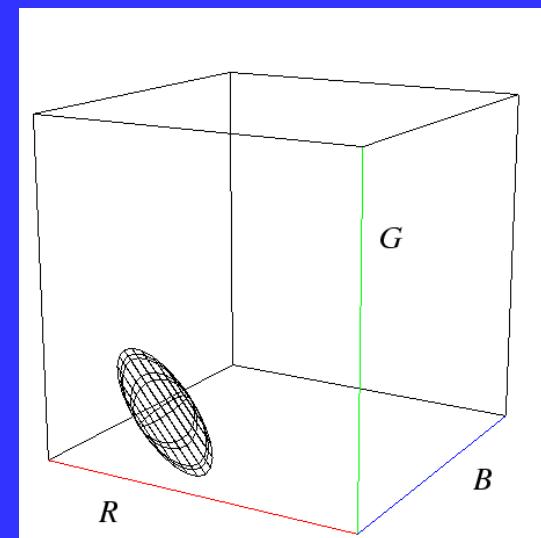
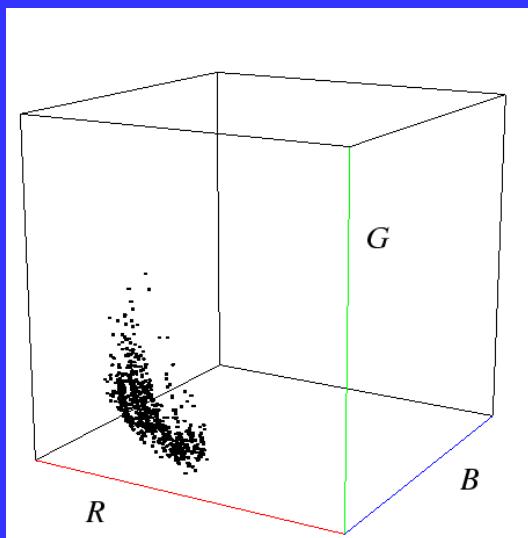
$$| S^{-1} \mathbf{R}^T (I(x, y) - T) | < 1 \text{ inside}$$

# Homogeneous Region: Photometry



Sample

1/22/2022



PCA-fitted  
ellipsoid

# Color Extension (Contd.)

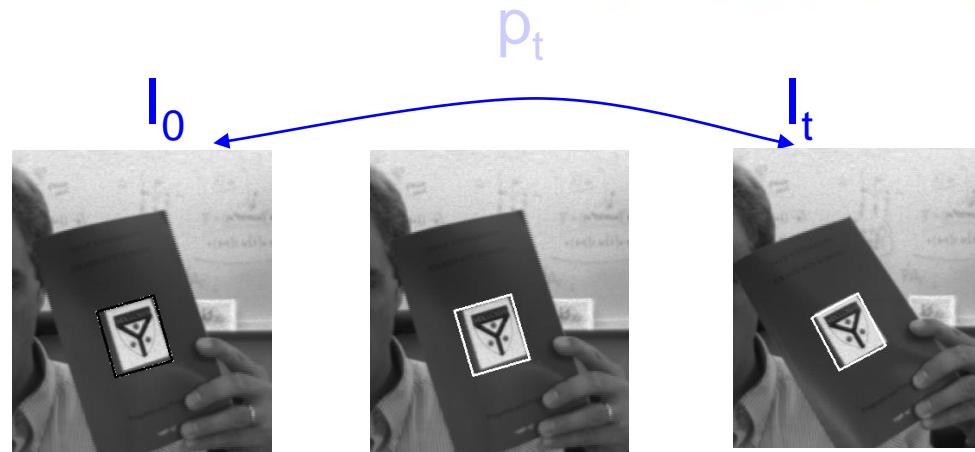
## Color Histograms (Swain et al.)

- $H_i$  = number of pixels in class  $i$
- Histograms are vectors of bin counts
- Given histograms  $H$  and  $G$ , compare by

$$H \cdot G = \frac{\sum H_i G_i}{\sqrt{\sum H_i^2} \sqrt{\sum G_i^2}}$$

- dense, stable signature
- relies on segmentation
- relative color and feature locations lost
- affine transformations preserve area ratios of planar objects

# Intro to Registration Tracking



Variability model:  $I_t = g(I_0, p_t)$

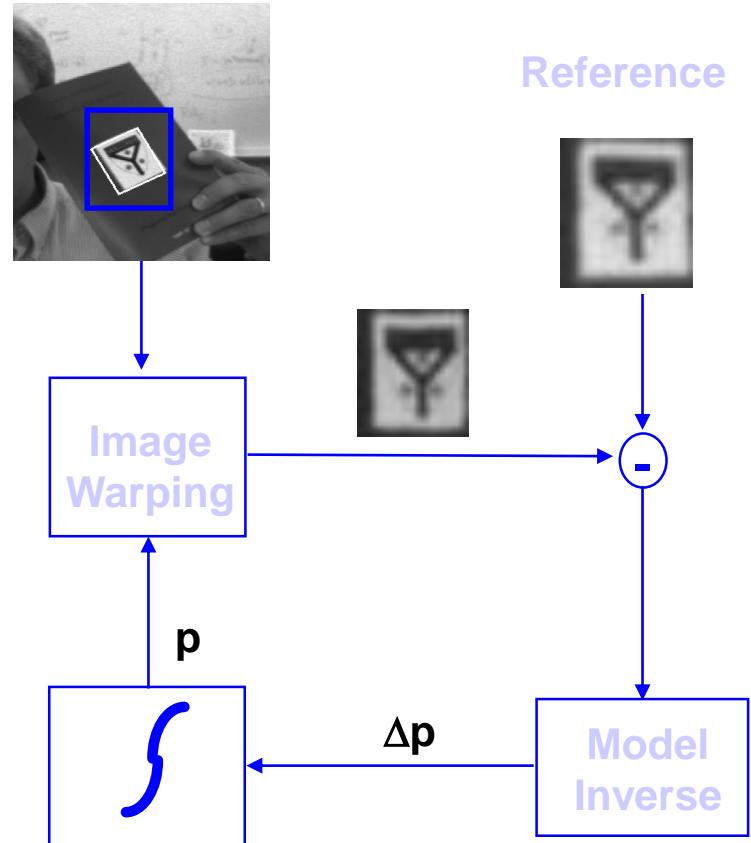
Incremental Estimation: From  $I_0$ ,  $I_{t+1}$  and  $p_t$  compute  $\Delta p_{t+1}$

$$\| I_0 - g(I_{t+1}, p_{t+1}) \|^2 \Rightarrow \min$$

Visual Tracking = Visual Stabilization

# Tracking Cycle

- Prediction
  - Prior states predict new appearance
- Image warping
  - Generate a “normalized view”
- Model inverse
  - Compute error from nominal
- State integration
  - Apply correction to state



# Stabilization tracking: Planar Case

Planar Object +

Linear camera => Affine motion model:

$$u'_i = A u_i + d$$

Subtract these:



= nonsense

Warping



But these:

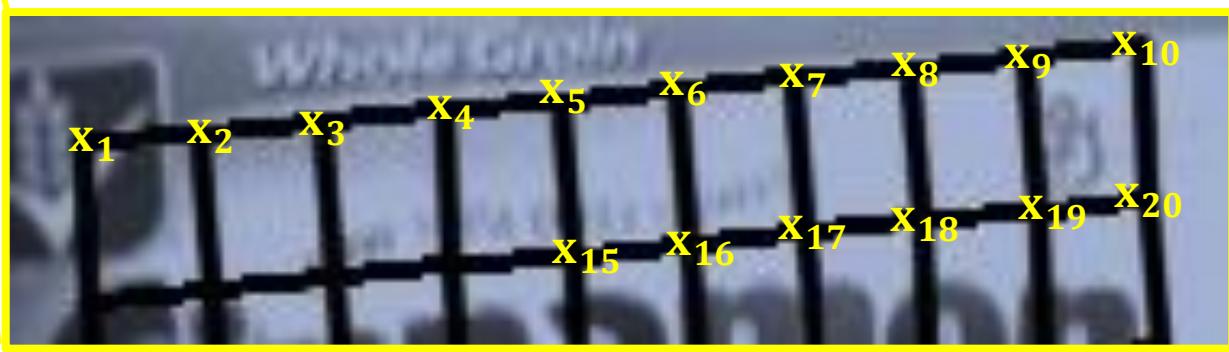
= small variation

$$I_t = g(p_t, I_0)$$

# Registration based Tracking

$$\mathbf{p}_t = \underset{\mathbf{p}}{\operatorname{argmax}} f(\mathbf{I}_0(\mathbf{x}), \mathbf{I}_t(\mathbf{w}(\mathbf{x}, \mathbf{p})))$$


 $I_0$ 

 $I_t$ 


# Motivation

- Learning/detection based trackers are not suitable for tasks requiring **fast** and **high precision** tracking
  - Visual Servoing
  - Virtual reality
  - SLAM

Registration (8DOF)



Learning (3DOF)

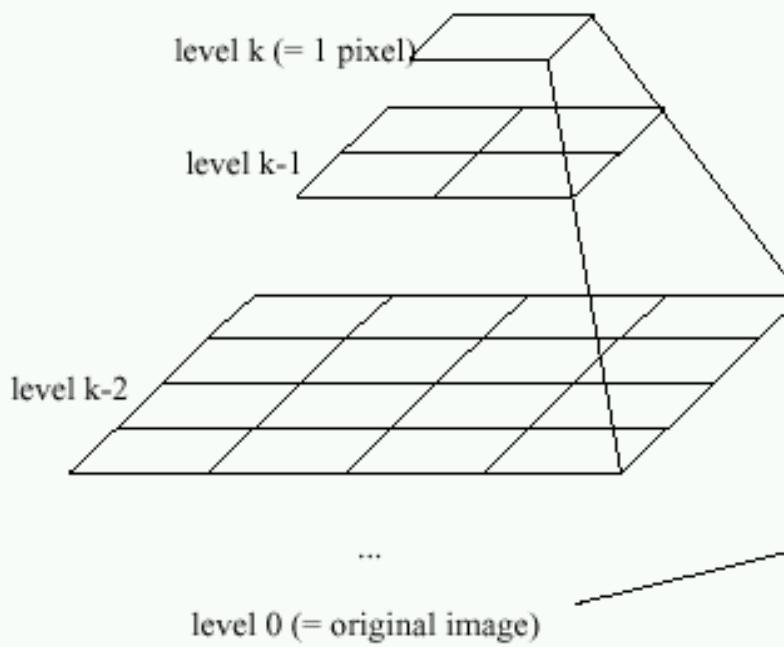


A photograph of a snowy mountain landscape. In the foreground, a person wearing a dark jacket and light-colored pants is walking towards the camera on a snow-covered slope. To the right, there are large, snow-laden evergreen trees. In the background, a range of majestic mountains with rugged peaks and patches of snow under a clear blue sky.

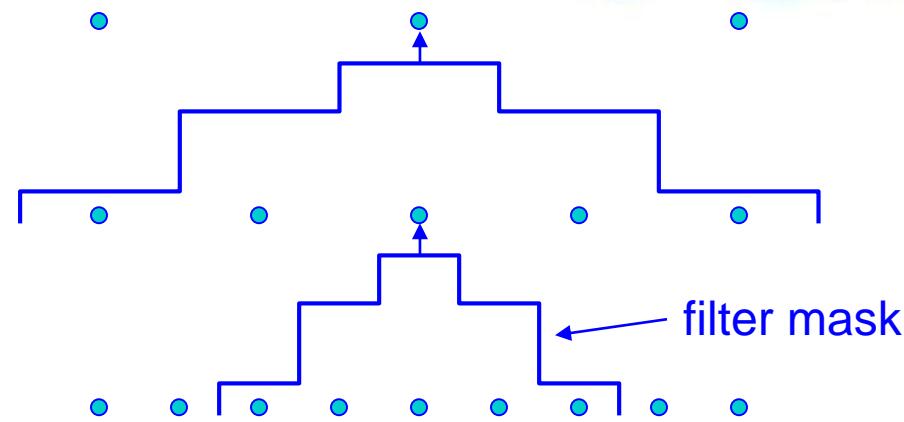
Image Pyramids

# Image Pyramids

Idea: Represent NxN image as a “pyramid” of  
1x1, 2x2, 4x4,...,  $2^k \times 2^k$  images (assuming N=2<sup>k</sup>)



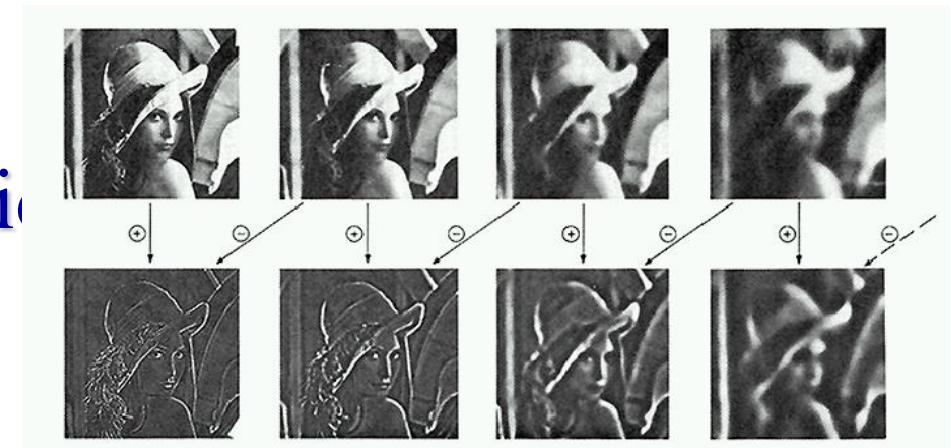
# Pyramid Creation



“Gaussian” Pyramid

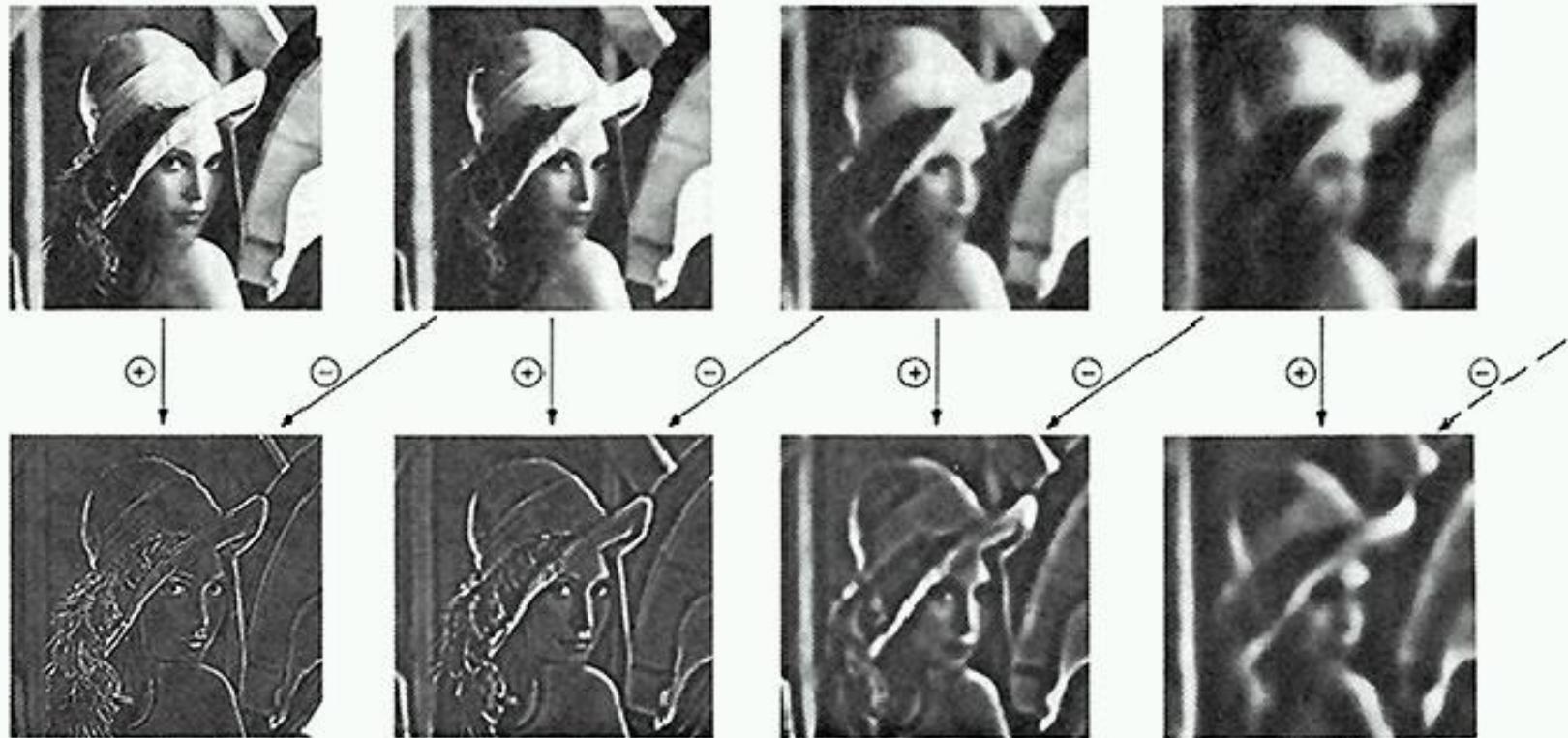
• “Laplacian” Pyramid

- Created from Gaussian pyramid by subtraction  
$$L_l = G_l - \text{expand}(G_{l+1})$$



# Octaves in the Spatial Domain

## Lowpass Images



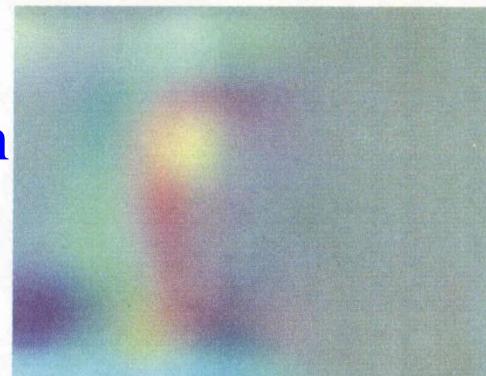
- Bandpass Images

# Pyramids

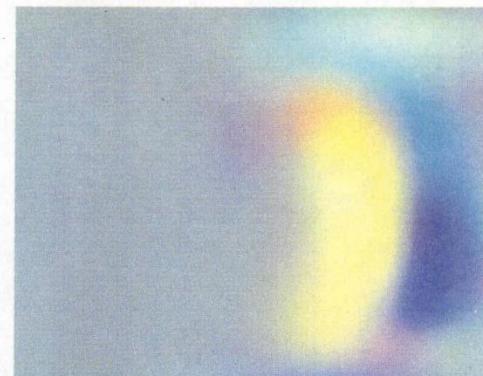
- Advantages of pyramids
  - Faster than Fourier transform
  - Avoids “ringing” artifacts
- Many applications
  - small images faster to process
  - good for multiresolution processing
  - compression
  - progressive transmission
- Known as “mip-maps” in graphics community
- Precursor to wavelets
  - Wavelets also have these advantages

# Mb/2005 estimation

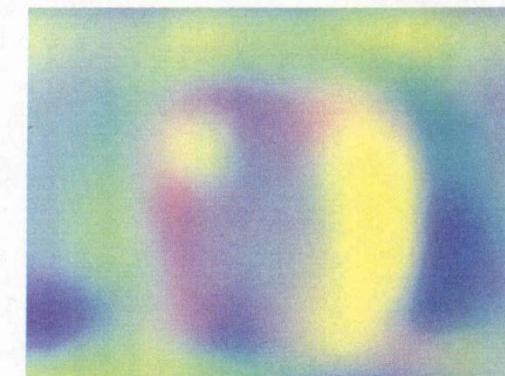
Laplacian  
level  
4



(c)



(g)

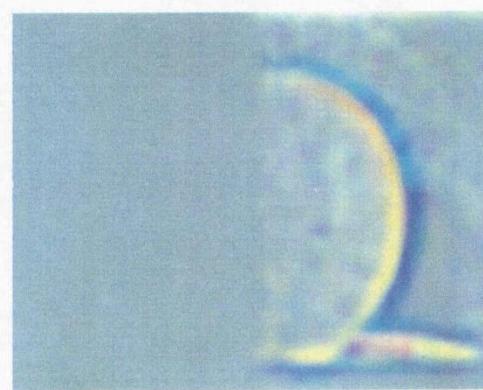


(k)

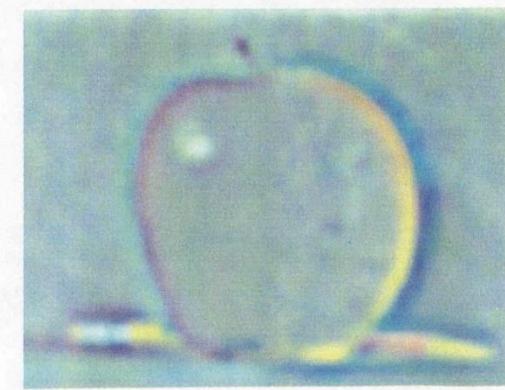
Laplacian  
level  
2



(b)

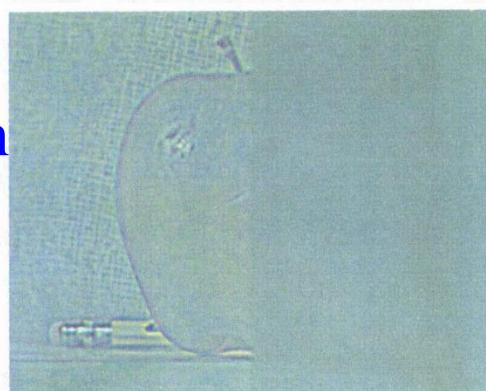


(f)

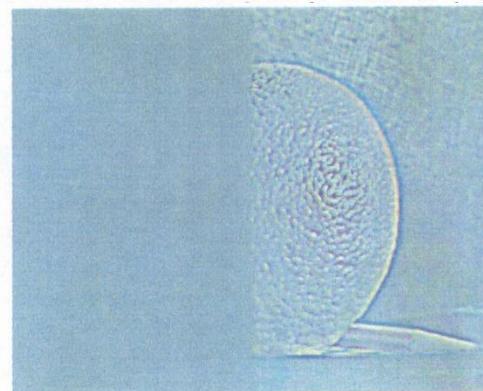


(j)

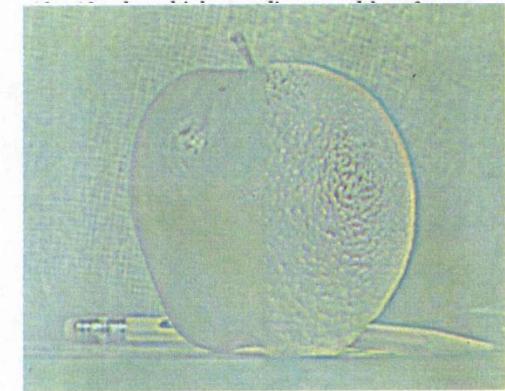
Laplacian  
level  
0



(a)



(e)



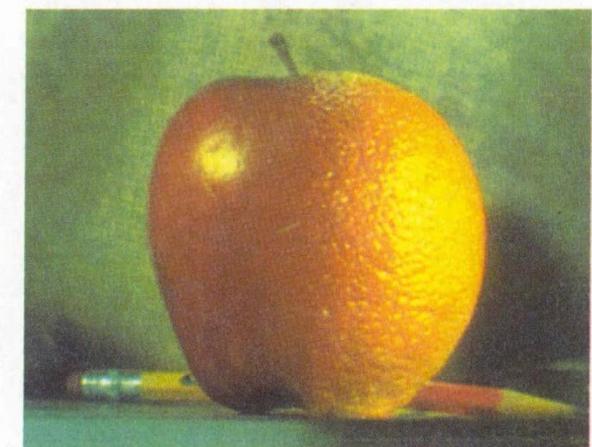
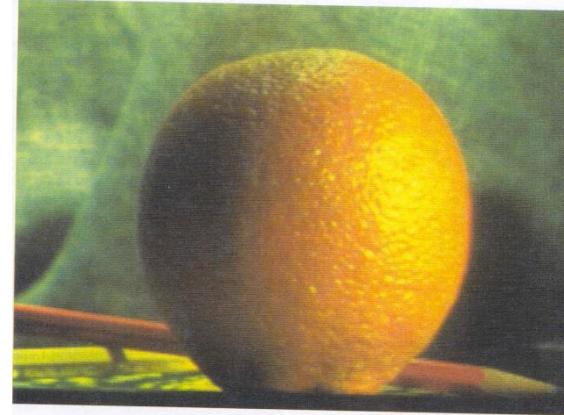
(i)

left pyramid

right pyramid

blended pyramid

# Pyramid Blending



A photograph of a snowy mountain landscape. In the foreground, a person wearing a dark jacket and light-colored pants is walking towards the camera on a snow-covered slope. To the right, there are large, snow-laden evergreen trees. In the background, a range of majestic mountains with snow-capped peaks stretches across the horizon under a clear blue sky.

Image Warping

# Parametric (global) warping

- Examples of parametric warps:



translation



rotation



aspect



affine



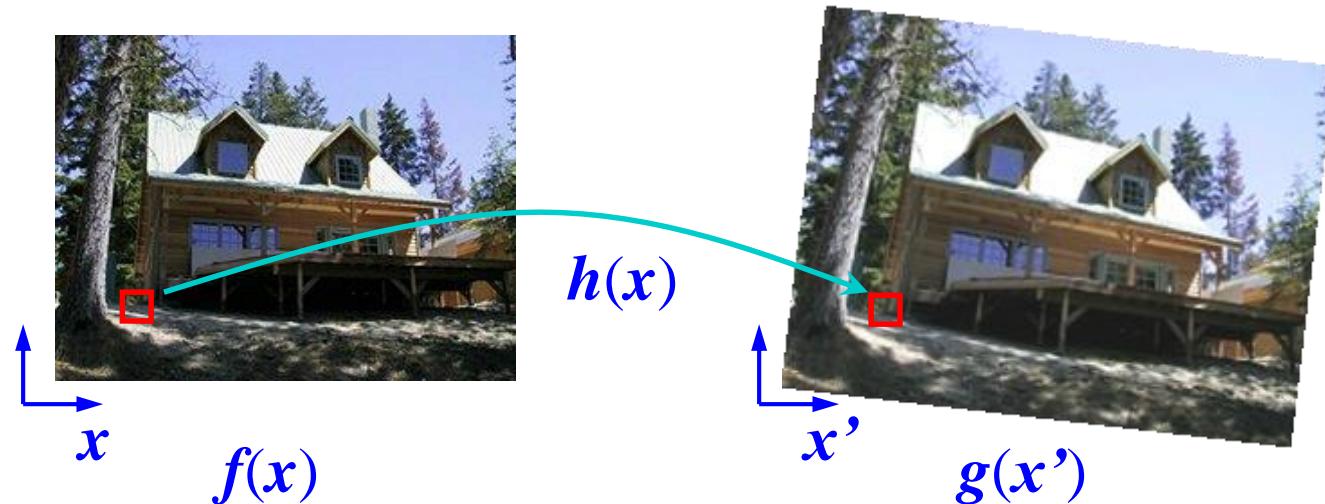
perspective



cylindrical

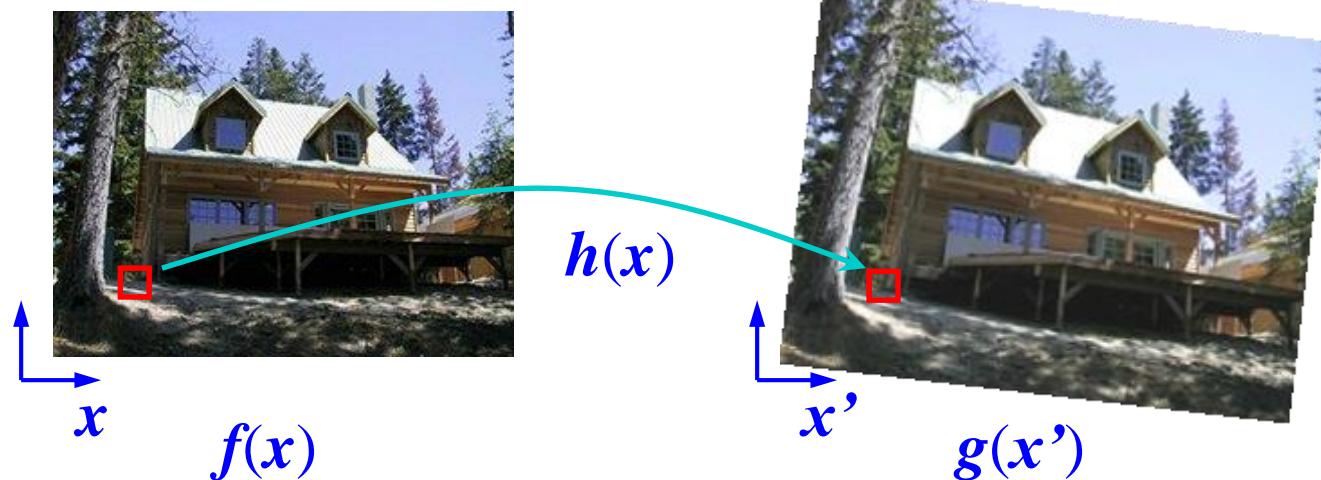
# Image Warping

- Given a coordinate transform  $\mathbf{x}' = \mathbf{h}(\mathbf{x})$  and a source image  $f(\mathbf{x})$ , how do we compute a transformed image  $g(\mathbf{x}') = f(\mathbf{h}(\mathbf{x}))$ ?



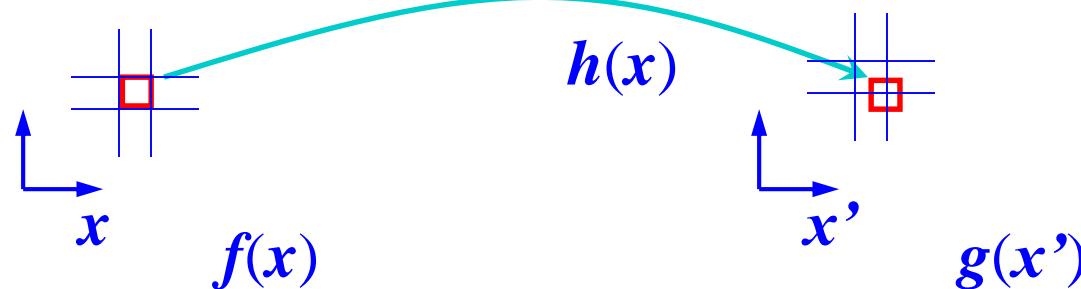
# Forward Warping

- Send each pixel  $f(x)$  to its corresponding location  $x' = h(x)$  in  $g(x')$ 
  - What if pixel lands “between” two pixels?



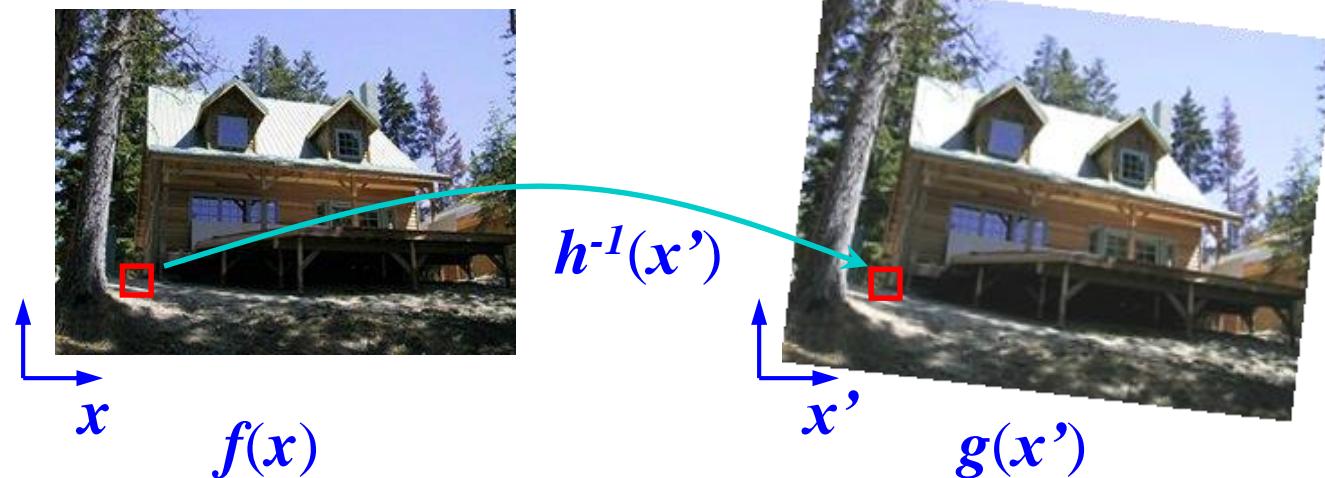
# Forward Warping

- Send each pixel  $f(x)$  to its corresponding location  $x' = h(x)$  in  $g(x')$ 
  - What if pixel lands “between” two pixels?
  - Answer: add “contribution” to several pixels, normalize later (*splatting*)



# Inverse Warping

- Get each pixel  $g(x')$  from its corresponding location  $x = h^{-1}(x')$  in  $f(x)$ 
  - What if pixel comes from “between” two pixels?



# Inverse Warping

- Get each pixel  $g(x')$  from its corresponding location  $x = h^{-1}(x')$  in  $f(x)$ 
  - What if pixel comes from “between” two pixels?
  - Answer: *resample* color value from *interpolated (prefiltered)* source image



# Interpolation

- Possible interpolation filters:

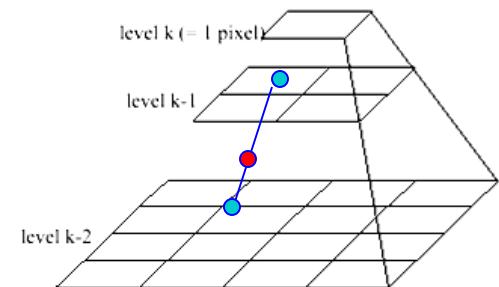
- nearest neighbor
- bilinear
- bicubic (interpolating)
- sinc / FIR

- Needed to prevent “jaggies” and “texture crawl” (see demo)



# Prefiltering

- Essential for *downsampling (decimation)* to prevent *aliasing*
- MIP-mapping [Williams'83]:
  1. build pyramid (but what decimation filter?)
    - block averaging
    - Burt & Adelson (5-tap binomial)
    - 7-tap wavelet-based filter (better)
  2. *trilinear* interpolation
    - bilinear within each 2 adjacent levels
    - linear blend *between* levels (determined by pixel size)

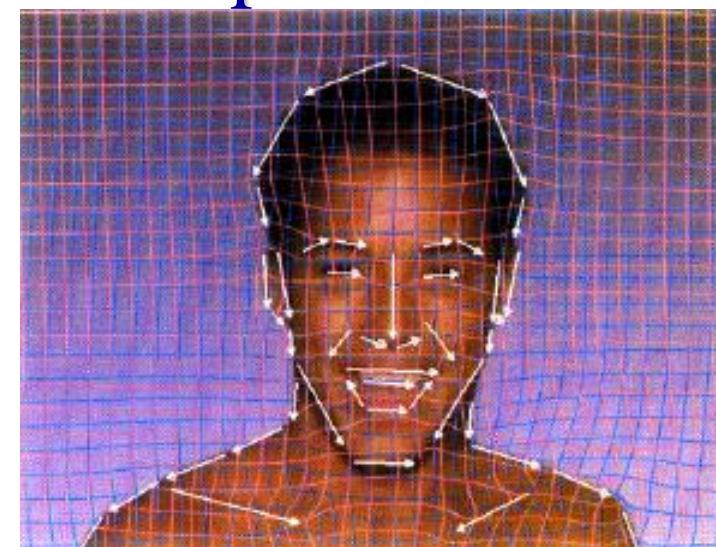


# Prefiltering

- Essential for *downsampling (decimation)* to prevent *aliasing*
- Other possibilities:
  - summed area tables
  - elliptically weighted Gaussians (EWA) [Heckbert'86]

# Image Warping – non-parametric

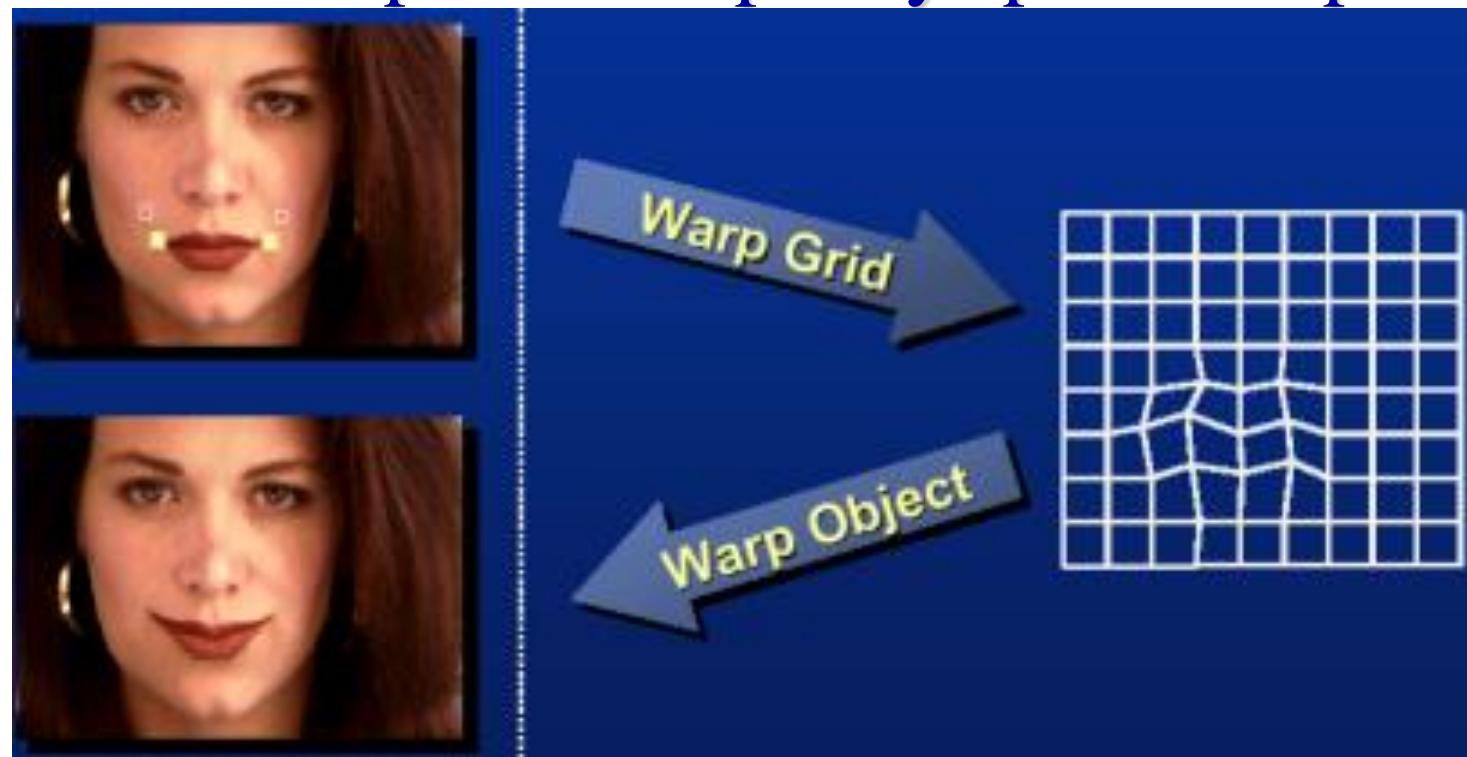
- Specify more detailed warp function



- Examples:
  - splines
  - triangles
  - optical flow (per-pixel motion)

# Image Warping – non-parametric

- Move control points to specify spline warp

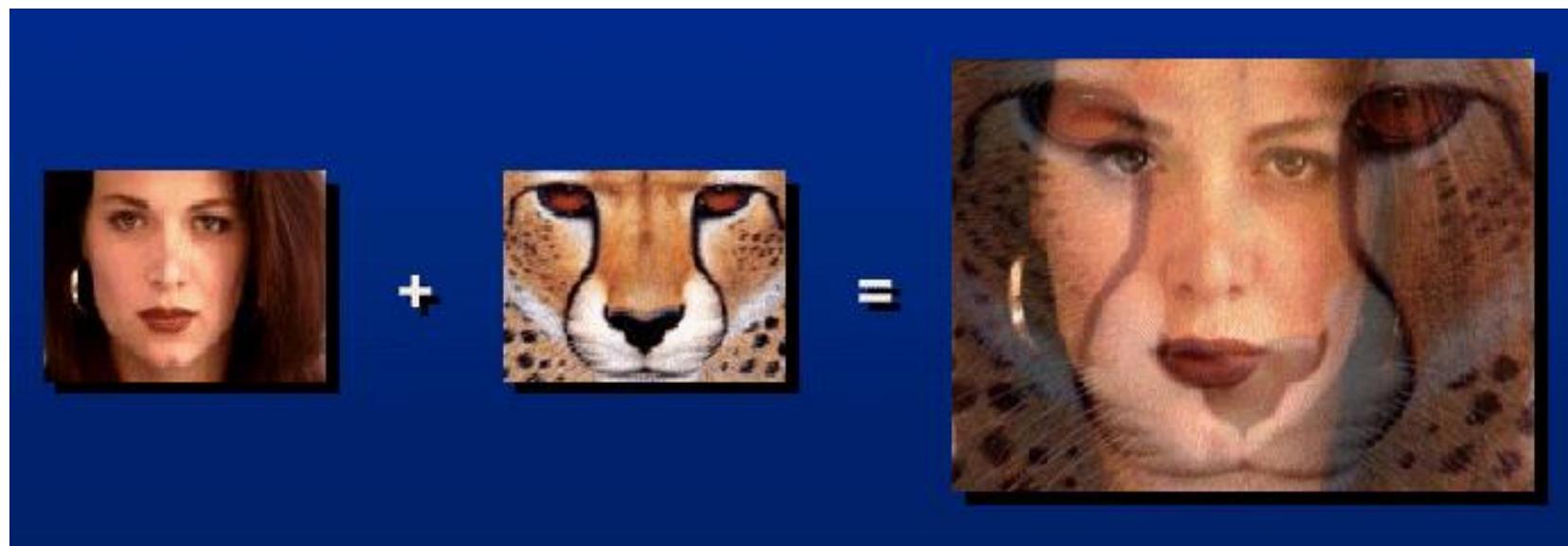


# Image Morphing



# Image Morphing

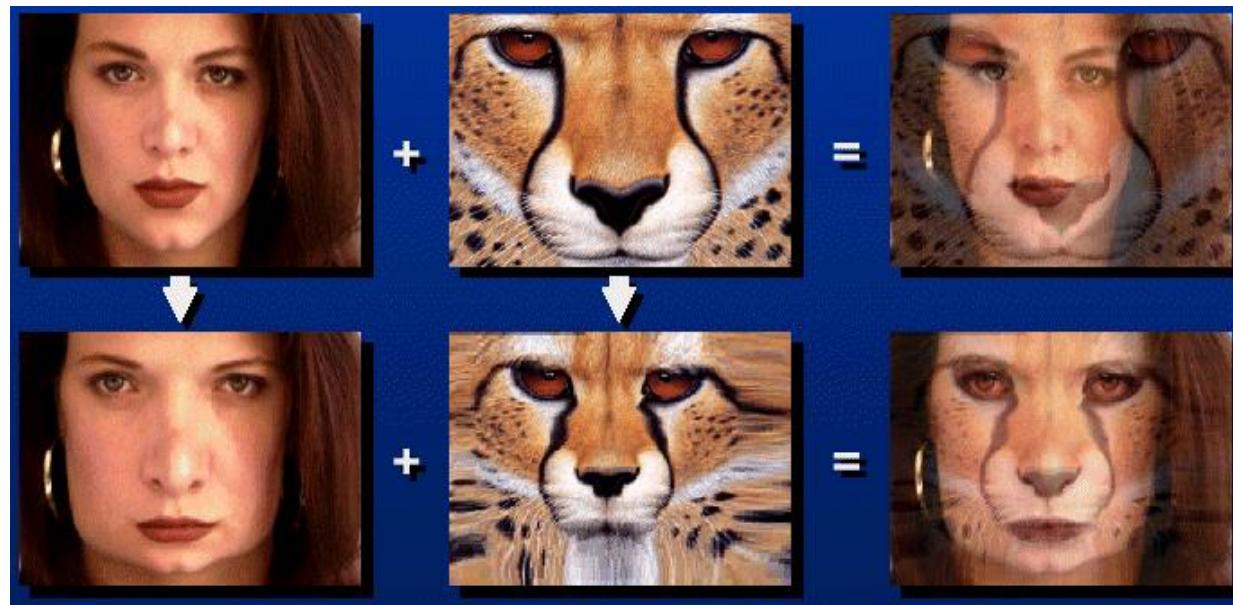
- How can we *in-between* two images?
  1. Cross-dissolve



(all examples from [Gomes *et al.* '99])

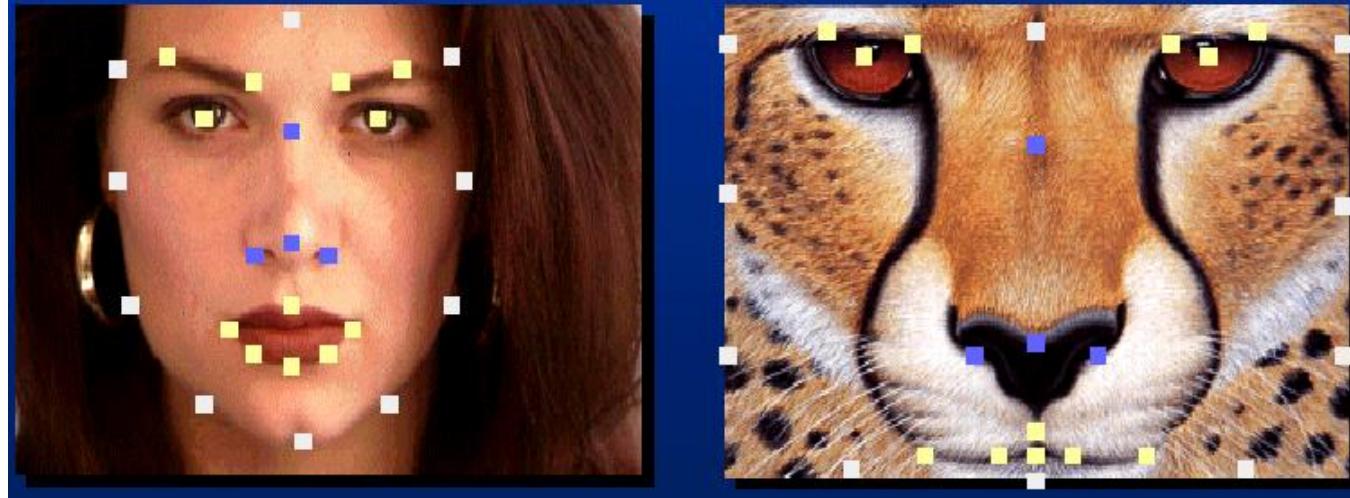
# Image Morphing

- How can we *in-between* two images?
  2. Warp then cross-dissolve = *morph*



# Warp specification

- How can we specify the warp?
  1. Specify corresponding *points*
    - *interpolate* to a complete warping function



- Nielson, *Scattered Data Modeling*, IEEE CG&A'93]

# Warp specification

- How can we specify the warp?
  2. Specify corresponding *vectors*
    - *interpolate* to a complete warping function

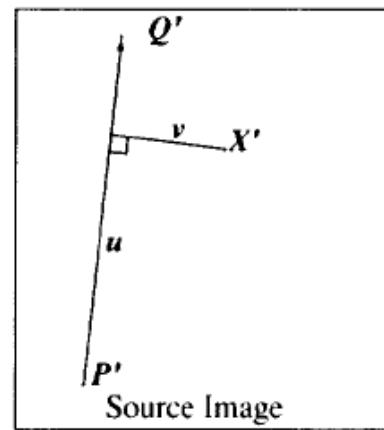
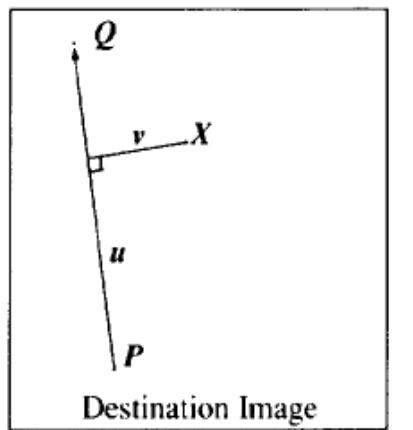


# Warp specification

- How can we specify the warp?

## 2. Specify corresponding *vectors*

- *interpolate* [Beier & Neely, SIGGRAPH'92]



For each pixel  $X$  in the destination

$$DSUM = (0,0)$$

$$weightsum = 0$$

For each line  $P_i Q_i$

calculate  $u, v$  based on  $P_i Q_i$

calculate  $X'_i$  based on  $u, v$  and  $P_i Q_i'$

calculate displacement  $D_i = X'_i - X_i$  for this line

$dist =$  shortest distance from  $X$  to  $P_i Q_i$

$$weight = (length^p / (a + dist))^b$$

$$DSUM += D_i * weight$$

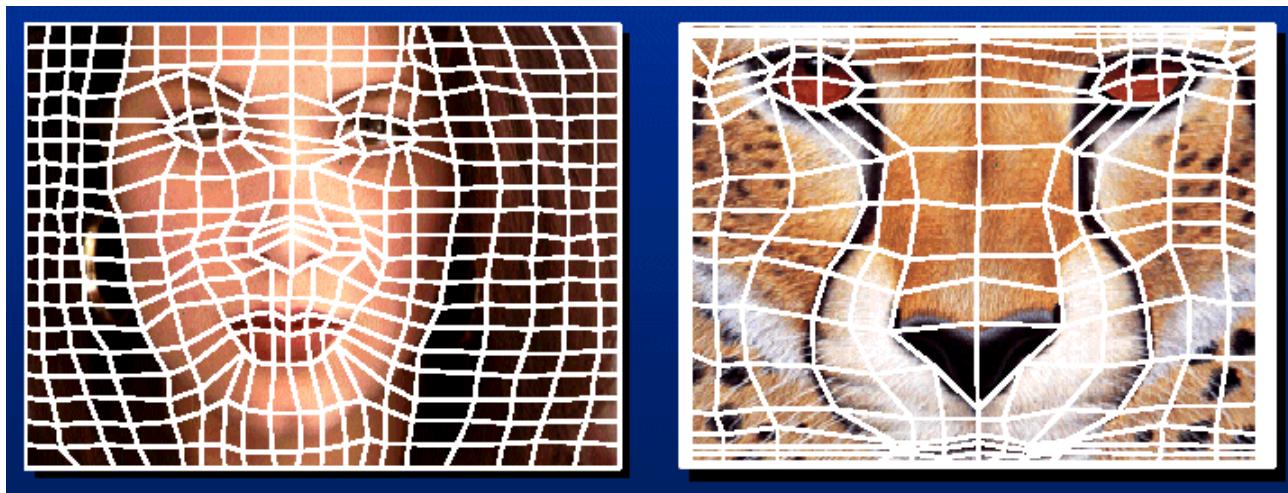
$$weightsum += weight$$

$$X' = X + DSUM / weightsum$$

$$\text{destinationImage}(X) = \text{sourceImage}(X')$$

# Warp specification

- How can we specify the warp?
  3. Specify corresponding *spline control points*
    - *interpolate* to a complete warping function

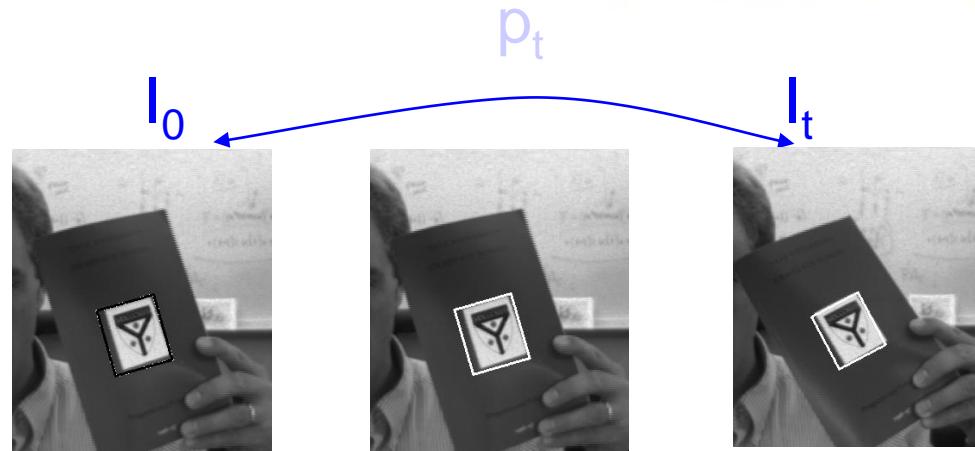


Mb/2005 estimation

# Final Morph Result



# Intro to Registration Tracking



Variability model:  $I_t = g(I_0, p_t)$

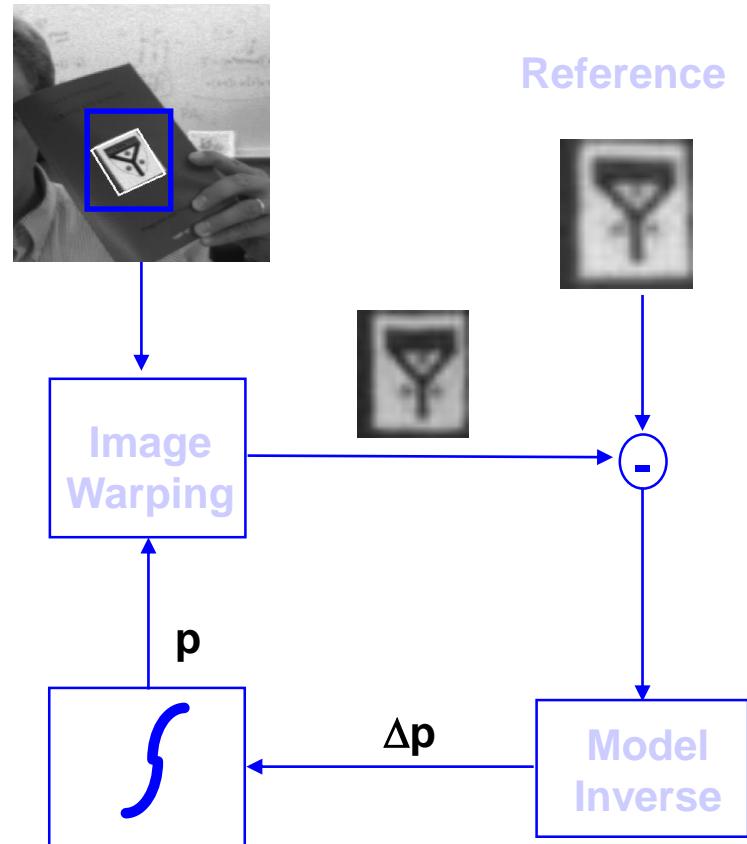
Incremental Estimation: From  $I_0$ ,  $I_{t+1}$  and  $p_t$  compute  $\Delta p_{t+1}$

$$\| I_0 - g(I_{t+1}, p_{t+1}) \|^2 \Rightarrow \min$$

Visual Tracking = Visual Stabilization

# Tracking Cycle

- Prediction
  - Prior states predict new appearance
- Image warping
  - Generate a “normalized view”
- Model inverse
  - Compute error from nominal
- State integration
  - Apply correction to state



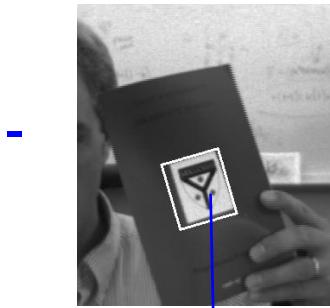
# Stabilization tracking: Planar Case

Planar Object +

Linear camera => Affine motion model:

$$u'_i = A u_i + d$$

Subtract these:



= nonsense

Warping



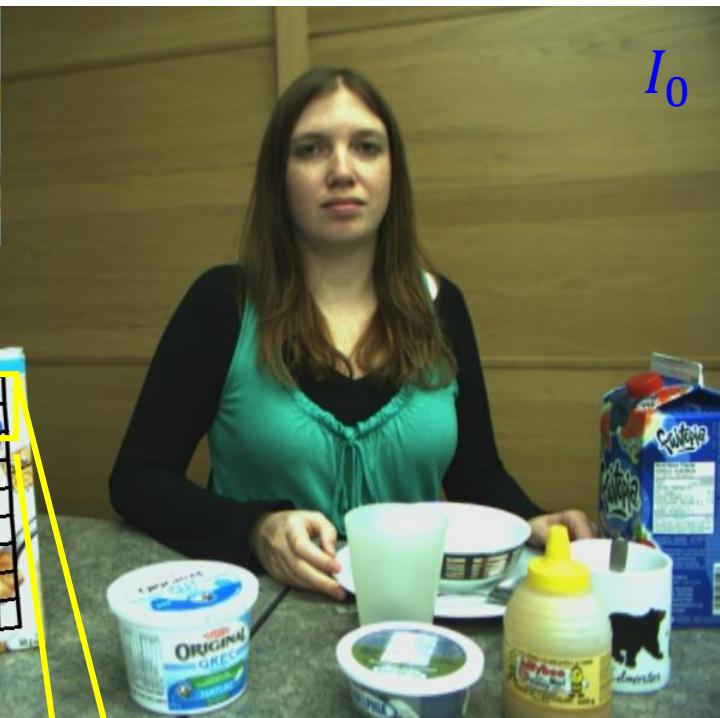
But these:

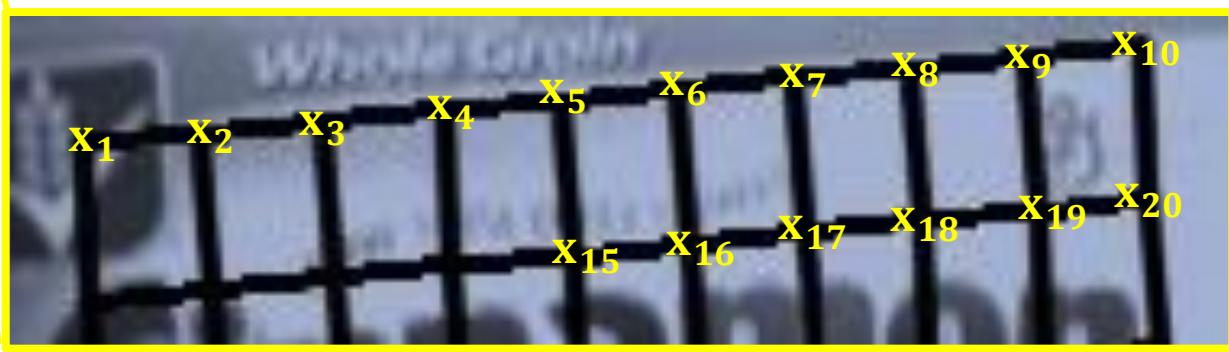
= small variation

$$I_t = g(p_t, I_0)$$

# Registration based Tracking

$$\mathbf{p}_t = \operatorname{argmax}_{\mathbf{p}} f(\mathbf{I}_0(\mathbf{x}), \mathbf{I}_t(\mathbf{w}(\mathbf{x}, \mathbf{p})))$$


 $I_0$ 

 $I_t$ 


# Motivation

- Learning/detection based trackers are not suitable for tasks requiring **fast** and **high precision** tracking
  - Visual Servoing
  - Virtual reality
  - SLAM

Registration (8DOF)



Learning (3DOF)



# Evaluation Benchmarks

## TMT



## PAMI



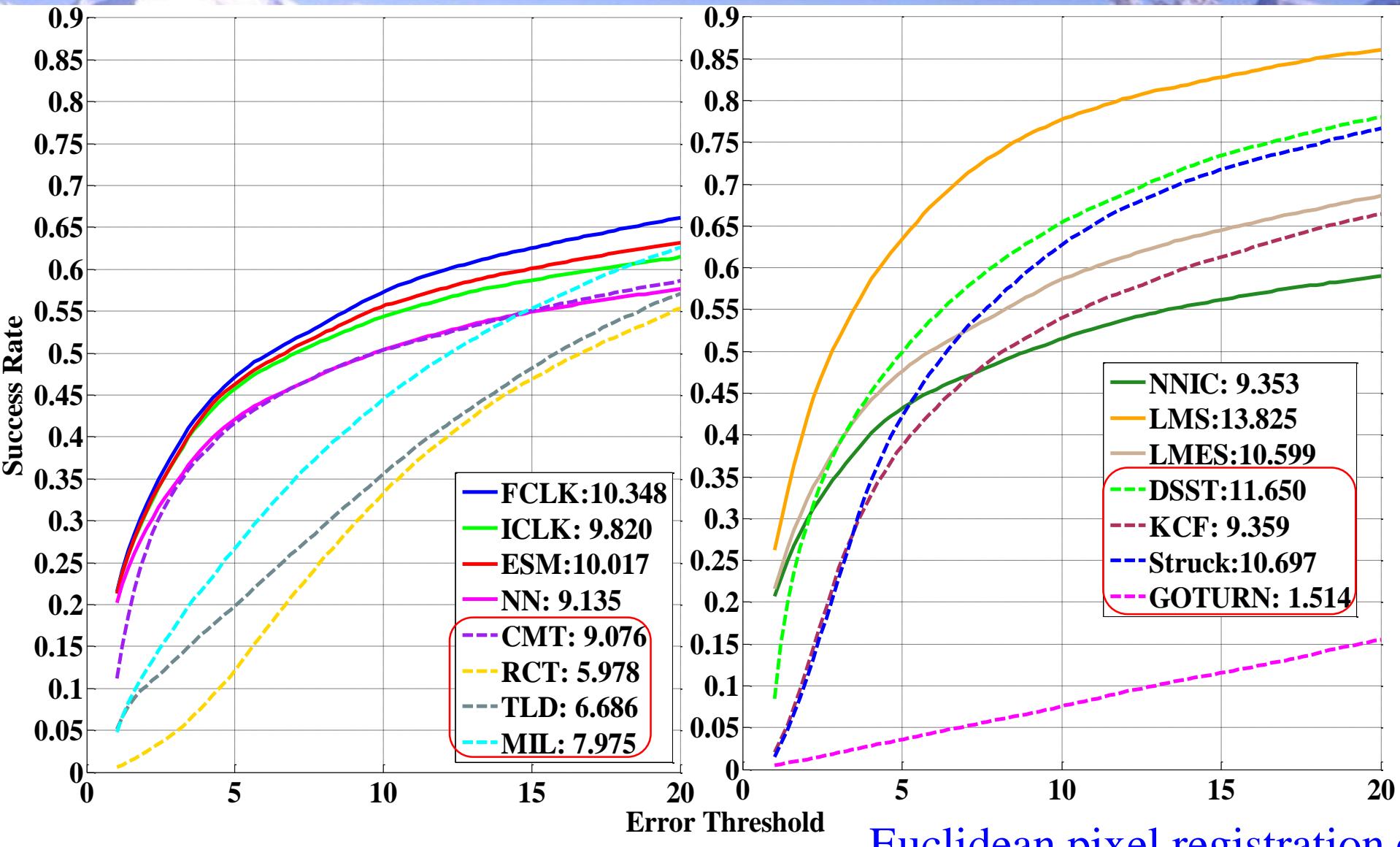
## UCSB



## LinTrack



# 62 Results: Learning vs. 2DOF Registration Based Trackers (Accuracy)



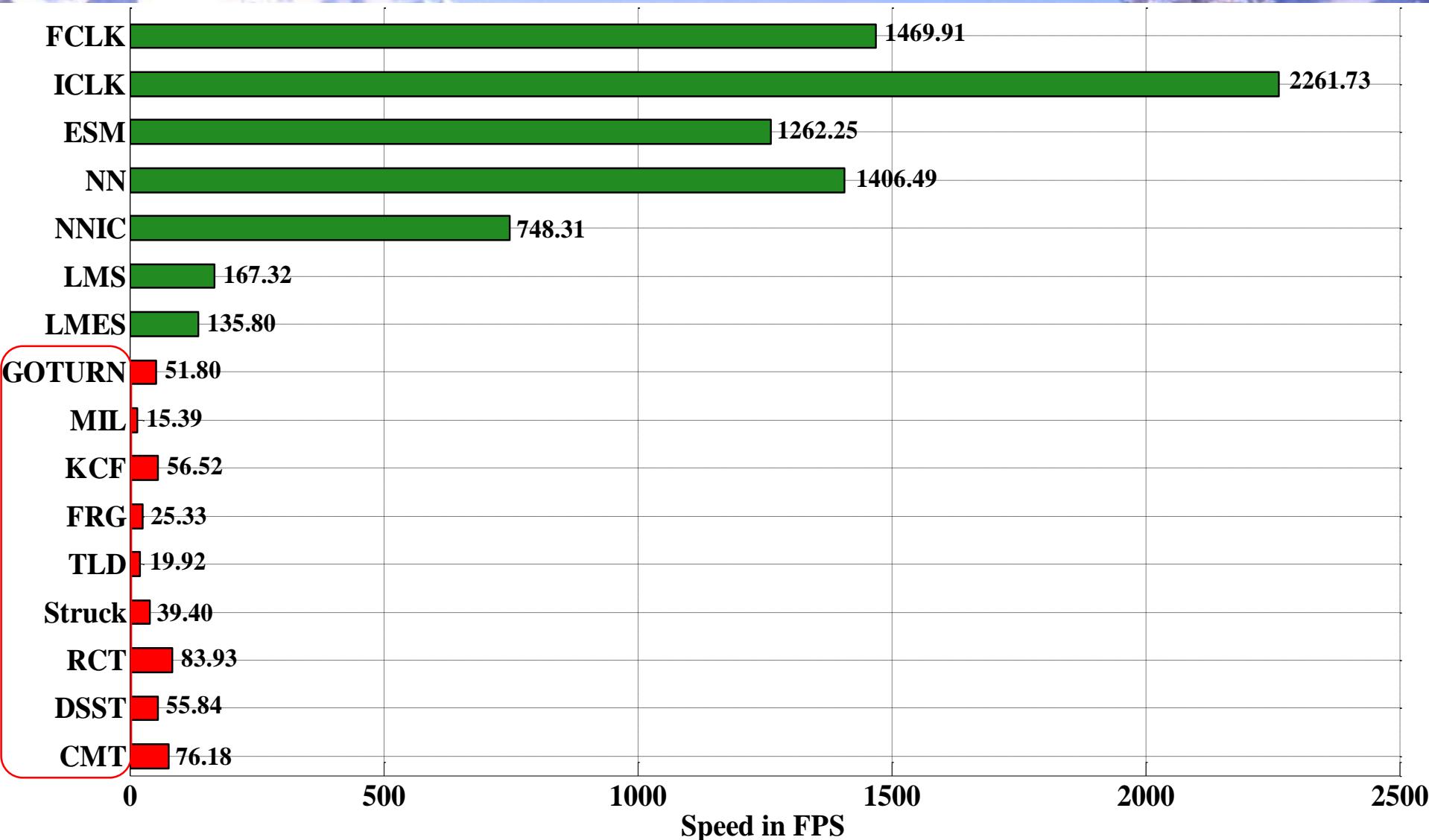
# RBT vs Learn vs GoTurn

nl\_cereal\_s3: frame 1

GOTURN DSST ESM LMES



# 65 Results: Learning vs. 2DOF Registration Based Trackers (Speed)



# References

- S. Baker and I. Matthews, “*Equivalence and Efficiency of Image Alignment Algorithms*”, CVPR 2001
- S. Benhimane and E. Malis, “*Real-time image-based tracking of planes using efficient second-order minimization* ”, IROS 2004
- A. Dame and E. Marchand, “*Accurate Real-time Tracking Using Mutual Information*”, ISMAR 2010
- T. Dick, C. Perez, A. Shademan and M. Jagersand, “*Realtime Registration-Based Tracking via Approximate Nearest Neighbor Search*”, RSS 2013
- S. Gauglitz, T. Hollerer and M. Turk. “*Evaluation of interest point detectors and feature descriptors for visual tracking*”, IJCV 2011
- J. Kwon, H. S. Lee, F. C. Park, and K. M. Lee, “*A Geometric Particle Filter for Template-Based Visual Tracking*”, TPAMI 2014
- R. Richa, R. Sznitman, R. Taylor and G. Hager, “*Visual Tracking Using the Sum of Conditional Variance*”, IROS 2011

# References (cont'd)

- R. Richa, R. Sznitman and G. Hager, “*Robust Similarity Measures for Gradient-based Direct Visual Tracking*”, CIRL Technical Report 2012
- R. Richa, M. Souza, G. Scandaroli, E. Comunello and A. Wangenheim, “*Direct visual tracking under extreme illumination variations using the sum of conditional variance*”, ICIP 2014
- A. Roy, X. Zhang, N. Wolleb, C. P. Quintero and M. Jagersand, “*Tracking Benchmark and Evaluation for Manipulation Tasks*”, ICRA 2015
- L. Ruthotto, “*Mass-preserving registration of medical images*”, Thesis 2010
- G. G. Scandaroli, M. Meilland and R. Richa, “*Improving NCC-Based Direct Visual Tracking*”, ECCV 2012
- A. Singh and M. Jagersand, “*Modular Tracking Framework: A Unified Approach to Registration based Tracking*”, arXiv:1602.09130, 2016