



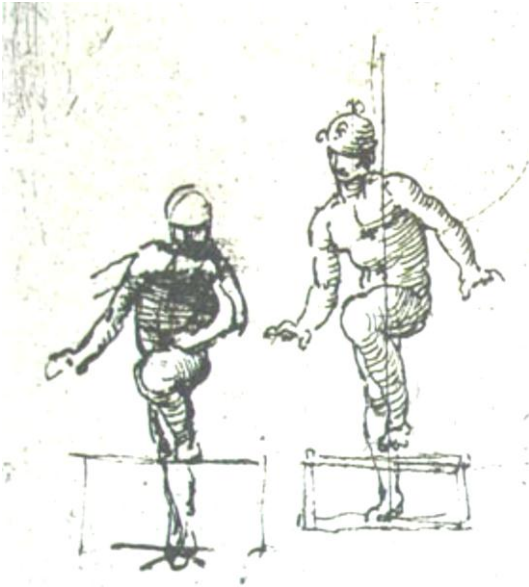
Motion and Human Actions

Ivan Laptev

ivan.laptev@inria.fr

INRIA, WILLOW, ENS/INRIA/CNRS UMR 8548
Laboratoire d'Informatique, Ecole Normale Supérieure, Paris

Class overview



Motivation

- Historic review
- Modern applications

Appearance-based methods

- Motion history images
- Active shape models
- Tracking and motion priors

Motion-based methods

- Generic and parametric Optical Flow
- Motion templates

Space-time methods

- Local space-time features
- Action classification and detection
- Weakly-supervised action learning

Motivation I: Artistic Representation

Early studies were motivated by human representations in Arts

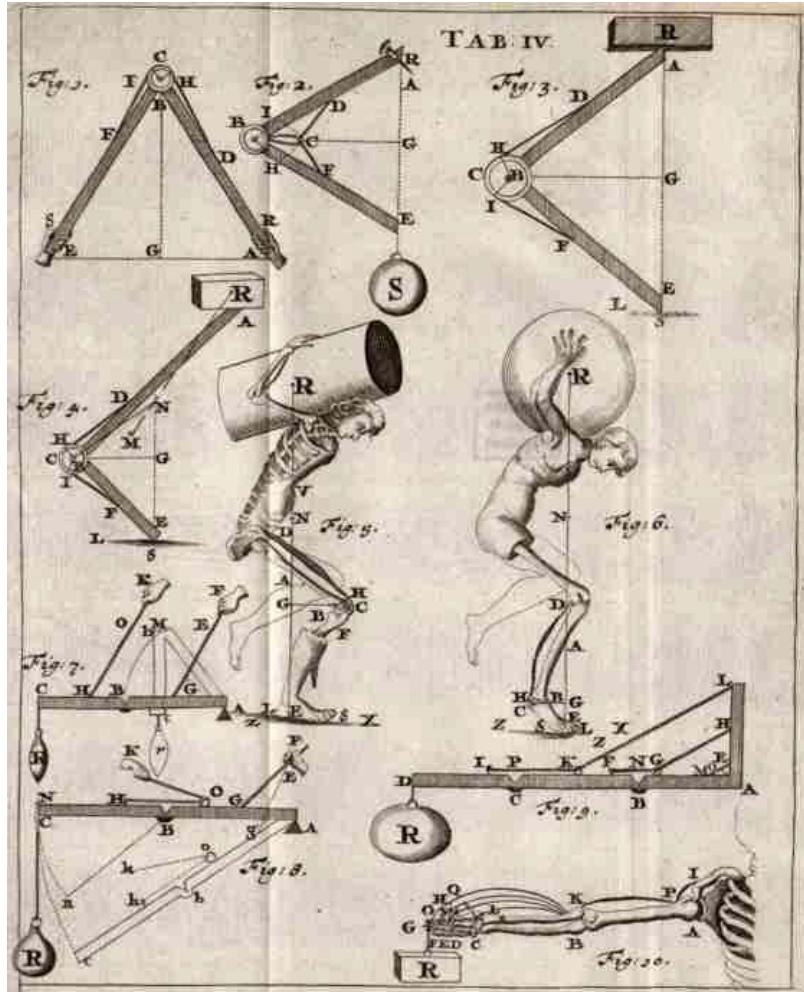
Da Vinci: “it is indispensable for a painter, to become totally familiar with the anatomy of nerves, bones, muscles, and sinews, such that he understands for their various motions and stresses, which sinews or which muscle causes a particular motion”

“I ask for the weight [pressure] of this man for every segment of motion when climbing those stairs, and for the weight he places on *b* and on *c*. Note the vertical line below the center of mass of this man.”



Leonardo da Vinci (1452–1519): A man going upstairs, or up a ladder.

Motivation II: Biomechanics

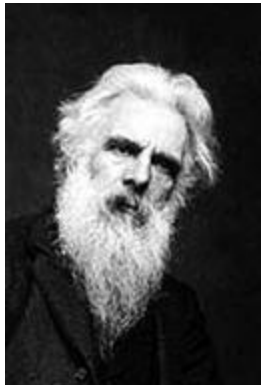
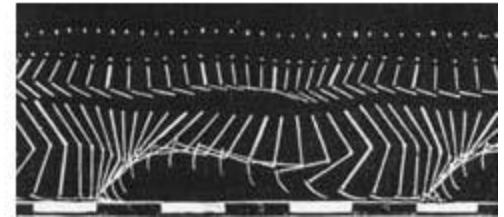
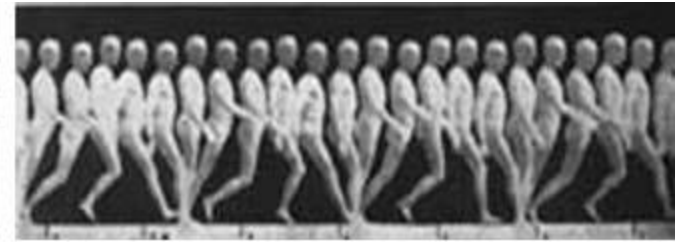


- The emergence of *biomechanics*
- Borelli applied to biology the analytical and geometrical methods, developed by Galileo Galilei
- He was the first to understand that bones serve as levers and muscles function according to mathematical principles
- His physiological studies included muscle analysis and a mathematical discussion of movements, such as running or jumping

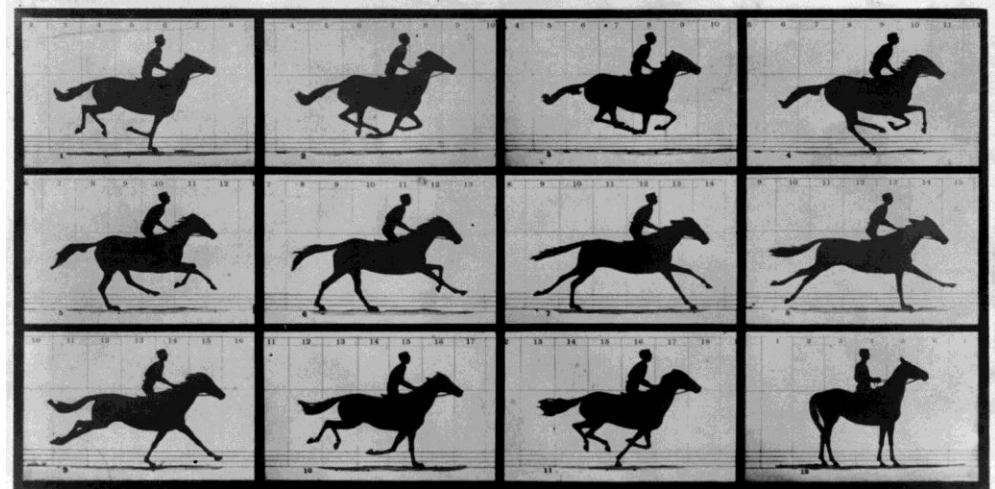
Motivation III: Motion perception



Etienne-Jules Marey:
(1830–1904) made Chronophotographic experiments influential for the emerging field of *cinematography*



Eadweard Muybridge
(1830–1904) invented a machine for displaying the recorded series of images. He pioneered motion pictures and applied his technique to movement studies



Copyright, 1888, by MUYBRIDGE.

THE HORSE IN MOTION.

Illustrated by
MUYBRIDGE.

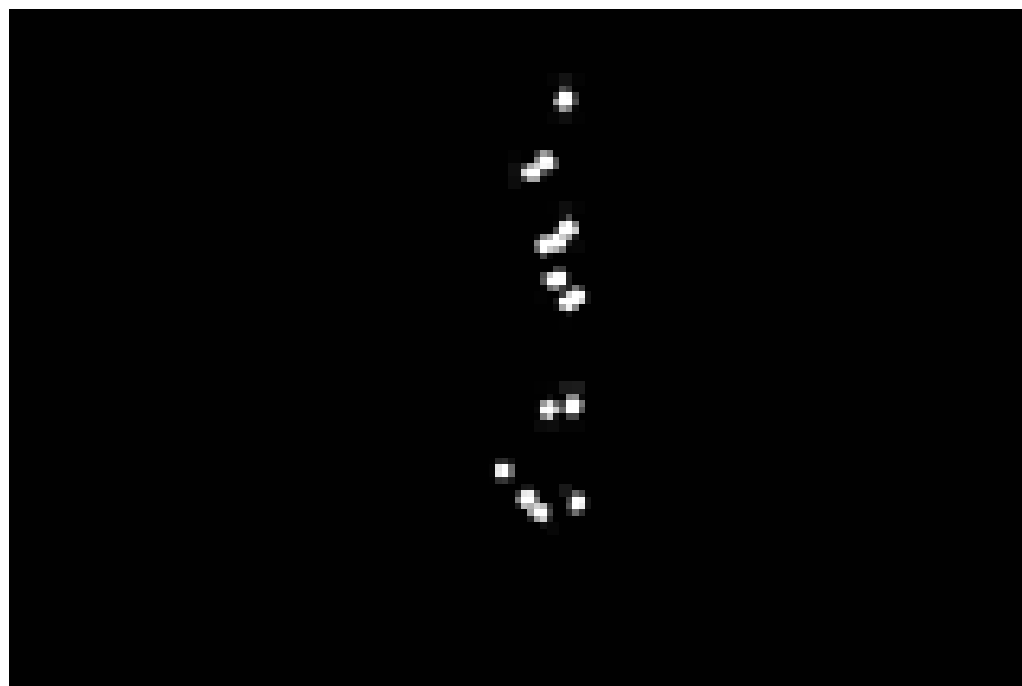
AUTOMATIC ELECTRO-PHOTOGRAPH.

"SALLIE GARDNER," owned by LELAND STANFORD; running at a 1.40 gait over the Palo Alto track, 10th June, 1878.
The negatives of these photographs were made at intervals of twenty-seven inches of distance, and about the twenty-fifth part of a second of time; they illustrate consecutive positions assumed in each twenty-seven inches of progress during a single stride of the horse. The vertical lines were twenty-seven inches apart; the horizontal lines represent elevations of four inches each. The exposure of each negative was less than the two-hundred-and-fiftieth part of a second.

MORSE'S Gallery, 417 Montgomery St., San Francisco.

Motivation III: Motion perception

- Gunnar Johansson [1973] pioneered studies on the use of image sequences for a programmed human motion analysis
- “Moving Light Displays” (LED) enable identification of familiar people and the gender and inspired many works in computer vision.



Gunnar Johansson, **Perception and Psychophysics**, 1973

Human actions: Historic overview



15th century
studies of
anatomy

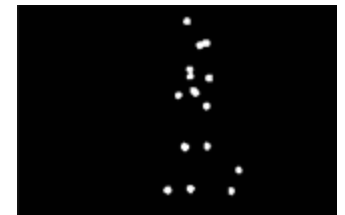


17th century
emergence of
biomechanics



19th century
emergence of
cinematography

1973
studies of human
motion perception



Modern computer vision



Modern applications: Motion capture and animation



Avatar (2009)

Modern applications: Motion capture and animation



Leonardo da Vinci (1452–1519)



Avatar (2009)

Modern applications: Video editing



Space-Time Video Completion

Y. Wexler, E. Shechtman and M. Irani, **CVPR** 2004

Modern applications: Video editing



Space-Time Video Completion

Y. Wexler, E. Shechtman and M. Irani, **CVPR** 2004

Modern applications: Video editing



Recognizing Action at a Distance

Alexei A. Efros, Alexander C. Berg, Greg Mori, Jitendra Malik, **ICCV** 2003

Modern applications: Video editing



Recognizing Action at a Distance

Alexei A. Efros, Alexander C. Berg, Greg Mori, Jitendra Malik, **ICCV** 2003

Why automatic video understanding?

- Huge amount of video is available and growing

BBC Motion Gallery



TV-channels recorded
since 60's



>34K hours of video
upload every day



~30M surveillance cameras in US
=> ~700K video hours/day



Movies



TV



YouTube



35%

Movies



34%

TV



40%

YouTube

Why action recognition

- Analyzing video archives



First appearance of
N. Sarkozy on TV



Sociology research:
Influence of character
smoking in movies



Education: How do I
make a pizza?

- Surveillance



Where is my cat?



Predicting crowd behavior
Counting people

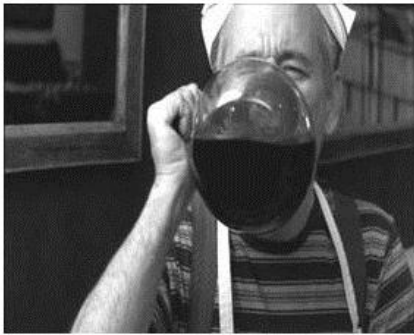
- Graphics



Motion capture and animation

Problem 1: Variability

- Need to deal with large appearance variations



Drinking



Smoking

- Large number of classes



falling



hugging



kicking



driving



Entering car



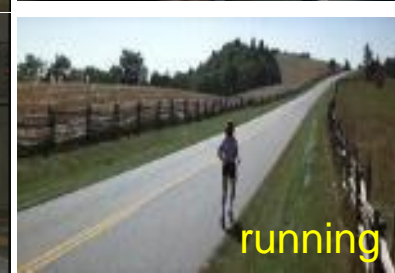
Answering phone



Hand-shaking



fighting



running



Standing up

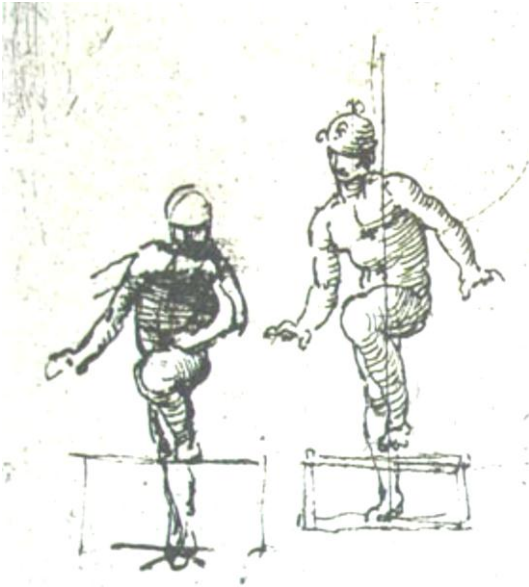
Problem 2: Granularity



Source: <http://www.youtube.com/watch?v=eYdUZdan5i8>

Do we want to learn *person-throws-cat-into-trash-bin* classifier?

Class overview



Motivation

- Historic review
- Modern applications

Appearance-based methods

- Motion history images
- Active shape models
- Tracking and motion priors

Motion-based methods

- Generic and parametric Optical Flow
- Motion templates

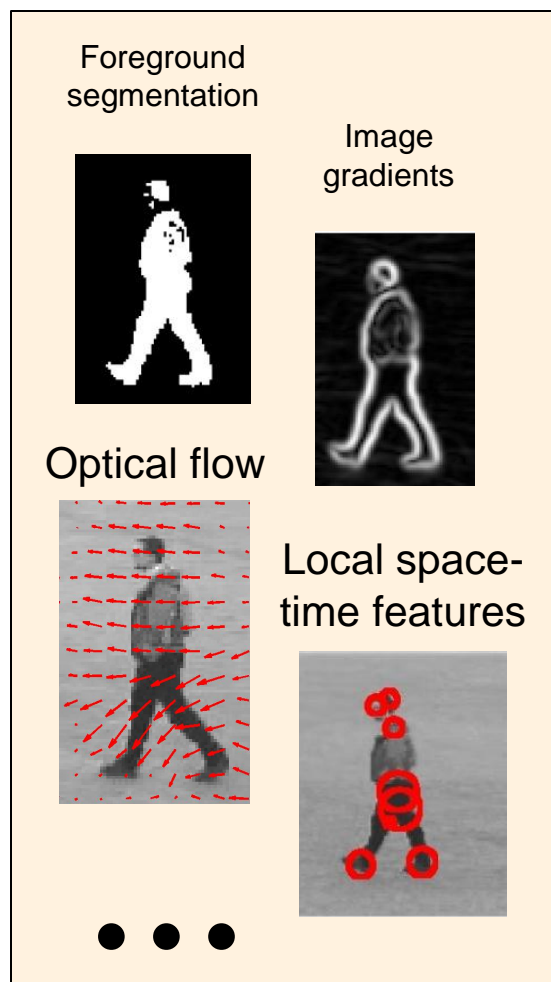
Space-time methods

- Local space-time features
- Action classification and detection
- Weakly-supervised action learning

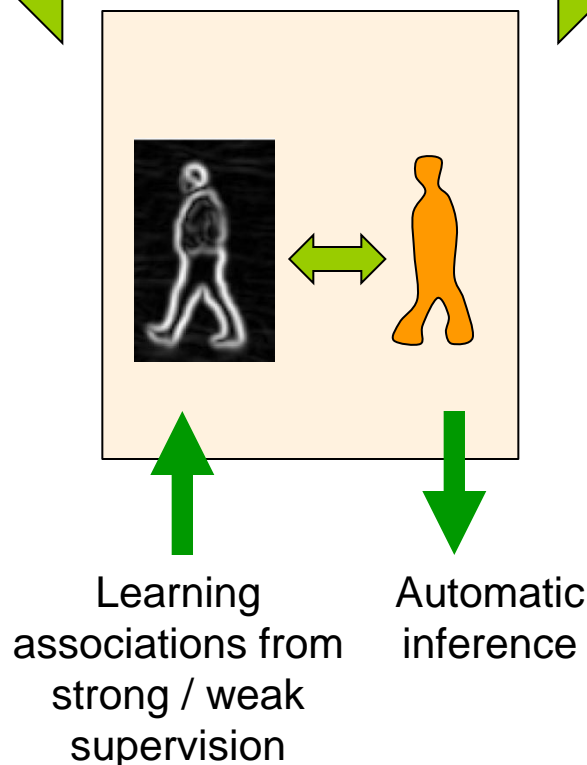
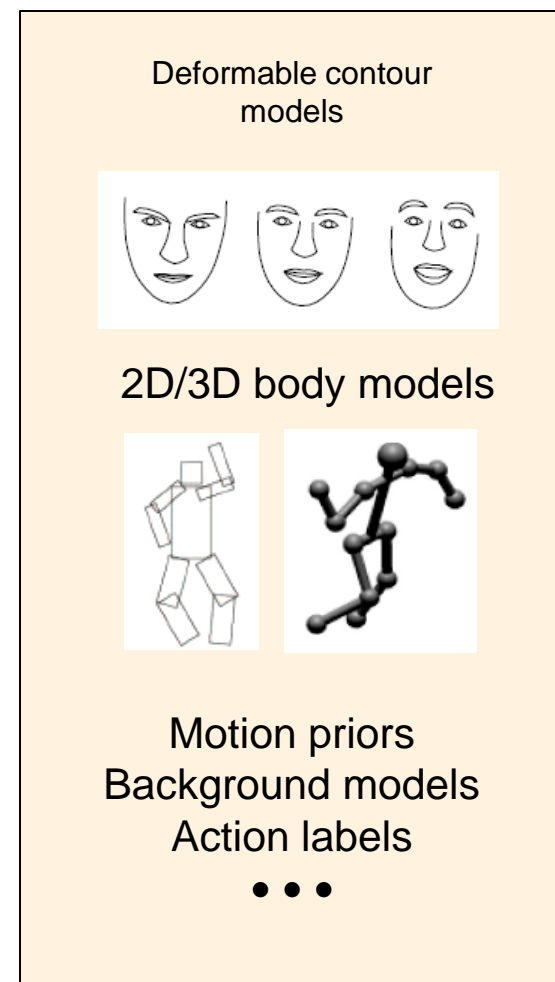
How to recognize actions?

Action understanding: Key components

Image measurements

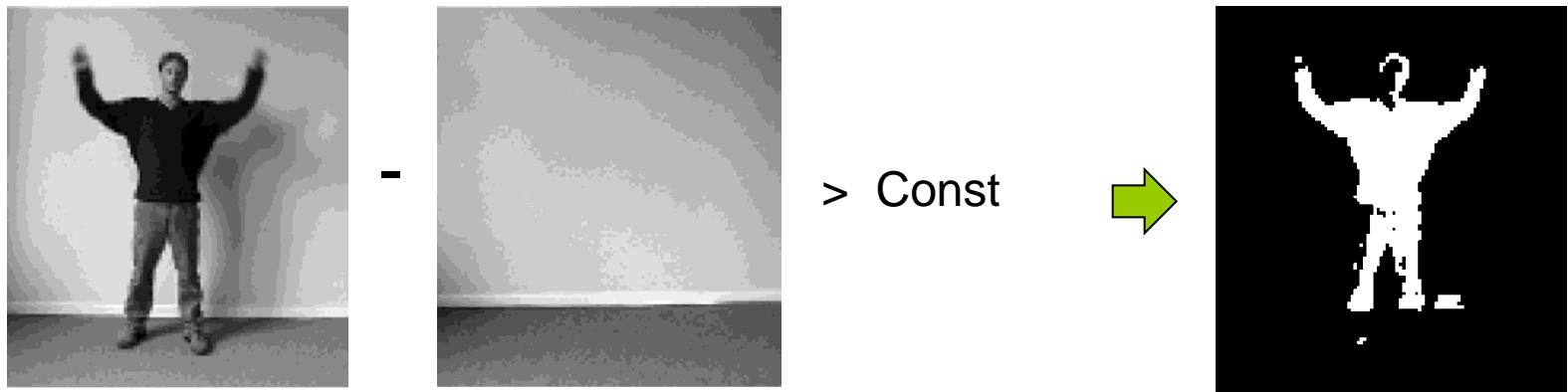


Prior knowledge



Foreground segmentation

Image differencing: a simple way to measure motion / temporal change



Better Background / Foreground separation methods exist:

- Modeling of color variation at each pixel with Gaussian Mixture
- Dominant motion compensation for sequences with moving camera
- Motion layer separation for scenes with non-static backgrounds

Temporal Templates

$$D(x, y, t) \quad t = 1, \dots, T$$



Idea: summarize motion in video in a
Motion History Image (MHI):

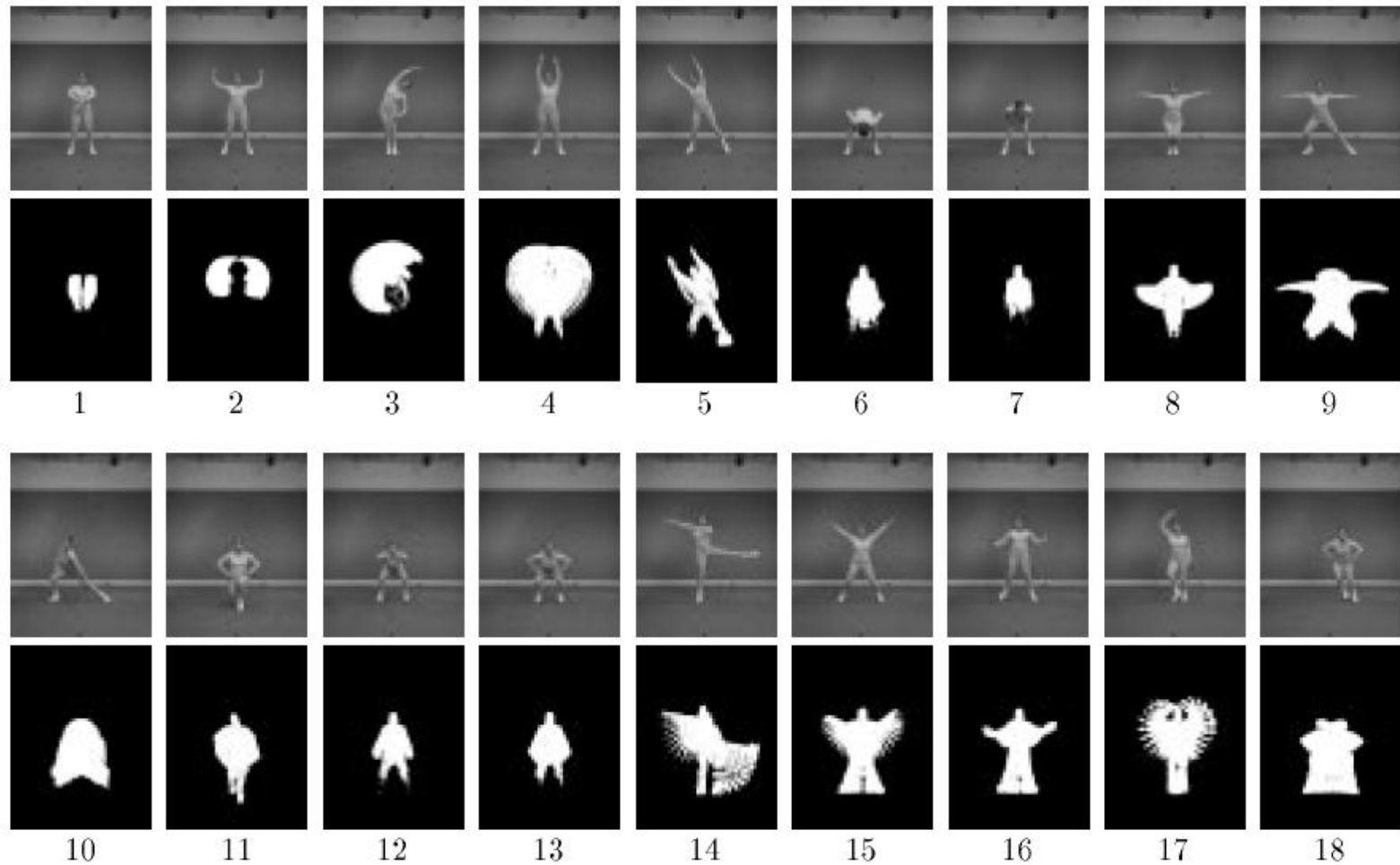
$$H_{\tau}(x, y, t) = \begin{cases} \tau & \text{if } D(x, y, t) = 1 \\ \max(0, H_{\tau}(x, y, t-1) - 1) & \text{otherwise} \end{cases}$$

Descriptor: Hu moments of different orders

$$m_{pq} = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} x^p y^q \rho(x, y) dx dy$$



Aerobics dataset



Nearest Neighbor classifier: 66% accuracy

Temporal Templates: Summary

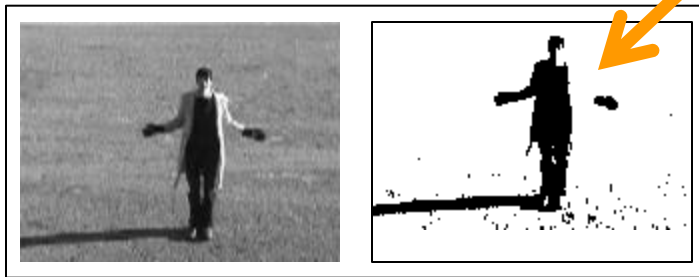
Pros:

- + Simple and fast
- + Works in controlled settings

Cons:

- Prone to errors of background subtraction

Not all shapes are valid
→ Restrict the space of admissible silhouettes



Variations in light, shadows, clothing...



What is the background here?

- Does not capture *interior* motion and shape



Silhouette tells little about actions

Active Shape Models of Cootes et al.

Point Distribution Model

- Represent the shape of samples by a set of corresponding points or *landmarks*

$$\mathbf{x} = (x_1, \dots, x_n, y_1, \dots, y_n)^T$$

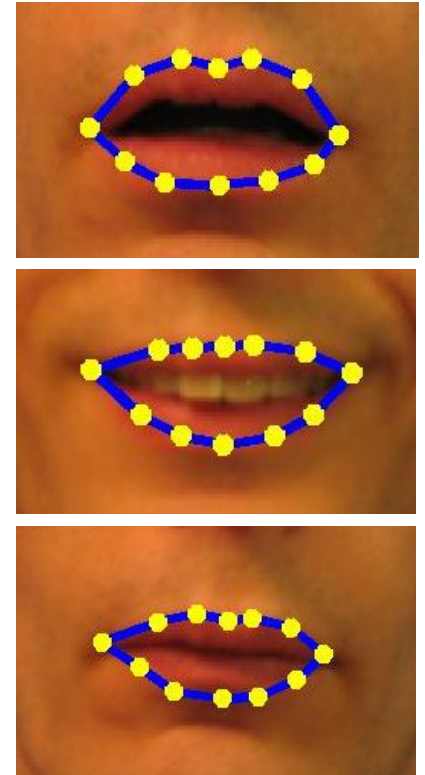
- Assume each shape can be represented by the linear combination of basis shapes

$$\Phi = (\phi_1 | \phi_2 | \dots | \phi_t)$$

such that $\mathbf{x} \approx \bar{\mathbf{x}} + \Phi \mathbf{b}$

for mean shape $\bar{\mathbf{x}} = \frac{1}{s} \sum_{i=1}^s \mathbf{x}_i$

and some parameters \mathbf{b}



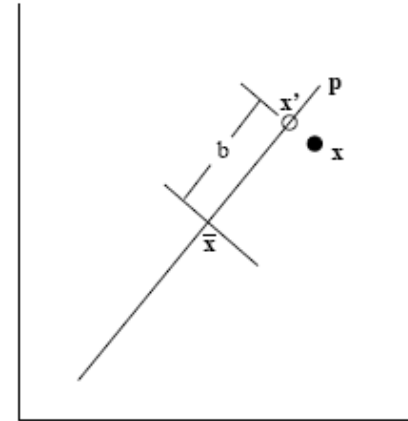
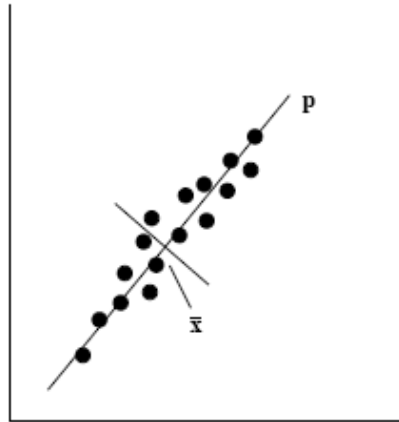
Active Shape Models of Cootes et al.

- Basis shapes can be found as the main modes of variation in the training data.

2D

Example:

(each point can be thought as a shape in N-Dim space)



Principle Component Analysis (PCA):

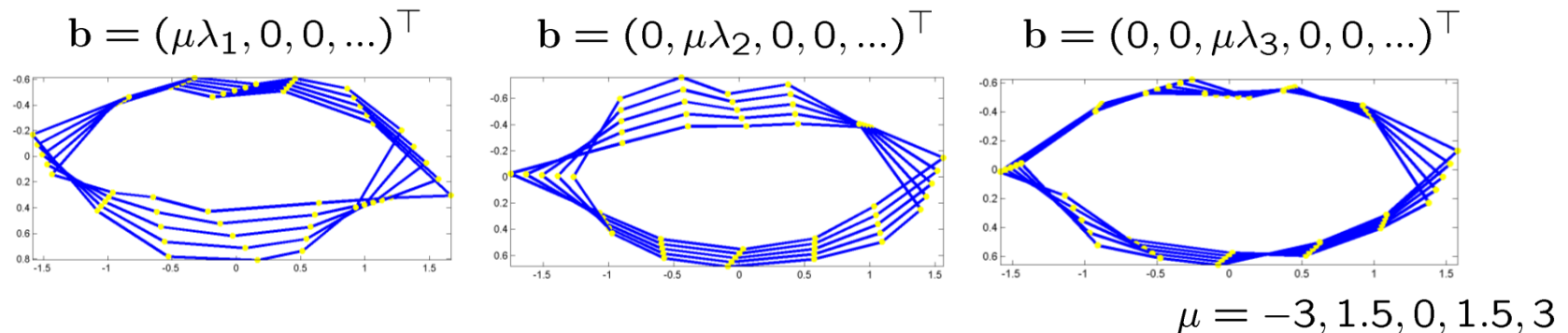
Covariance matrix
$$\mathbf{S} = \frac{1}{s-1} \sum_{i=1}^s (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T$$

Eigenvectors $\Phi = (\phi_1 | \phi_2 | \dots | \phi_t)$ eigenvalues $\lambda_1, \dots, \lambda_t$

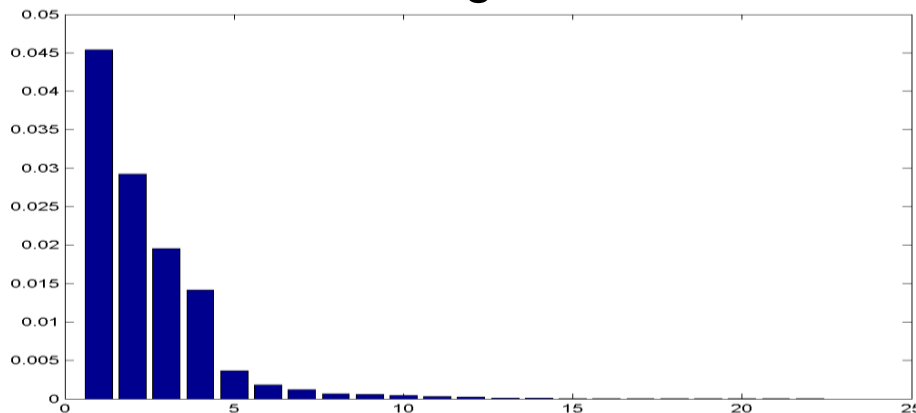
Active Shape Models of Cootes et al.

- Back-project from shape-space \mathbf{b} to image space $\mathbf{x} = \bar{\mathbf{x}} + \Phi \mathbf{b}$

➡ Three main modes of lips-shape variation:



Distribution of eigenvalues: $\lambda_1, \lambda_2, \lambda_3, \dots$



A small fraction of basis shapes (eigenvectors) accounts for the most of shape variation (\Rightarrow landmarks are redundant)

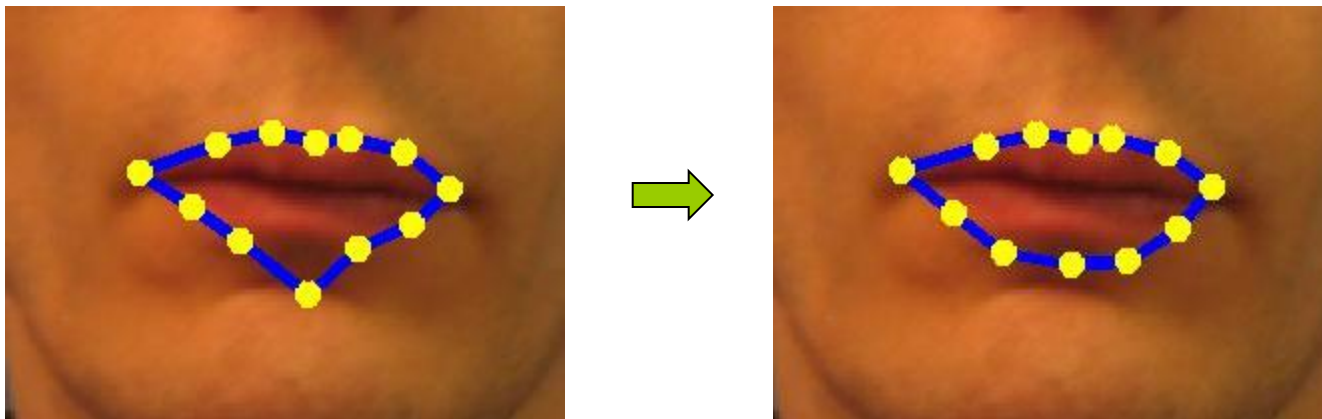
Active Shape Models of Cootes et al.

- Φ is orthonormal basis, therefore $\Phi^{-1} = \Phi^T$
➡ Given estimate of \mathbf{x} we can recover shape parameters \mathbf{b}

$$\mathbf{b} = \Phi^T (\mathbf{x} - \bar{\mathbf{x}})$$

- Projection onto the shape-space serves as a *regularization*

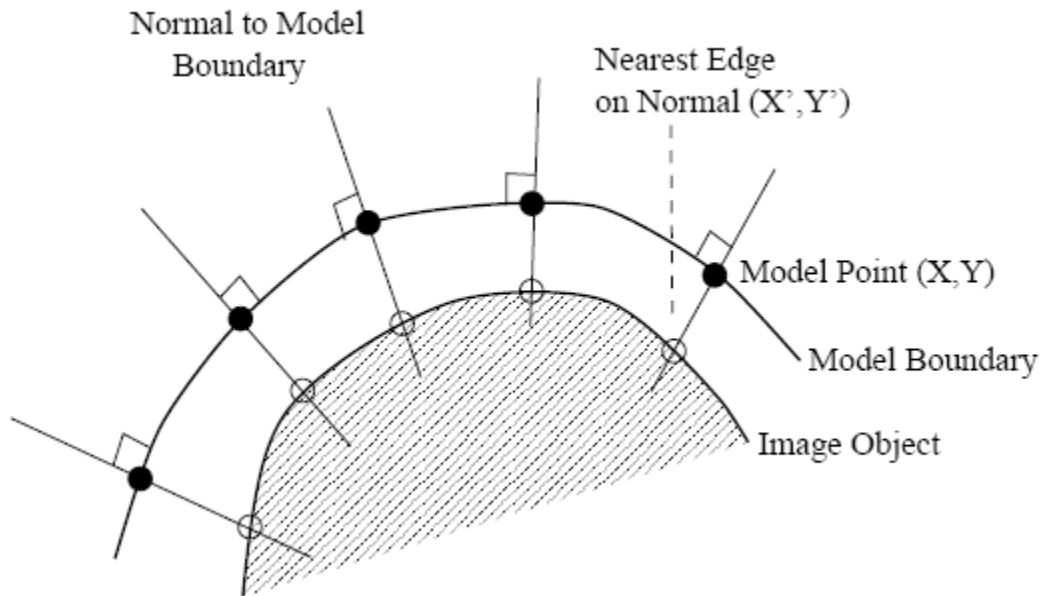
$$\mathbf{x} \quad \text{➡} \quad \mathbf{b} = \Phi^T (\mathbf{x} - \bar{\mathbf{x}}) \quad \text{➡} \quad \mathbf{x}_{\text{reg}} = \bar{\mathbf{x}} + \Phi \mathbf{b}$$



Active Shape Models of Cootes et al.

How to use Active Shape Models for shape estimation?

- Given initial guess of model points \bar{x} estimate new positions x' using local image search, e.g. locate the closest edge point



- Re-estimate shape parameters

$$b' = \Phi^T (x' - \bar{x})$$

Active Shape Models of Cootes et al.

- Iterative ASM alignment algorithm
 1. Initialize with the reasonable guess of \mathbf{T} and $\mathbf{b} = \mathbf{0}^\top$
 2. Estimate \mathbf{x}' from image measurements
 3. Re-estimate \mathbf{T}, \mathbf{b}
 4. Unless \mathbf{T}, \mathbf{b} converged, repeat from step 2

Example: face alignment

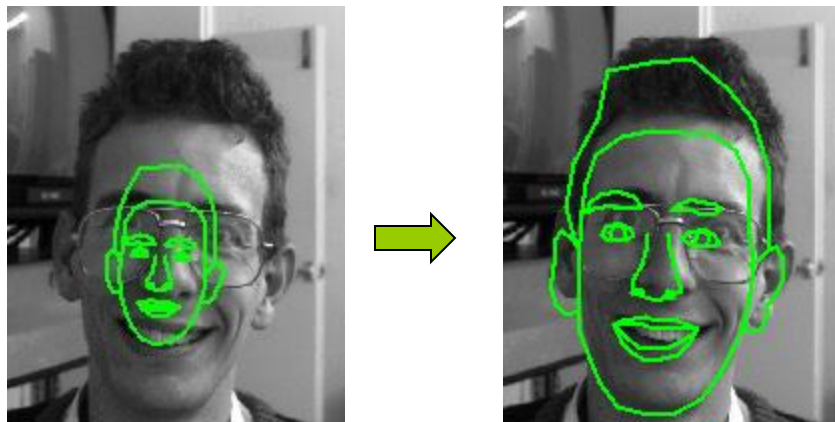
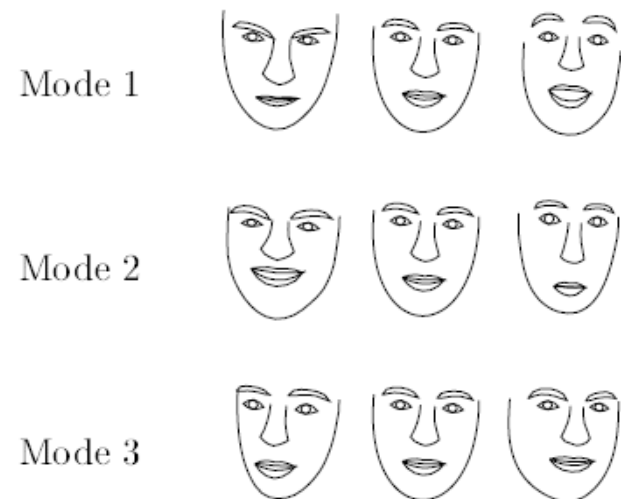


Illustration of face shape space



Active Shape Models: Their Training and Application

T.F. Cootes, C.J. Taylor, D.H. Cooper, and J. Graham, **CVIU** 1995

Active Shape Model tracking

Aim: to track ASM of time-varying shapes, e.g. human silhouettes

- Impose time-continuity constraint on model parameters.
For example, for shape parameters b :

$$b_i^{(k)} = b_i^{(k-1)} + w_i^{k-1}$$

$$w_i \sim \mathcal{N}(0, \mu\lambda_i) \quad \text{Gaussian noise}$$

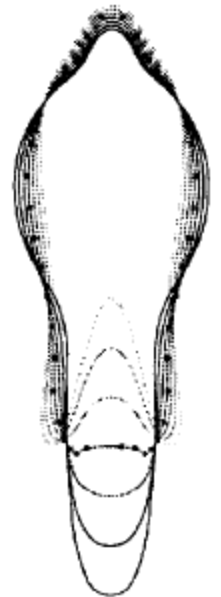
For similarity transformation T

$$a^{(k)} = a^{(k-1)} + w_a^{k-1}, \quad w_a = \mathcal{N}(0, \sigma_a)$$

$$t_{x|y}^{(k)} = t_{x|y}^{(k-1)} + v_{x|y}^{(k-1)} + w_{x|y}^{k-1}, \quad w_{x|y} = \mathcal{N}(0, \sigma_{x|y})$$

More complex dynamical models possible

- Update model parameters at each time frame using e.g. Kalman filter

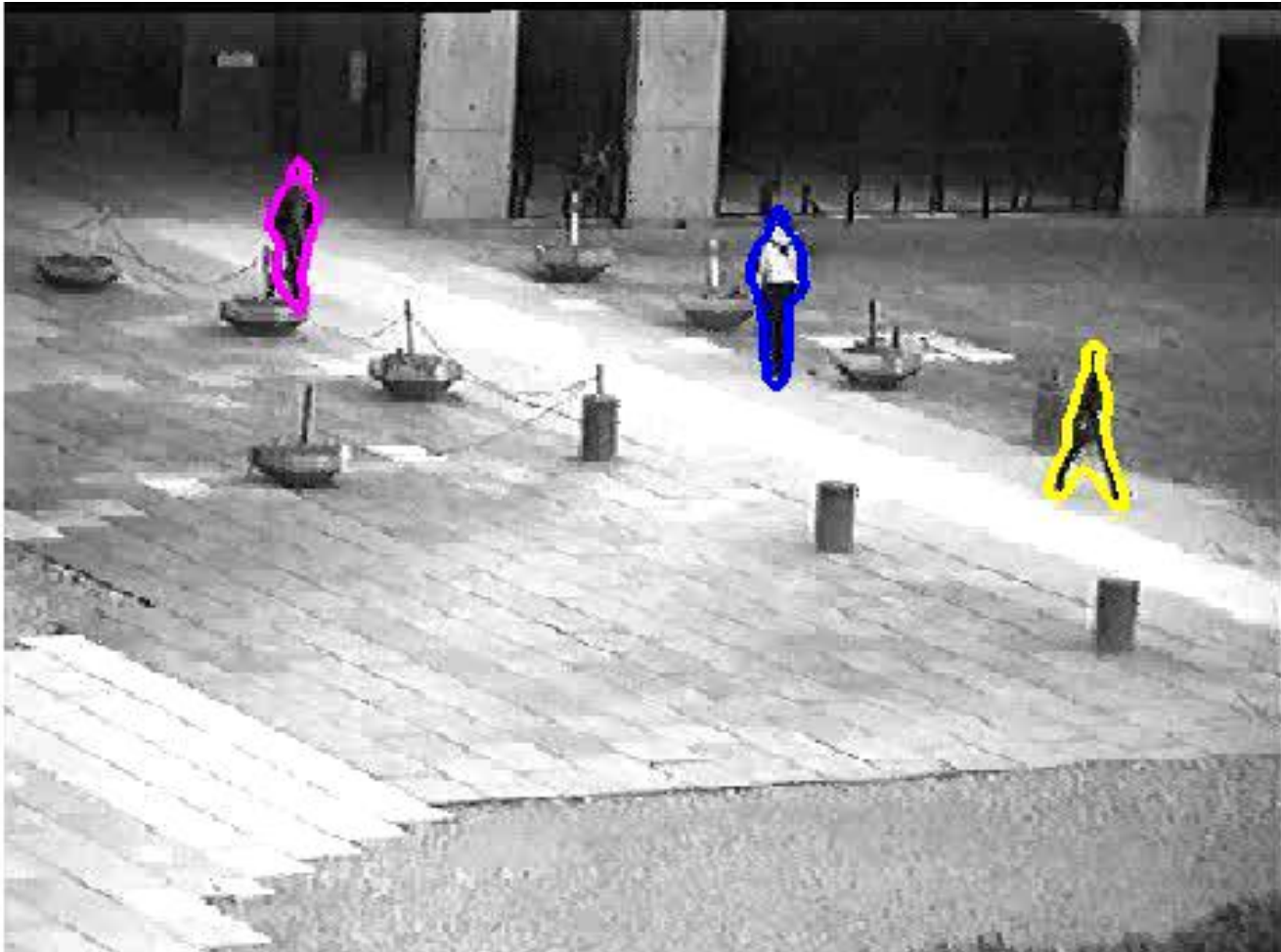


Person Tracking



Learning flexible models from image sequences
A. Baumberg and D. Hogg, **ECCV** 1994

Person Tracking



Learning flexible models from image sequences
A. Baumberg and D. Hogg, **ECCV** 1994

Active Shape Models: Summary

Pros:

- + Shape prior helps overcoming segmentation errors
- + Fast optimization
- + Can handle interior/exterior dynamics

Cons:

- Optimization gets trapped in local minima
- Re-initialization is problematic

Possible improvements:


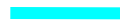

- Learn and use motion priors, possibly specific to different actions

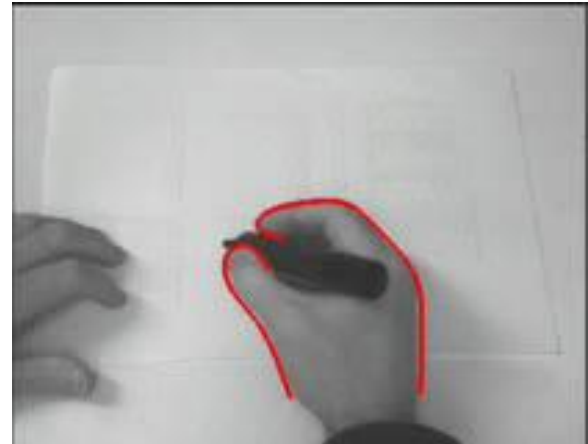
Motion priors

- Accurate motion models can be used both to:
 - ❖ Help accurate tracking
 - ❖ Recognize actions
- Goal: formulate motion models for different types of actions and use such models for action recognition

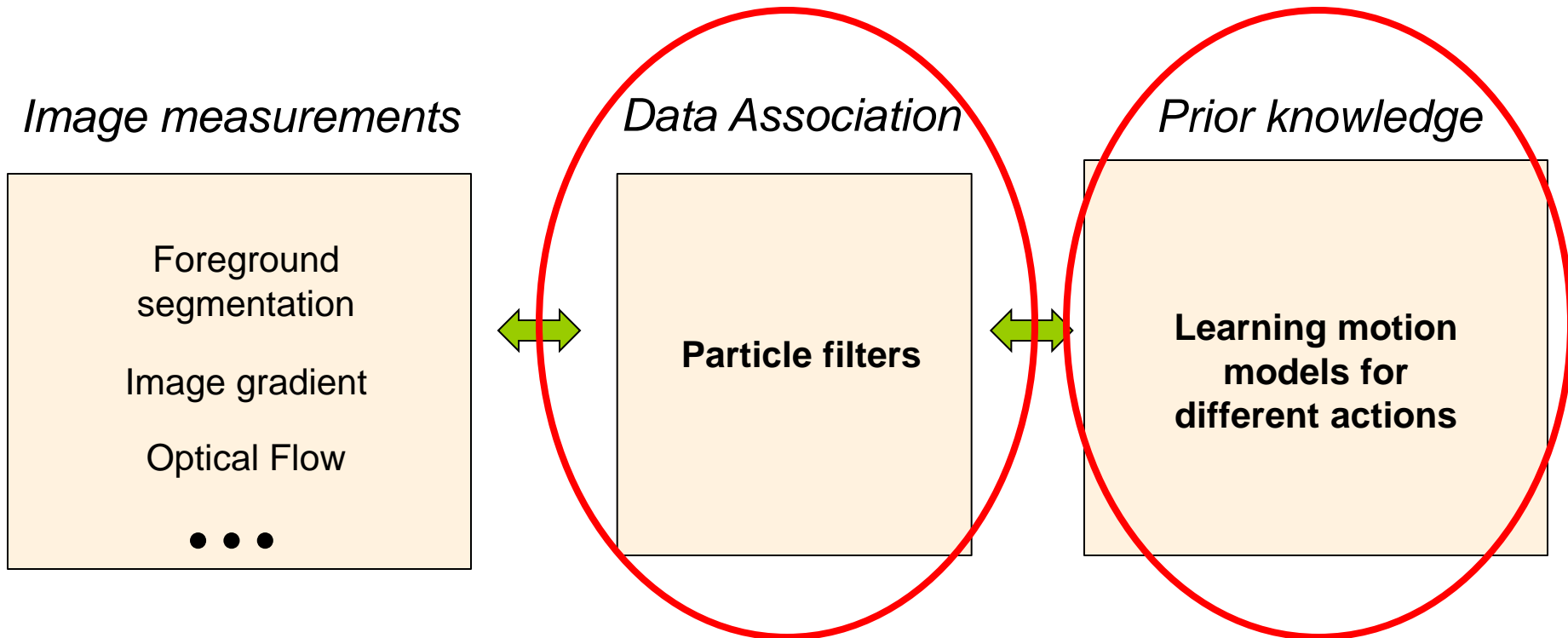
Example:

Drawing with 3 action modes

-  line drawing
-  scribbling
-  idle



Incorporating motion priors



Bayesian Tracking

General framework: recognition by synthesis;
 generative models;
 finding best explanation of the data

Notation:

Z_i image data at time i

X_i model parameters at time i (e.g. shape and its dynamics)

$p(X_i)$ prior density for X_i

$p(Z_i|X_i)$ likelihood of data for the given model configuration

We search posterior defined by the Bayes' rule

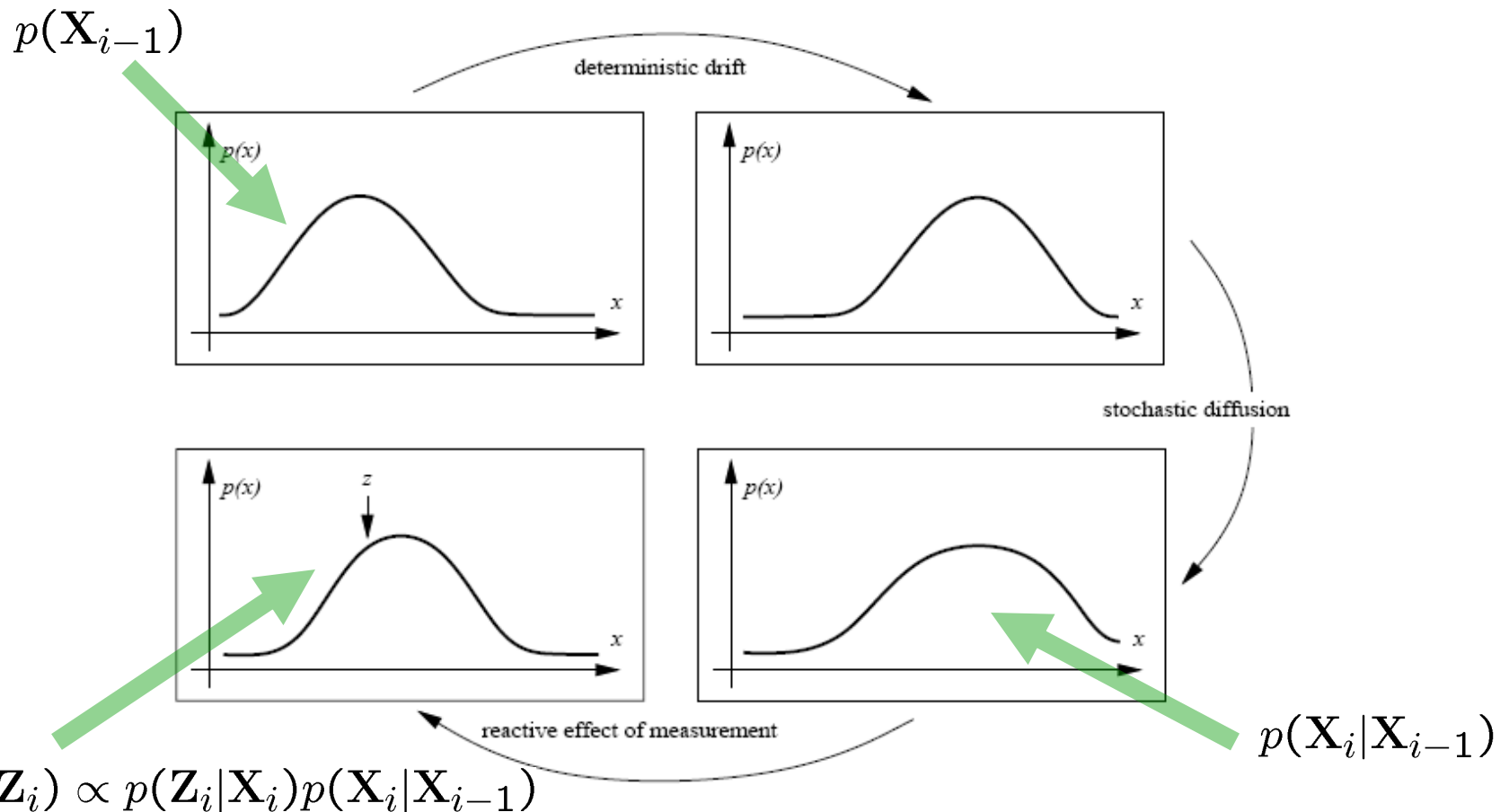
$$p(X|Z) \propto p(Z|X)p(X)$$

For tracking the Markov assumption gives the prior $p(X_i|X_{i-1})$

Temporal update rule: $p(X_i|Z_i) \propto p(Z_i|X_i)p(X_i|X_{i-1})$

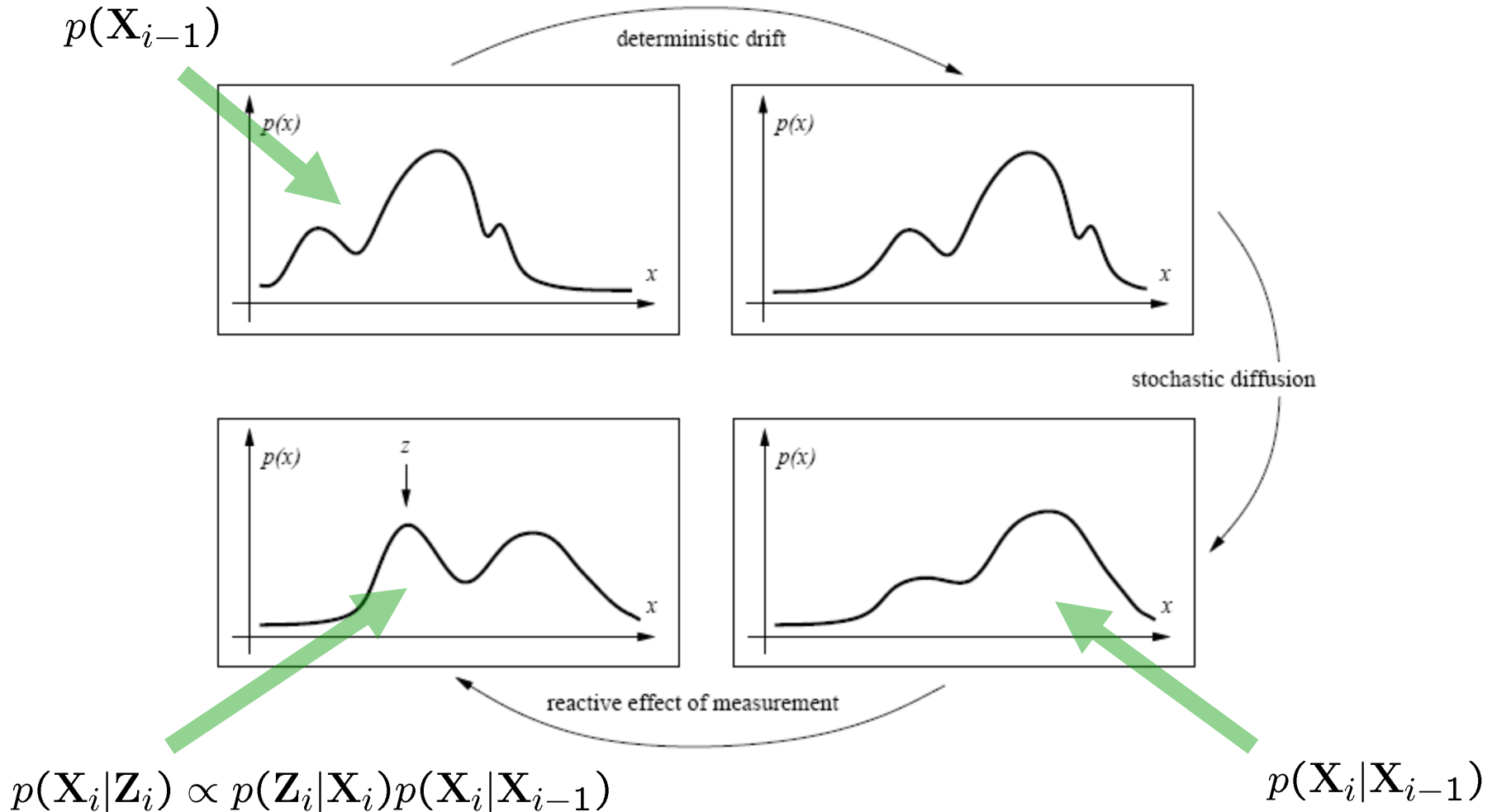
Kalman Filtering

If all probability densities are uni-modal, specifically Gaussians, the posterior can be evaluated in the closed form



Particle Filtering

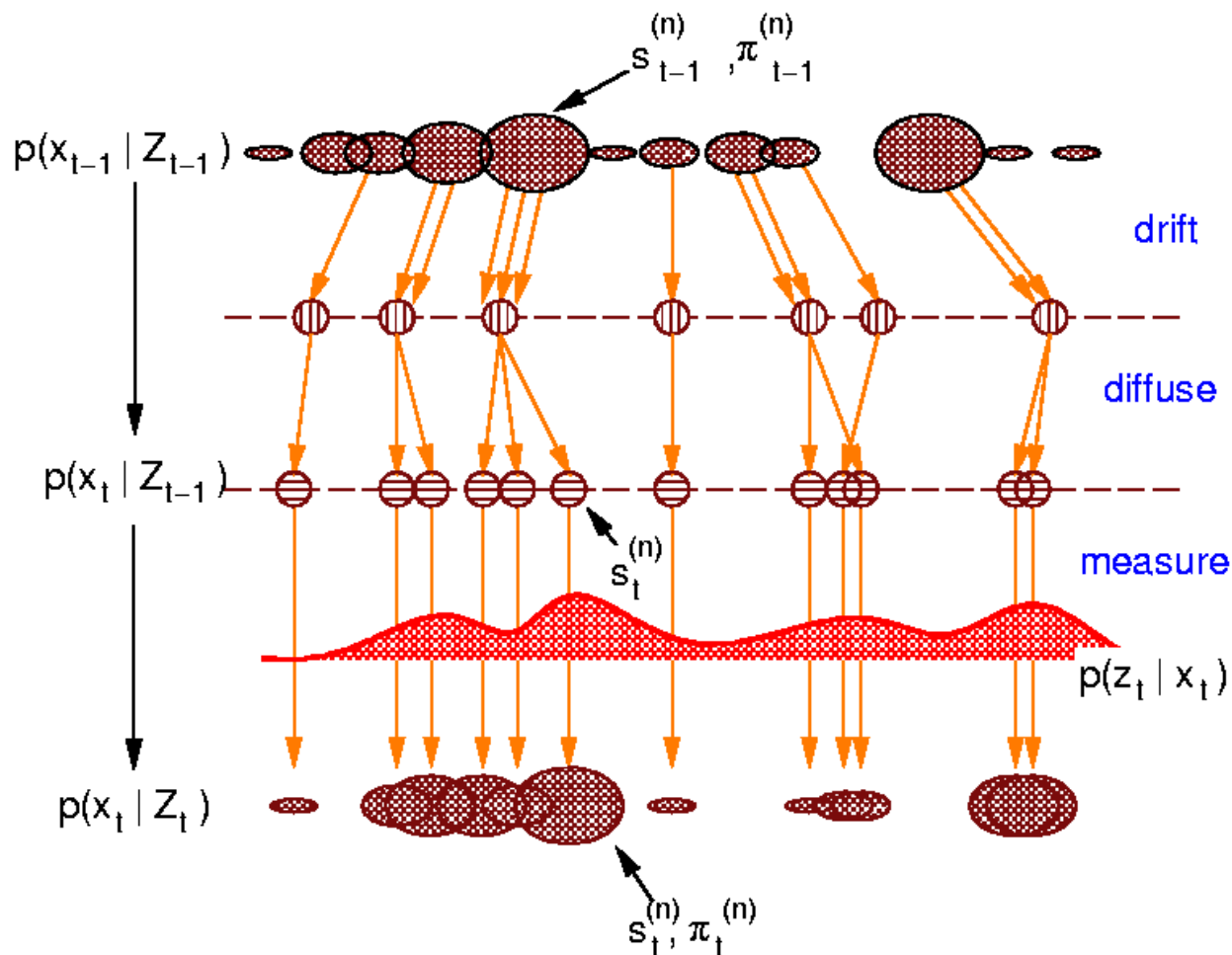
In reality probability densities are almost always *multi-modal*



Particle Filtering

In reality probability densities are almost always *multi-modal*

➡ Approximate distributions with weighted particles



Particle Filtering

Tracking examples:

X describes leave shape



X describes head shape



CONDENSATION - conditional density propagation for visual tracking
A. Blake and M. Isard **IJCV** 1998

Learning dynamic prior

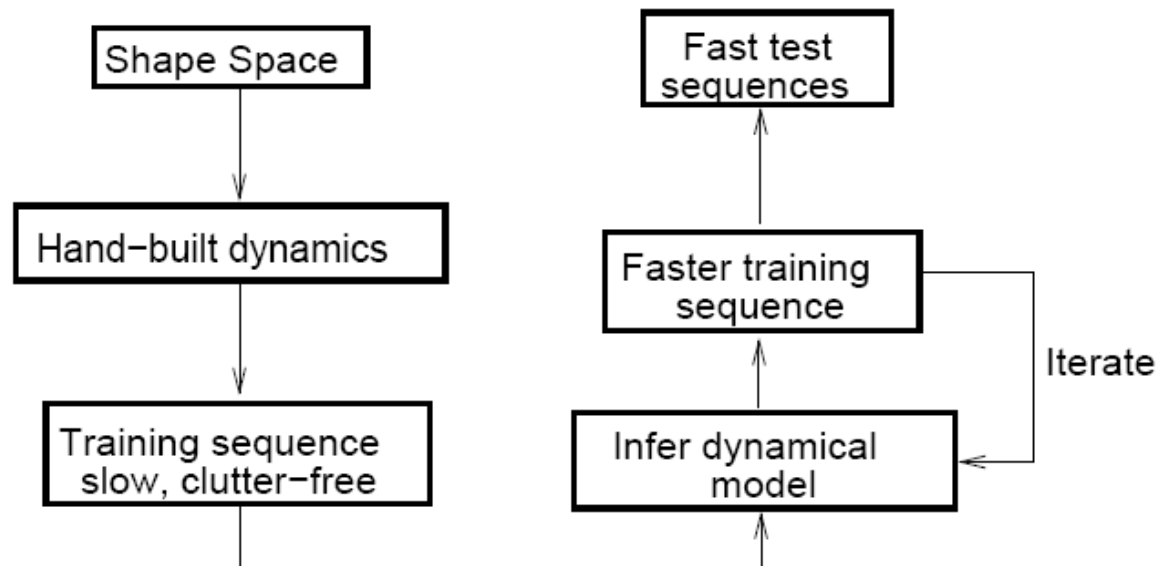
- Dynamic model: 2nd order Auto-Regressive Process

State $\mathcal{X}_k = \begin{pmatrix} \mathbf{X}_{k-1} \\ \mathbf{X}_k \end{pmatrix}$

Update rule: $\mathcal{X}_k - \overline{\mathcal{X}} = A(\mathcal{X}_{k-1} - \overline{\mathcal{X}}) + B\mathbf{w}_k$

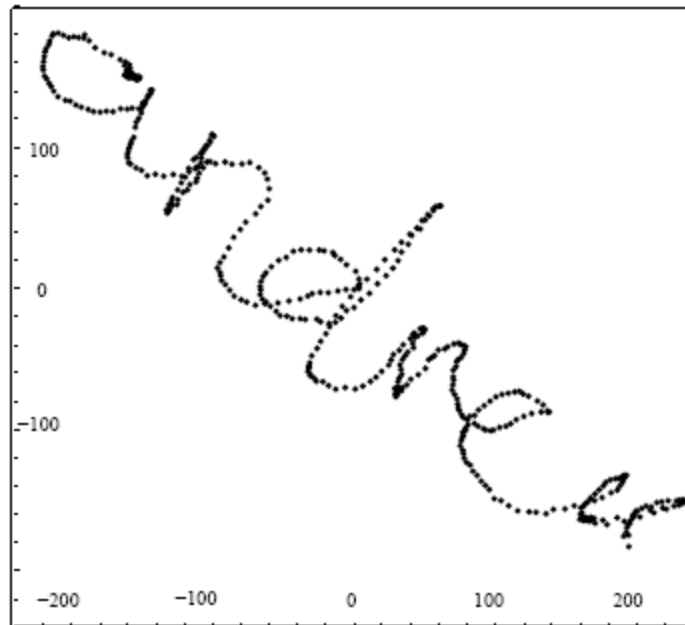
Model parameters: $A = \begin{pmatrix} 0 & I \\ A_2 & A_1 \end{pmatrix}$, $\overline{\mathcal{X}} = \begin{pmatrix} \overline{\mathbf{X}} \\ \overline{\mathbf{X}} \end{pmatrix}$ and $B = \begin{pmatrix} 0 \\ B_0 \end{pmatrix}$

Learning scheme:

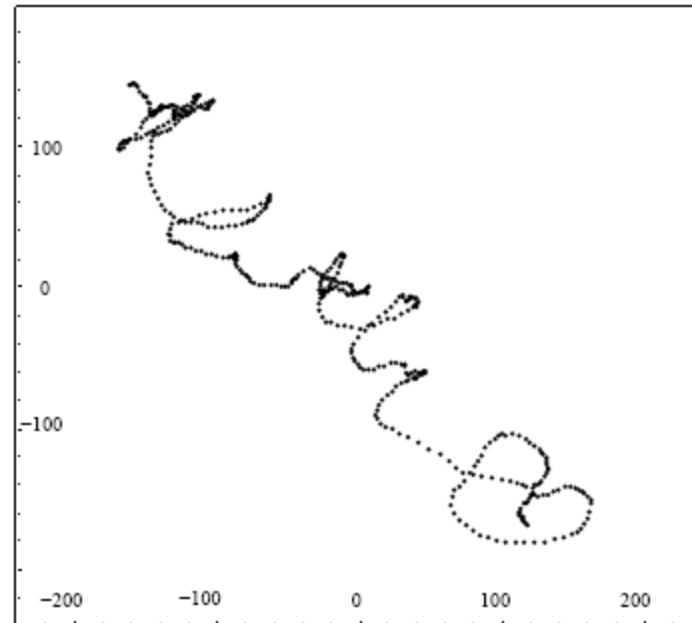


Learning dynamic prior

Learning point sequence



Random simulation of the learned dynamical model

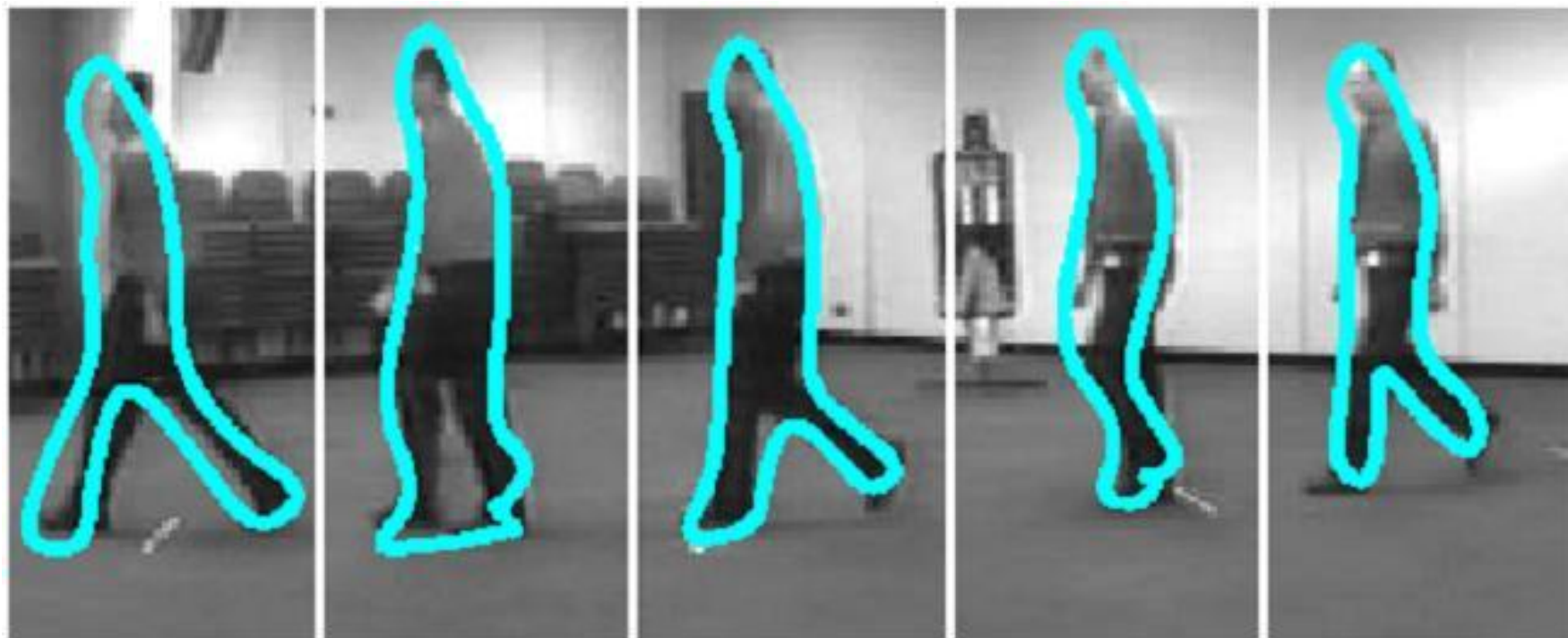


Statistical models of visual shape and motion

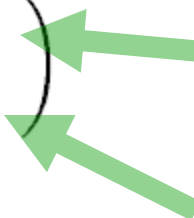
A. Blake, B. Bascle, M. Isard and J. MacCormick, **Phil.Trans.R.Soc. 1998**

Learning dynamic prior

Random simulation of the learned gate dynamics



Dynamics with discrete states

Introduce “mixed” state $\mathcal{X}_k^+ = \begin{pmatrix} \mathcal{X}_k \\ y_k \end{pmatrix}$  Continuous state space (as before)

Transition probability matrix

$$P(y_k = j | y_{k-1} = i) = T_{i,j},$$

or more generally $P(y_k = j | y_{k-1} = i, \mathcal{X}_{k-1}) = T_{i,j}(\mathcal{X}_{k-1})$

Incorporation of the mixed-state model into a particle filter is straightforward, simply use \mathcal{X}_k^+ instead of \mathcal{X}_k and the corresponding update rules

Dynamics with discrete states

Example: Drawing

Transition probability matrix

	line	idle	scribbling
line	0.9800	0.0015	0.0185
idle	0.0850	0.9000	0.0150
scribbling	0.0050	0.0150	0.9800

$$T = \begin{pmatrix} 0.9800 & 0.0015 & 0.0185 \\ 0.0850 & 0.9000 & 0.0150 \\ 0.0050 & 0.0150 & 0.9800 \end{pmatrix}$$

Result: simultaneously improved tracking and gesture recognition

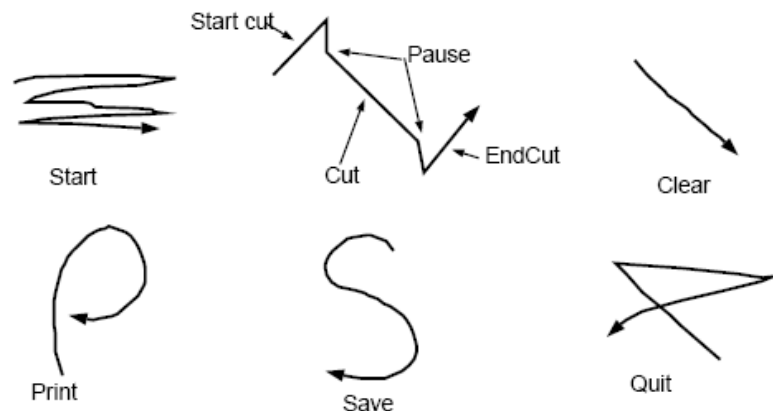


— line drawing
— scribbling
— idle

A mixed-state Condensation tracker with automatic model-switching
M. Isard and A. Blake, **ICCV** 1998

Dynamics with discrete states

Similar illustrated on
gesture recognition in
the context of a visual
black-board interface



Motion priors & Tracking: Summary

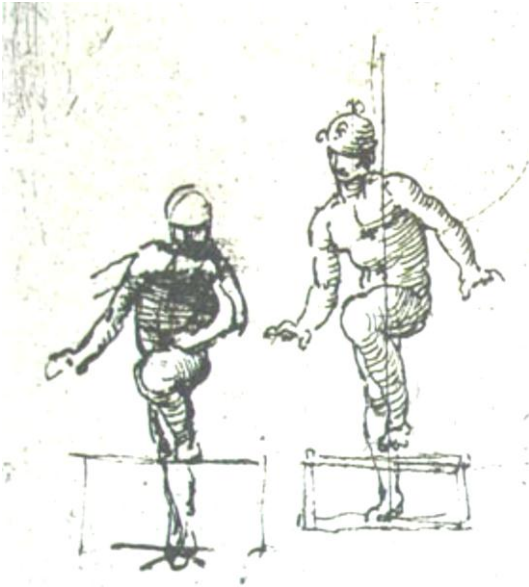
Pros:

- + more accurate tracking using specific motion models
- + Simultaneous tracking and motion recognition with discrete state dynamical models

Cons:

- Local minima is still an issue
- Re-initialization is still an issue

Class overview



Motivation

Historic review

Modern applications

Appearance-based methods

Motion history images

Active shape models

Tracking and motion priors

Motion-based methods

Generic and parametric Optical Flow

Motion templates

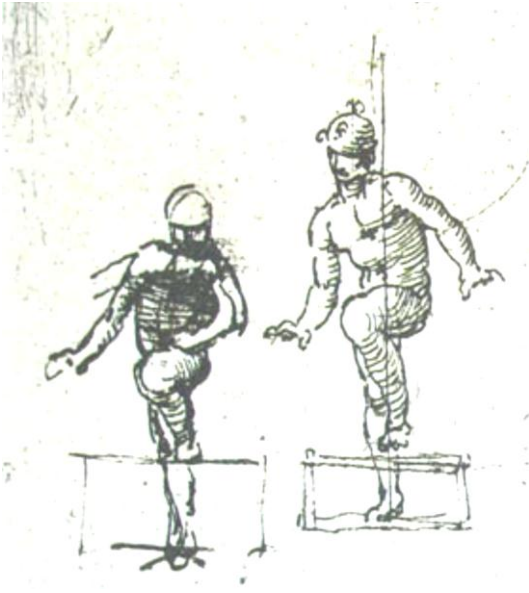
Space-time methods

Local space-time features

Action classification and detection

Weakly-supervised action learning

Class overview



Motivation

Historic review

Modern applications

Appearance-based methods

Motion history images

Active shape models

Tracking and motion priors

Motion-based methods

Generic and parametric Optical Flow

Motion templates

Space-time methods

Local space-time features

Action classification and detection

Weakly-supervised action learning

Shape and Appearance vs. Motion

- Shape and appearance in images depends on many factors: clothing, illumination contrast, image resolution, etc...



[Efros et al. 2003]

- Motion field (in theory) is invariant to shape and can be used directly to describe human actions



Motion estimation: Optical Flow

- Classic problem of computer vision [Gibson 1955]

- Goal: estimate **motion field**

How? We only have access to image pixels



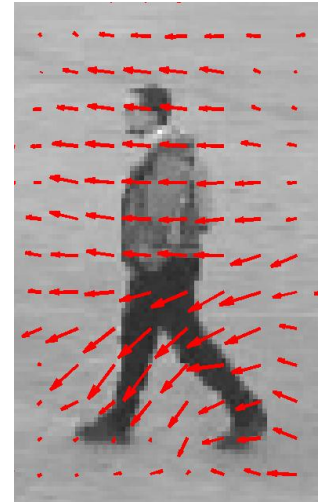
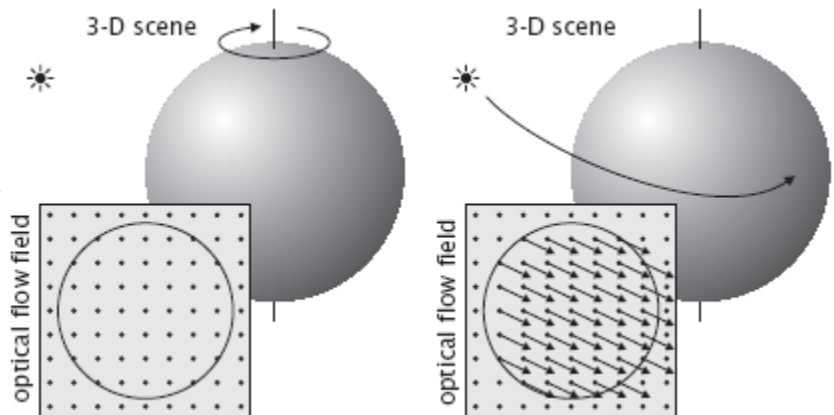
Estimate pixel-wise correspondence
between frames = **Optical Flow**

- **Brightness Change** assumption: corresponding pixels preserve their intensity (color)

❖ Useful assumption in many cases

❖ Breaks at occlusions and illumination changes

❖ Physical and visual motion may be different



Generic Optical Flow

- Brightness Change Constraint Equation (BCCE)

$$(\nabla I)^\top \mathbf{v} + I_t = 0$$

$\mathbf{v} = (v_x, v_y)^\top$ Optical flow
 $\nabla I = (I_x, I_y)^\top$ Image gradient

One equation, two unknowns \Rightarrow cannot be solved directly

➡ Integrate several measurements in the local neighborhood and obtain a *Least Squares Solution* [Lucas & Kanade 1981]

$$\langle \nabla I (\nabla I)^\top \rangle \mathbf{v} = - \langle \nabla I I_t \rangle$$

Second-moment matrix, the same one used to compute Harris interest points!

$$\begin{pmatrix} \langle I_x^2 \rangle & \langle I_x I_y \rangle \\ \langle I_x I_y \rangle & \langle I_y^2 \rangle \end{pmatrix} \mathbf{v} = - \begin{pmatrix} \langle I_x I_t \rangle \\ \langle I_y I_t \rangle \end{pmatrix}$$

$\langle \cdot \rangle$ Denotes integration over a spatial (or spatio-temporal) neighborhood of a point

Generic Optical Flow

- The solution of $\langle \nabla I (\nabla I)^\top \rangle \mathbf{v} = - \langle \nabla I I_t \rangle$ assumes
 1. Brightness change constraint holds in $\langle \cdot \rangle$
 2. Sufficient variation of image gradient in $\langle \cdot \rangle$
 3. Approximately constant motion in $\langle \cdot \rangle$

Motion estimation becomes *inaccurate* if any of assumptions 1-3 is violated.

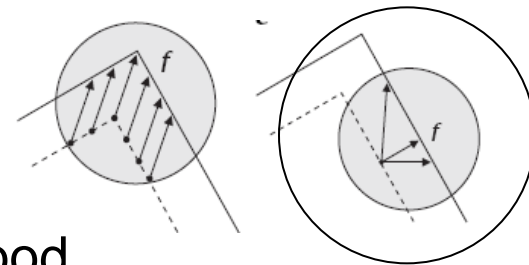
- Solutions:

- (2) Insufficient gradient variation known as *aperture problem*

➡ Increase integration neighborhood

- (3) Non-constant motion in $\langle \cdot \rangle$

➡ Use more sophisticated motion model

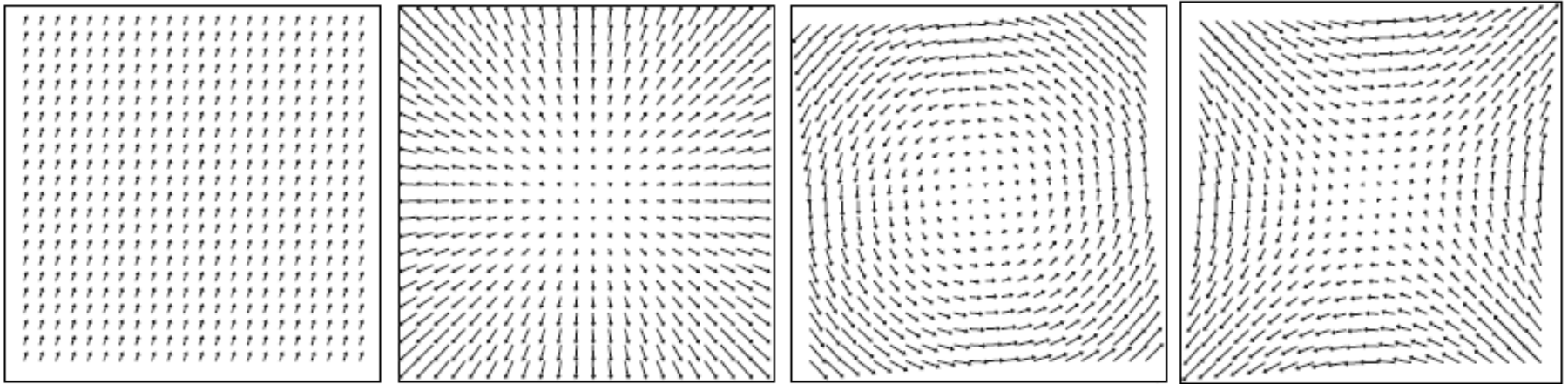


Parameterized Optical Flow

- Constant velocity model: $\mathbf{v} = \begin{pmatrix} v_x \\ v_y \end{pmatrix}$
- Upgrade to affine motion model: $\mathbf{v} = \begin{pmatrix} a_1 & a_2 \\ a_3 & a_4 \end{pmatrix} \begin{pmatrix} x \\ y \end{pmatrix} + \begin{pmatrix} v_x \\ v_y \end{pmatrix}$

Now motion depends on the position $(x, y)^\top$ inside the neighborhood

Examples of Affine motion models for different parameters:



- Can be formulated as Least Squares approach to estimate \mathbf{v} as before!

Parameterized Optical Flow

- Another extension of the constant motion model is to compute PCA basis flow fields from training examples

1. Compute standard Optical Flow for many examples
2. Put velocity components into one vector

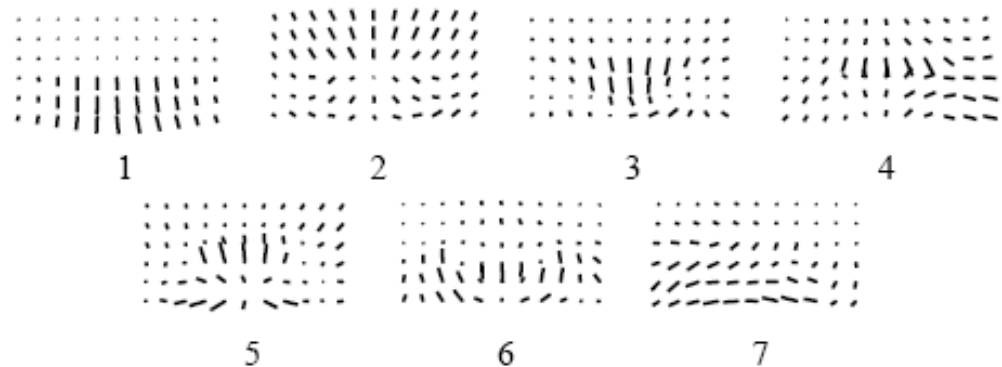
$$\mathbf{w} = (v_x^1, v_y^1, v_x^2, v_y^2, \dots, v_x^n, v_y^n)^\top$$

3. Do PCA on \mathbf{w} and obtain most informative PCA flow basis vectors

Training samples



PCA flow bases



Learning Parameterized Models of Image Motion

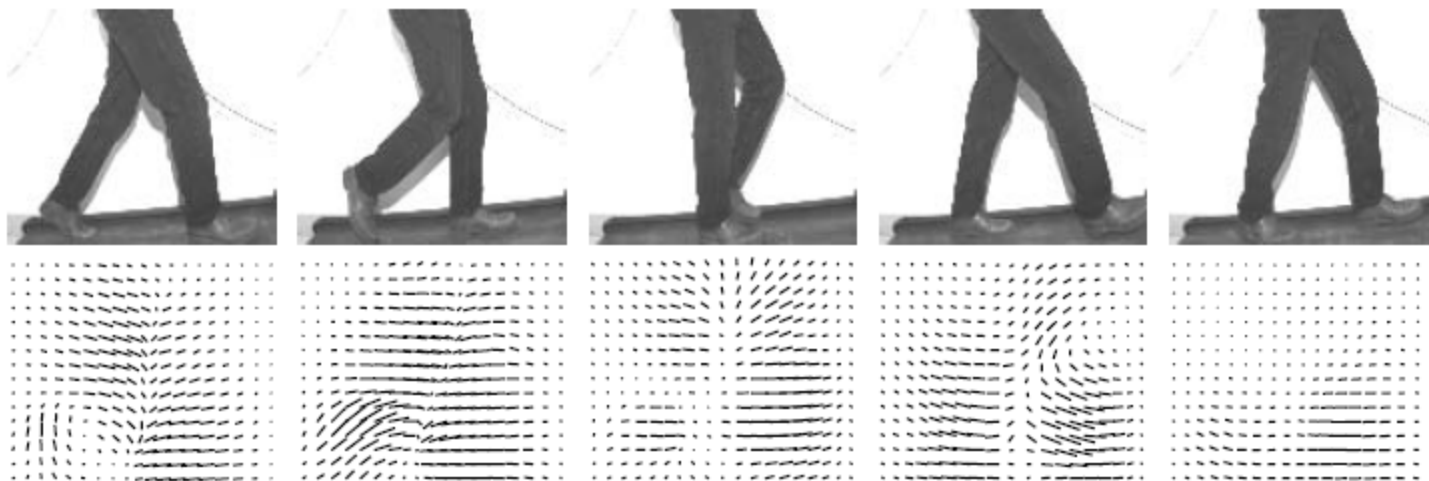
M.J. Black, Y. Yacoob, A.D. Jepson and D.J. Fleet, **CVPR 1997**

Parameterized Optical Flow

- Use PCA flow bases to *regularize* solution of motion estimation
- Motion estimation for test samples can be computed *without* explicit computation of optical flow!

Solution formulation e.g. in terms of Least Squares

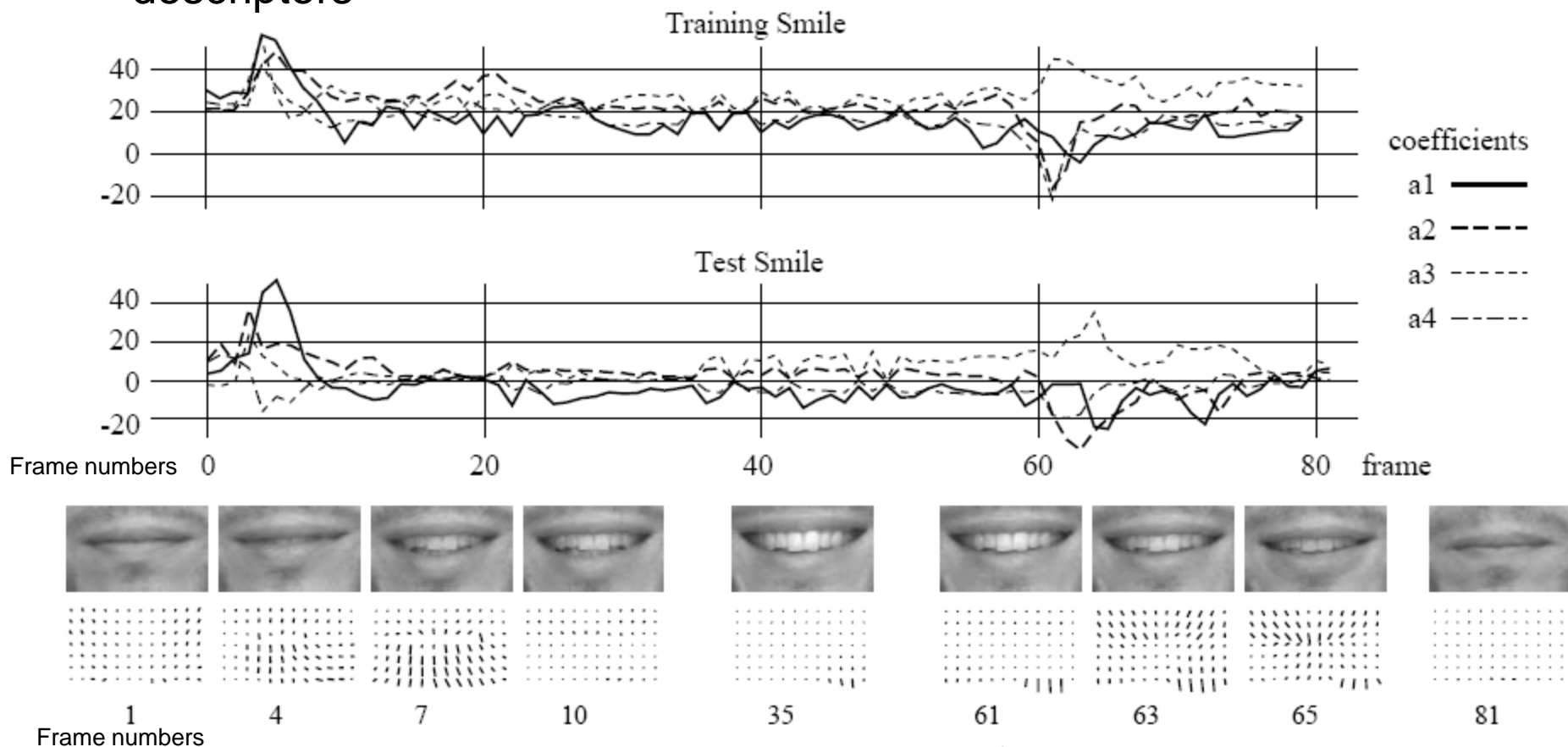
Direct flow recovery:



Learning Parameterized Models of Image Motion
M.J. Black, Y. Yacoob, A.D. Jepson and D.J. Fleet, **CVPR 1997**

Parameterized Optical Flow

- Estimated coefficients of PCA flow bases can be used as action descriptors

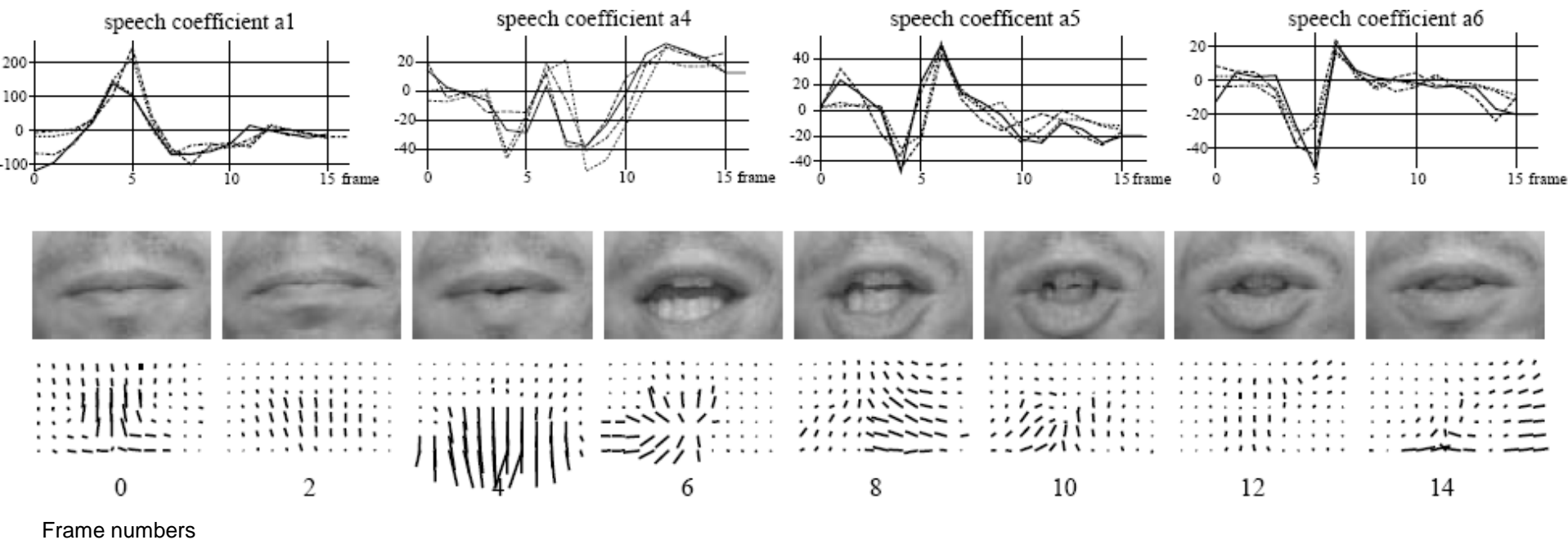


Learning Parameterized Models of Image Motion

M.J. Black, Y. Yacoob, A.D. Jepson and D.J. Fleet, **CVPR 1997**

Parameterized Optical Flow

- Estimated coefficients of PCA flow bases can be used as action descriptors



Optical flow seems to be an interesting descriptor for motion/action recognition

Spatial Motion Descriptor

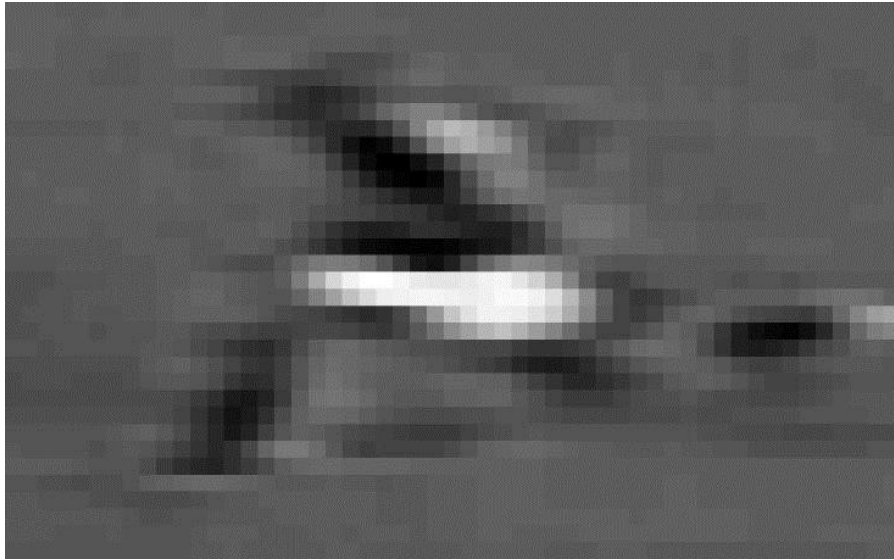
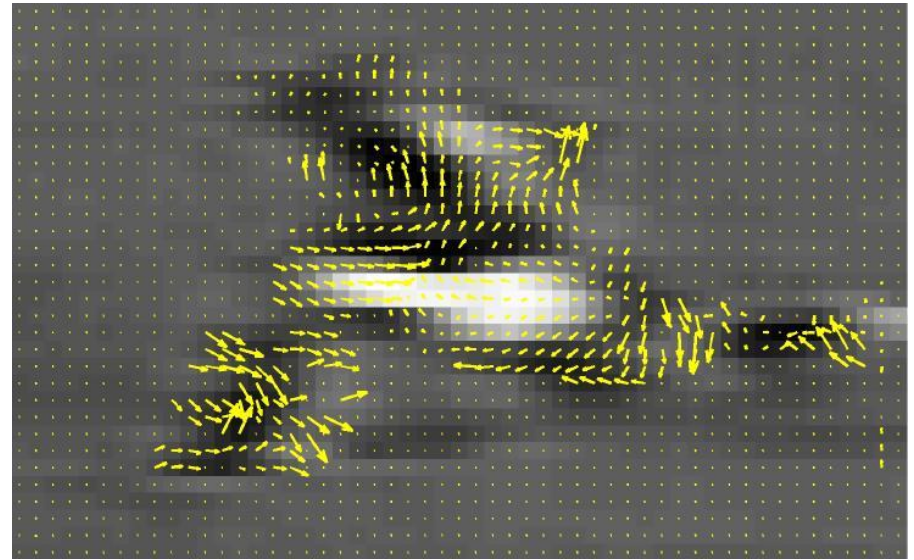
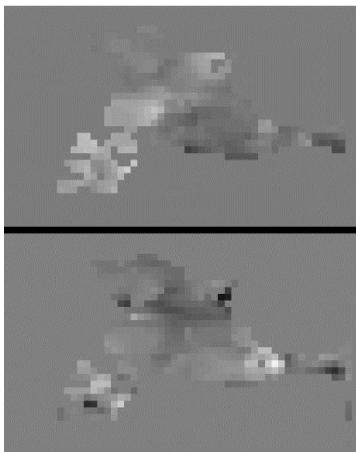


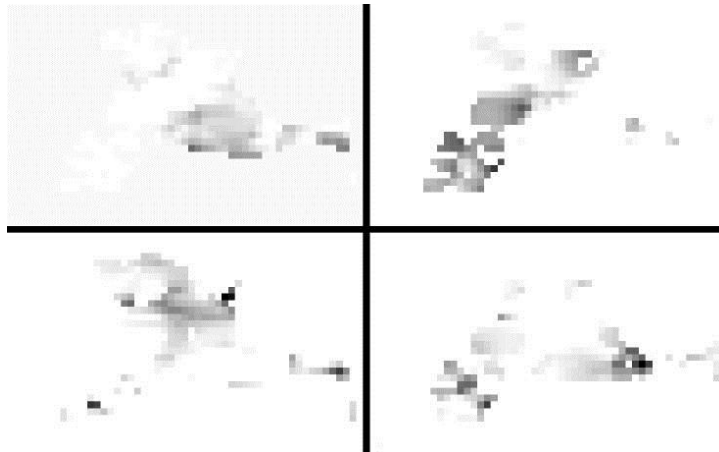
Image frame



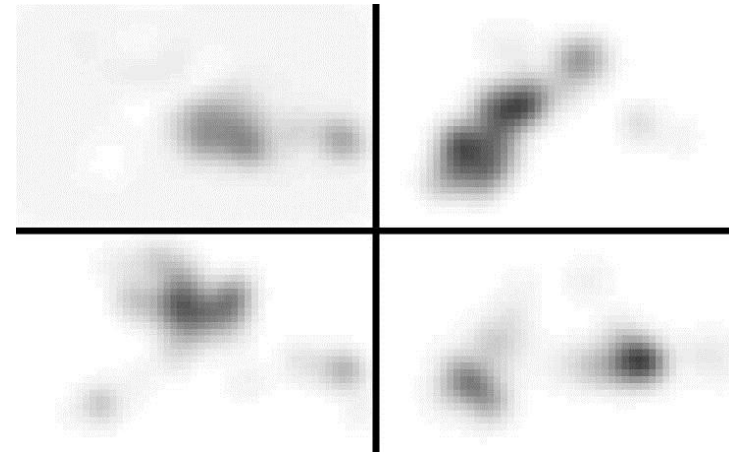
Optical flow $F_{x,y}$



F_x, F_y

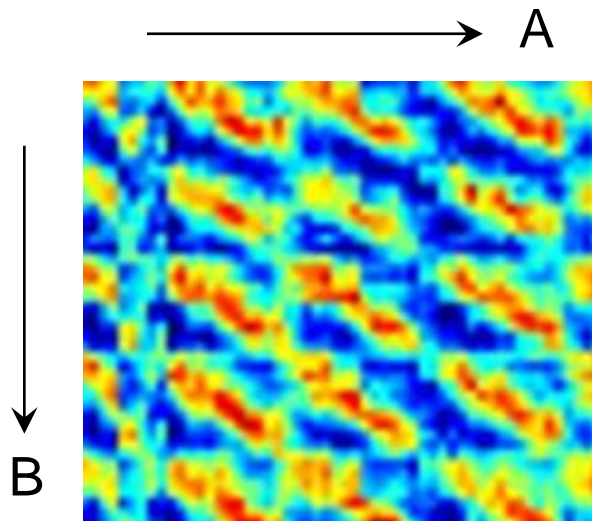
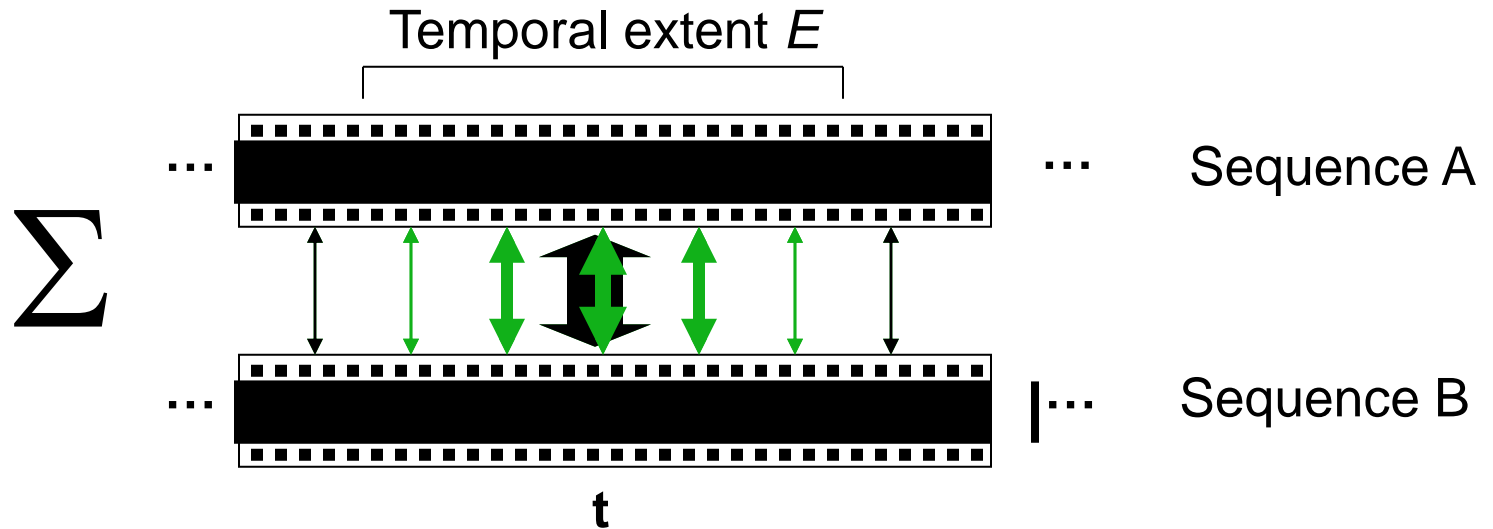


$F_x^-, F_x^+, F_y^-, F_y^+$



blurred $F_x^-, F_x^+, F_y^-, F_y^+$

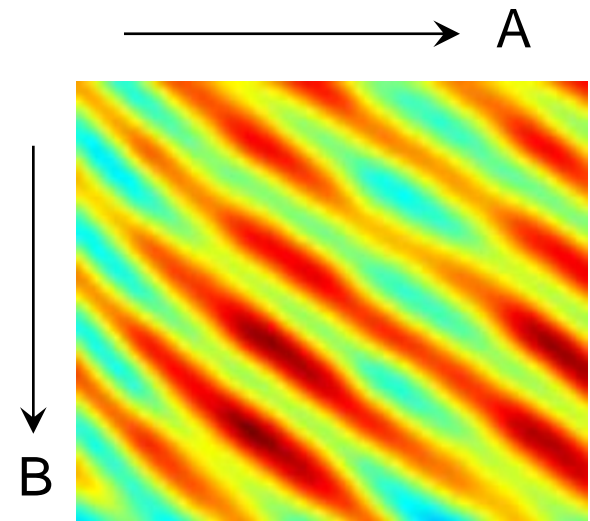
Spatio-Temporal Motion Descriptor



frame-to-frame
similarity matrix



blurry I



motion-to-motion
similarity matrix

Football Actions: matching

Input
Sequence



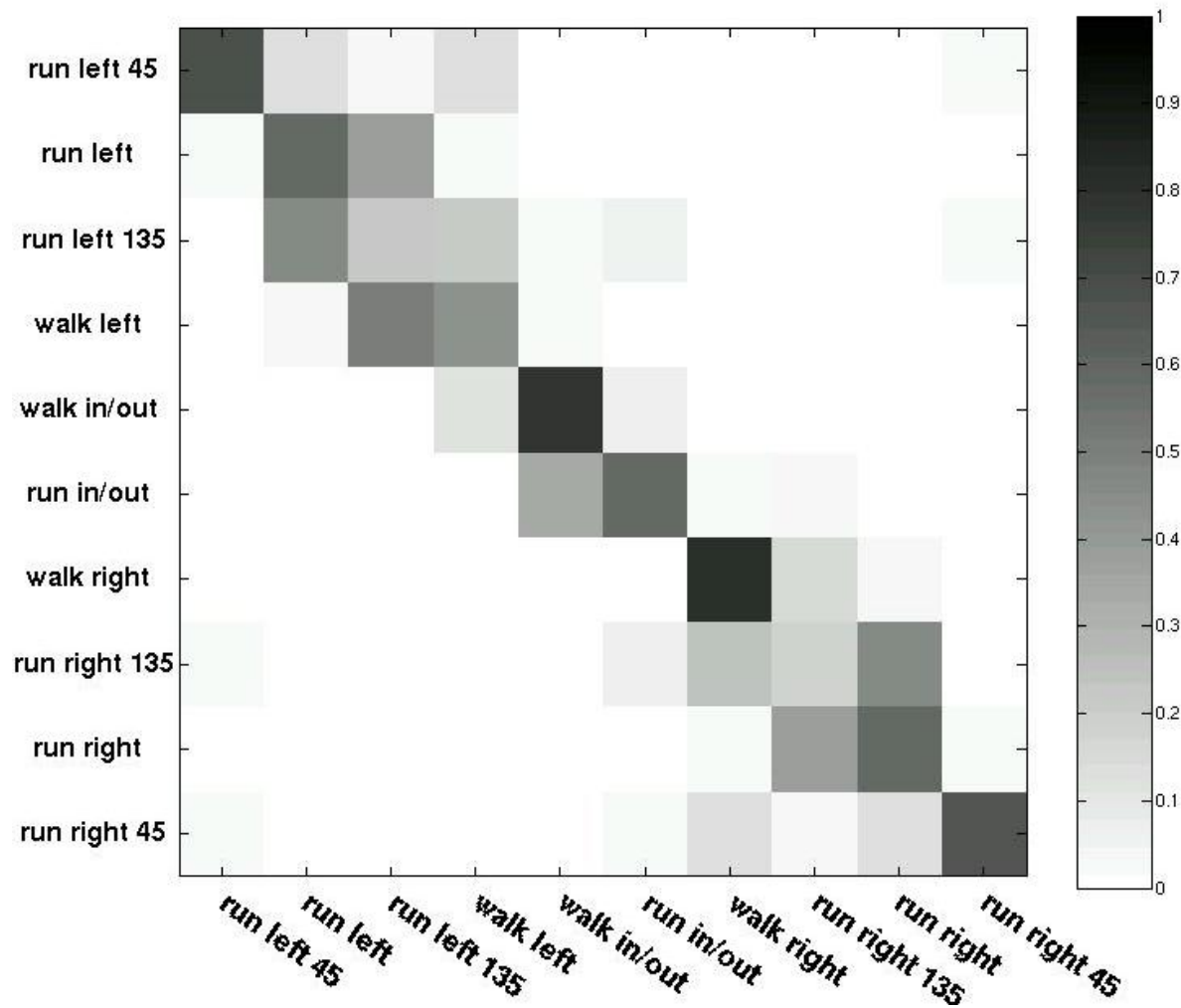
Matched
Frames



input

matched

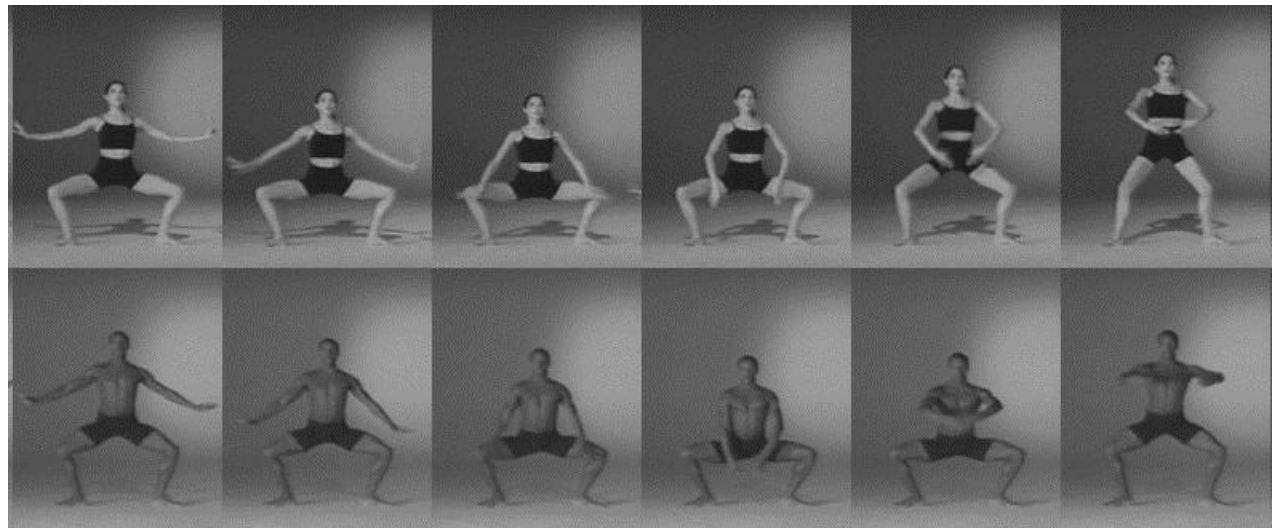
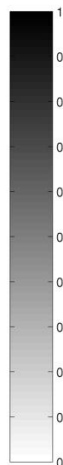
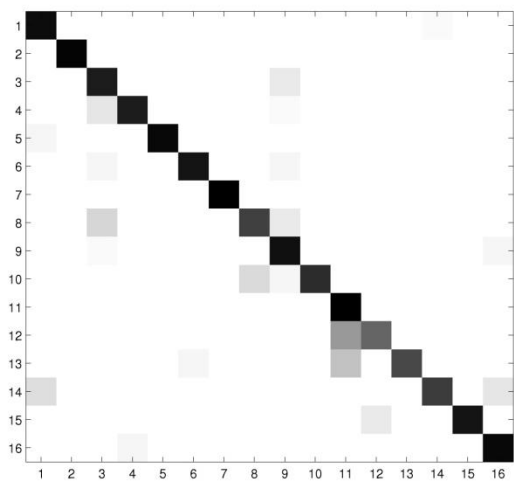
Football Actions: classification



10 actions; 4500 total frames; 13-frame motion descriptor

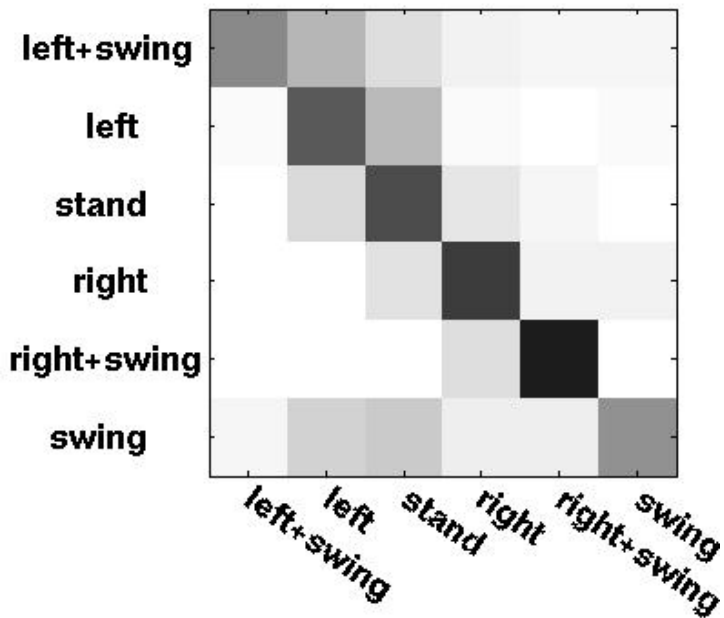
Classifying Ballet Actions

16 Actions; 24800 total frames; 51-frame motion descriptor. Men used to classify women and vice versa.



Classifying Tennis Actions

6 actions; 4600 frames; 7-frame motion descriptor
Woman player used as training, man as testing.

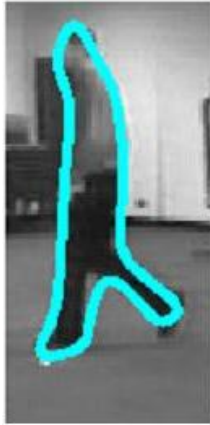


Where are we so far ?



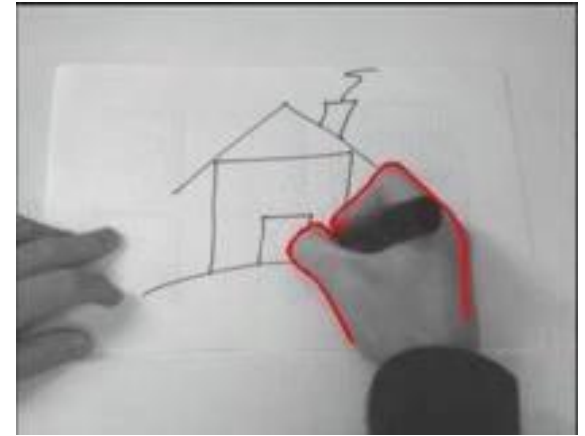
Temporal templates:

- + simple, fast
- sensitive to segmentation errors



Active shape models:

- + shape regularization
- sensitive to initialization and tracking failures



Tracking with motion priors:

- + improved tracking and simultaneous action recognition
- sensitive to initialization and tracking failures

Motion-based recognition:

- + generic descriptors; less depends on appearance
- sensitive to localization/tracking errors

