

# Human Detection and Tracking

CS 543 - D.A. Forsyth

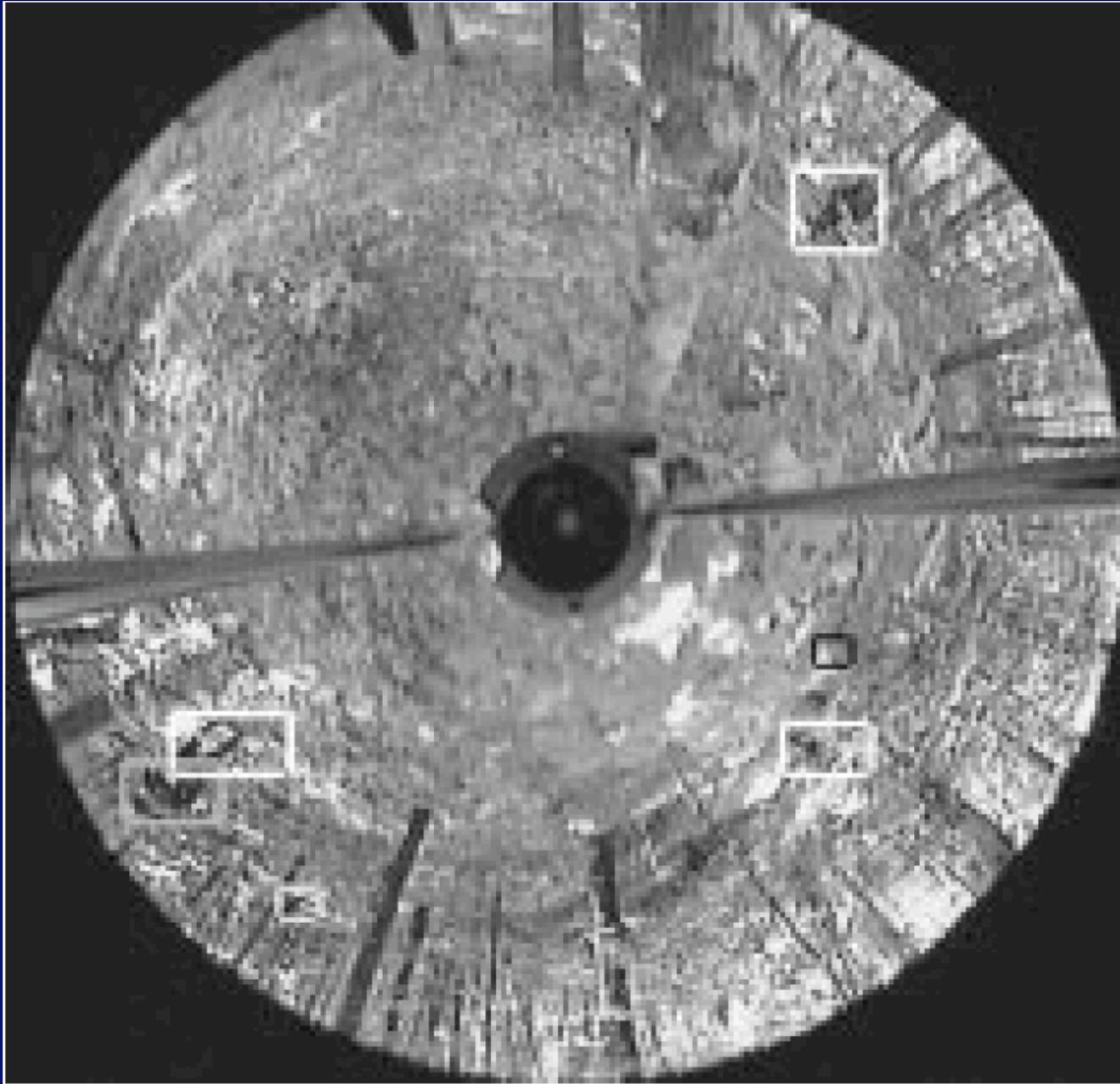
# Why are humans important?

- **Surveillance**
  - prosecution; intelligence gathering; crime prevention
  - HCI; architecture;
- **Synthesis**
  - games; movies;
- **Safety applications**
  - pedestrian detection
- **People are interesting**
  - movies; news

# Surveillance: Where you are can tell what you are doing



Intille et al 95, 97



And can suggest you are  
doing something you  
shouldn't be  
Boult 2001



Bill Freeman flies a magic carpet.

Orientation histograms detect body configuration to control bank, raised arm to fire magic spell.

Freeman et al, 98.



**9** An example of a user playing a Decathlon event, the javelin throw. The computer's timing of the set and release for the javelin is based on when the integrated downward and upward motion exceeds predetermined thresholds.

Motion fields set javelin timing  
Freeman et al 98

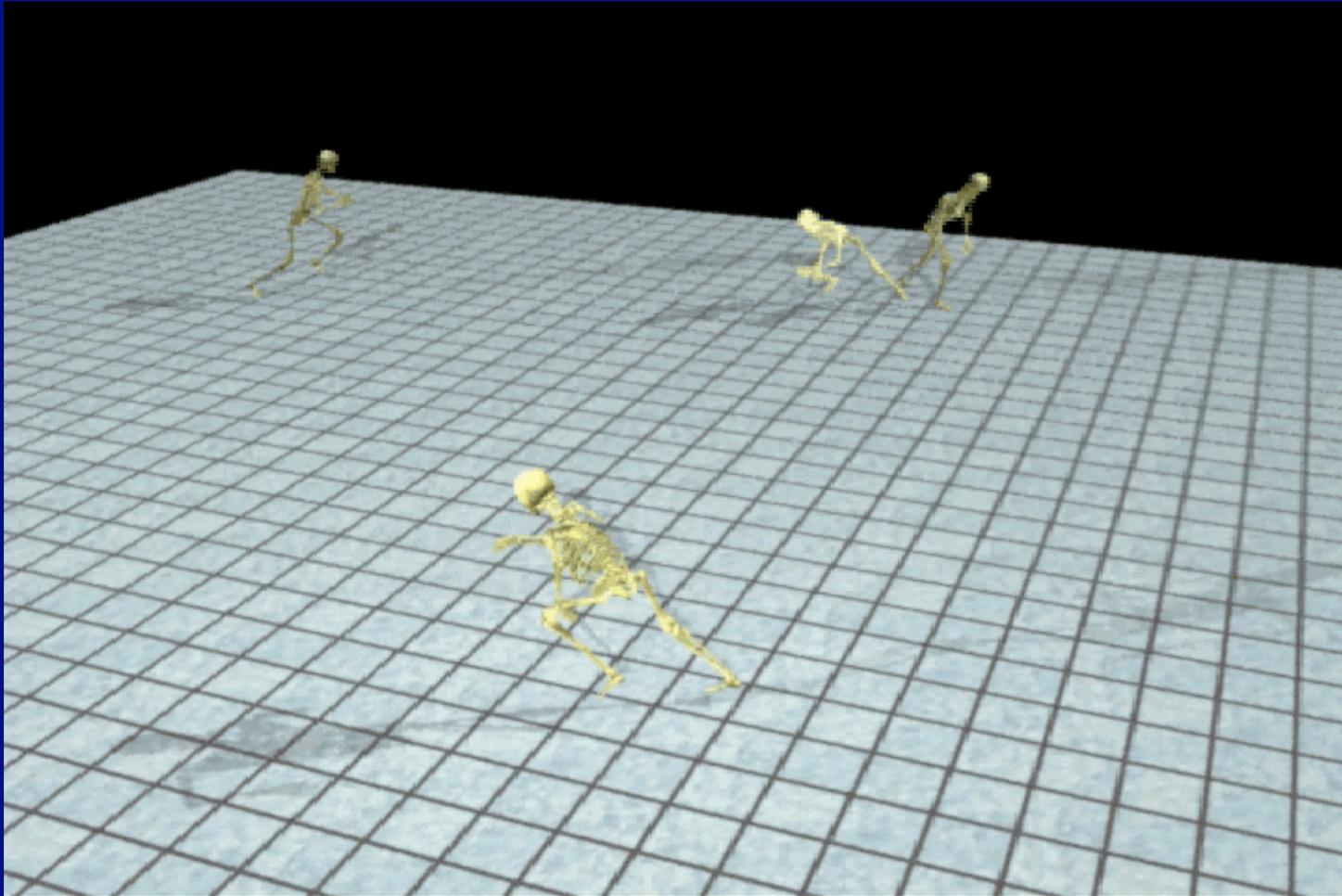


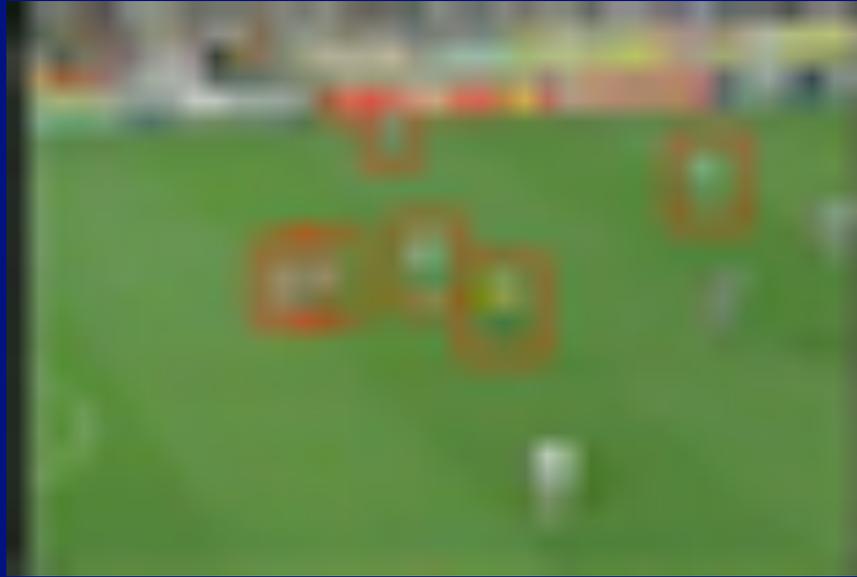
Sony's eyetoy estimates motion fields,  
links these to game inputs.  
Huge hit in EU, well received in US



# Why are humans important?

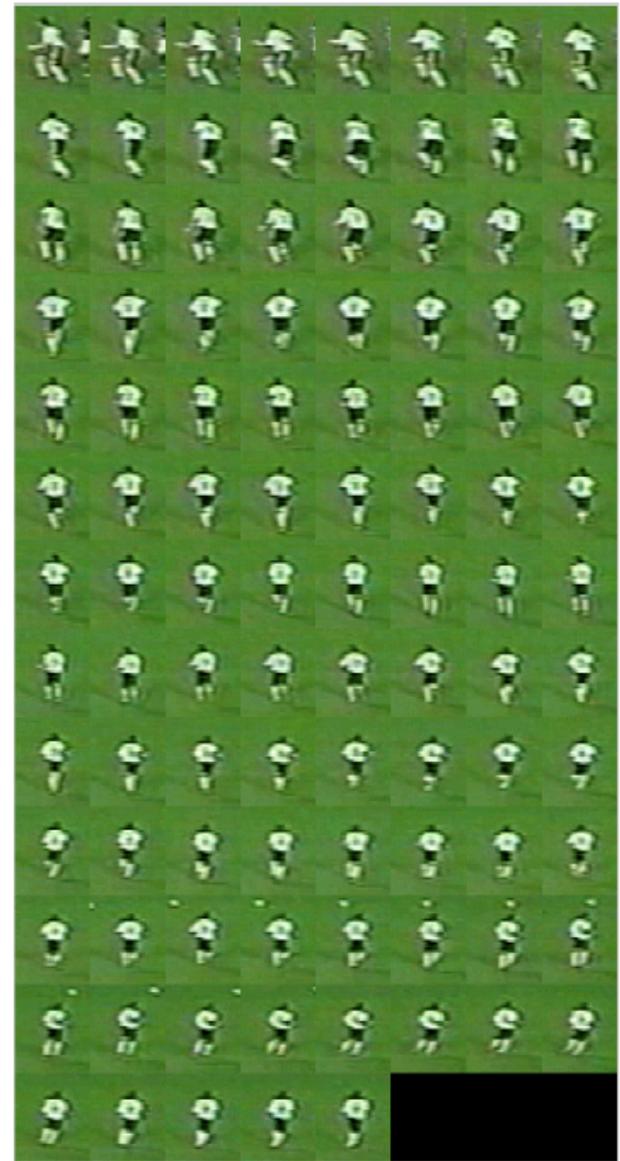
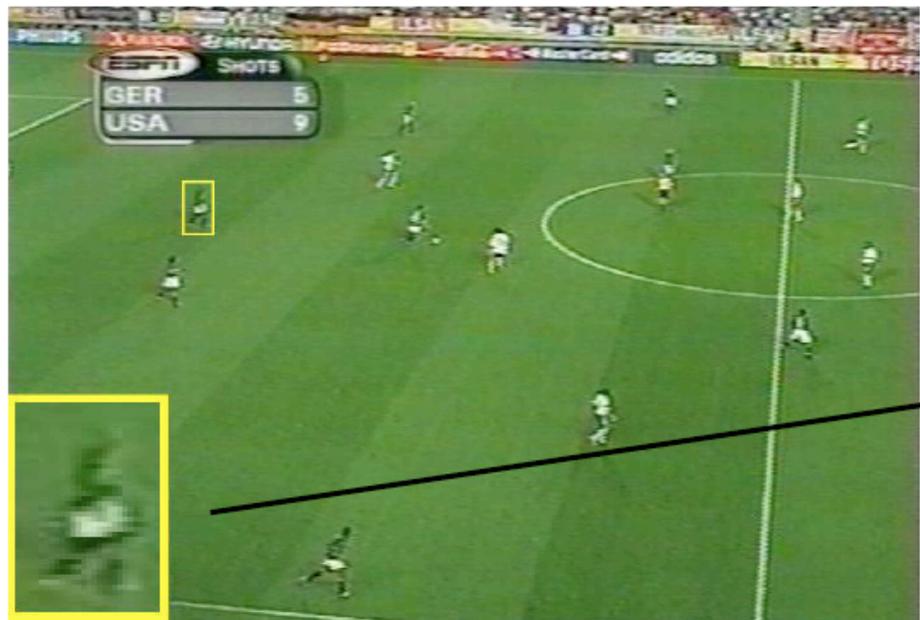
- Surveillance
  - prosecution; intelligence gathering; crime prevention
  - HCI; architecture;
- **Synthesis**
  - games; movies;
- Safety applications
  - pedestrian detection
- People are interesting
  - movies; news





Efros et al, 03

# Motion is a powerful cue at low resolution



Efros et al 03

# Motion Descriptor

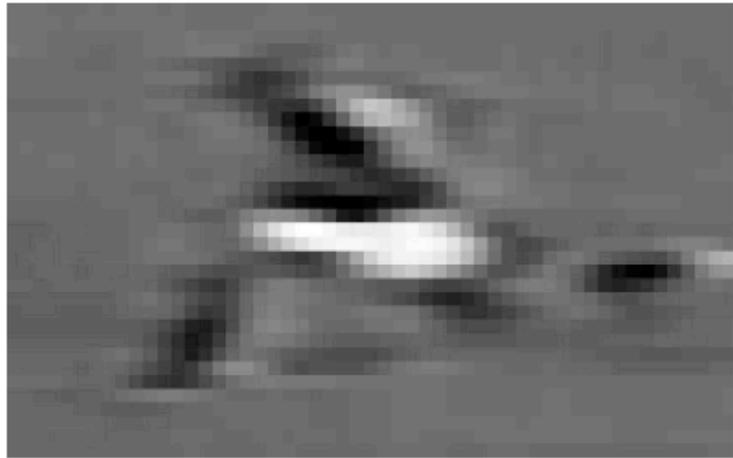
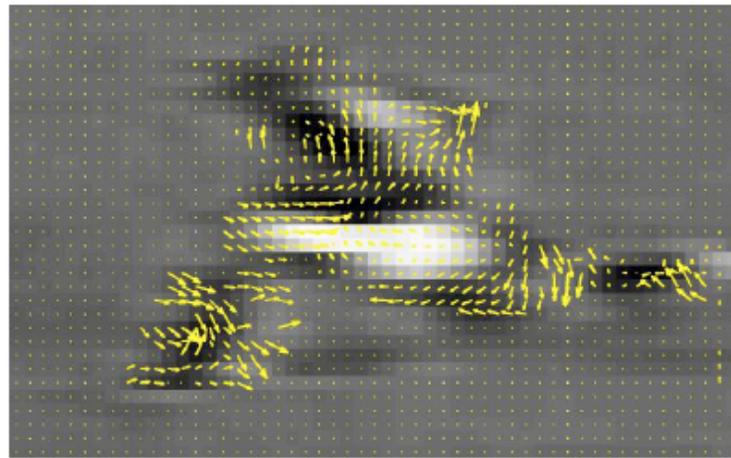
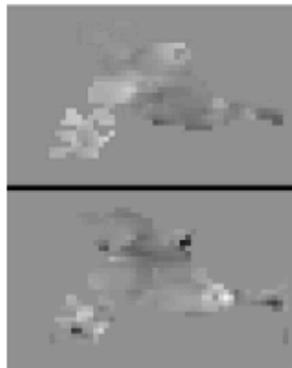


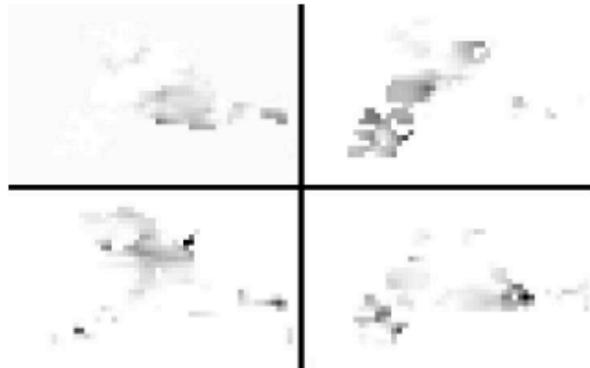
Image frame



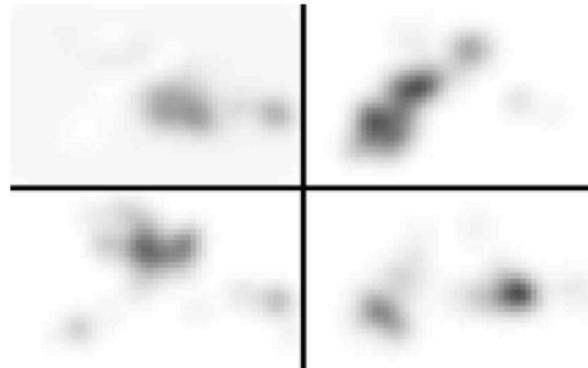
Optical flow



Components

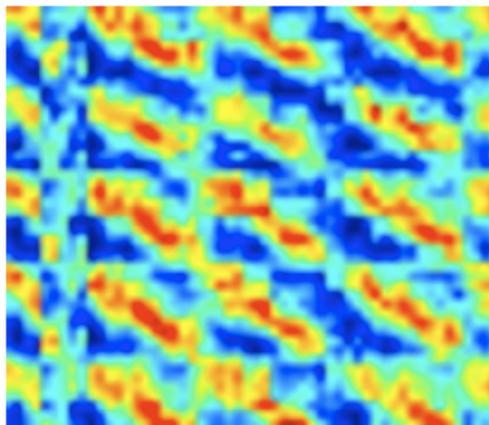
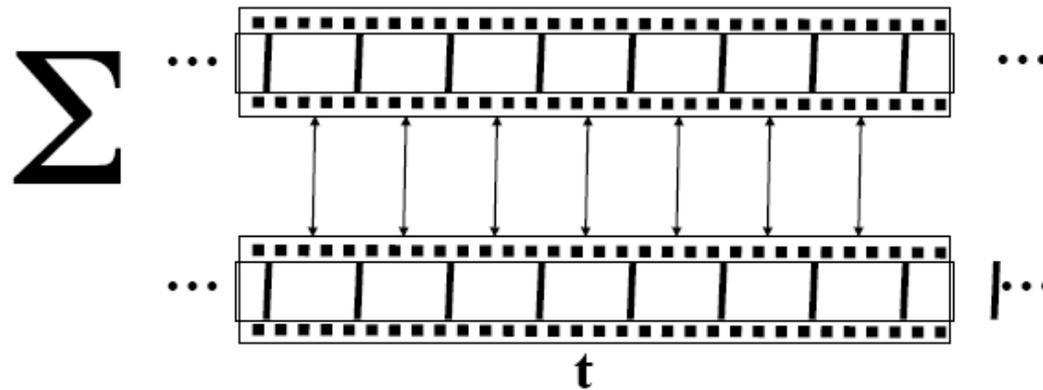


Rectified components

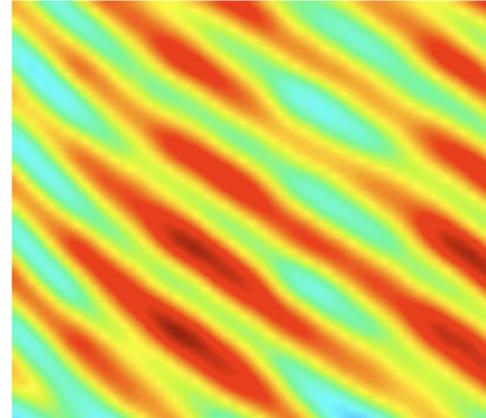


Blurred

# Comparing motion descriptors



frame-to-frame  
similarity matrix



motion-to-motion  
similarity matrix

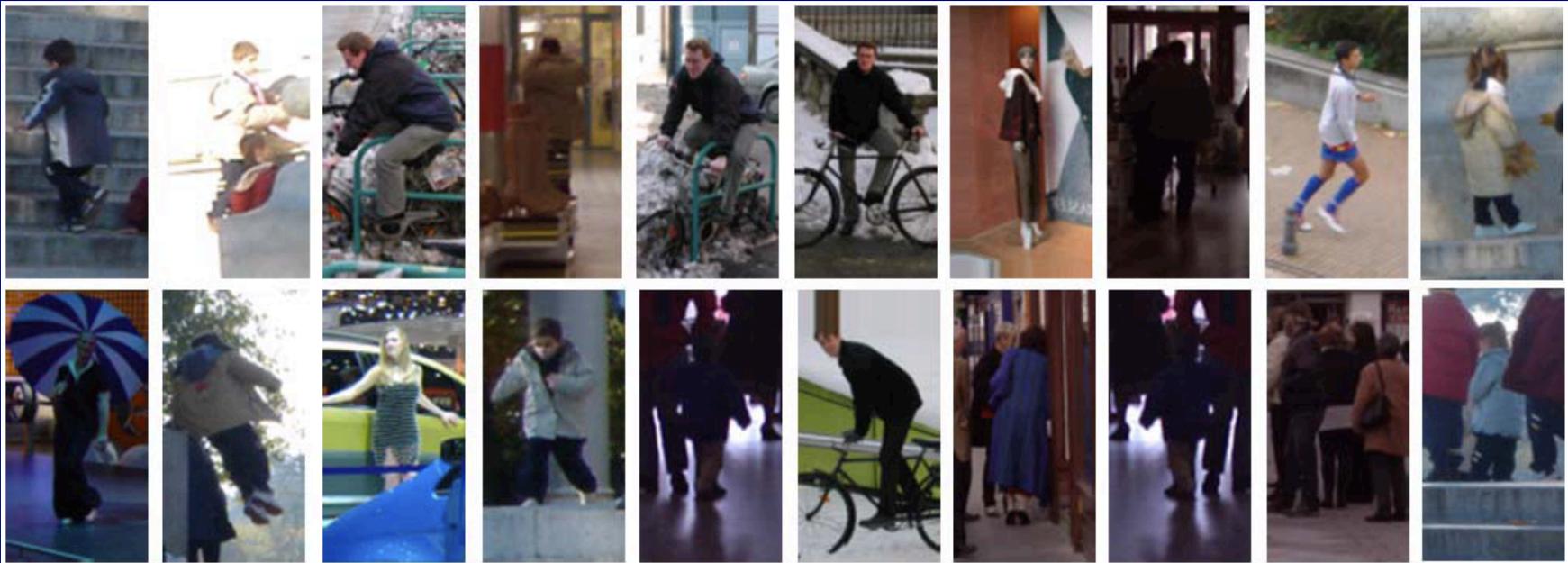


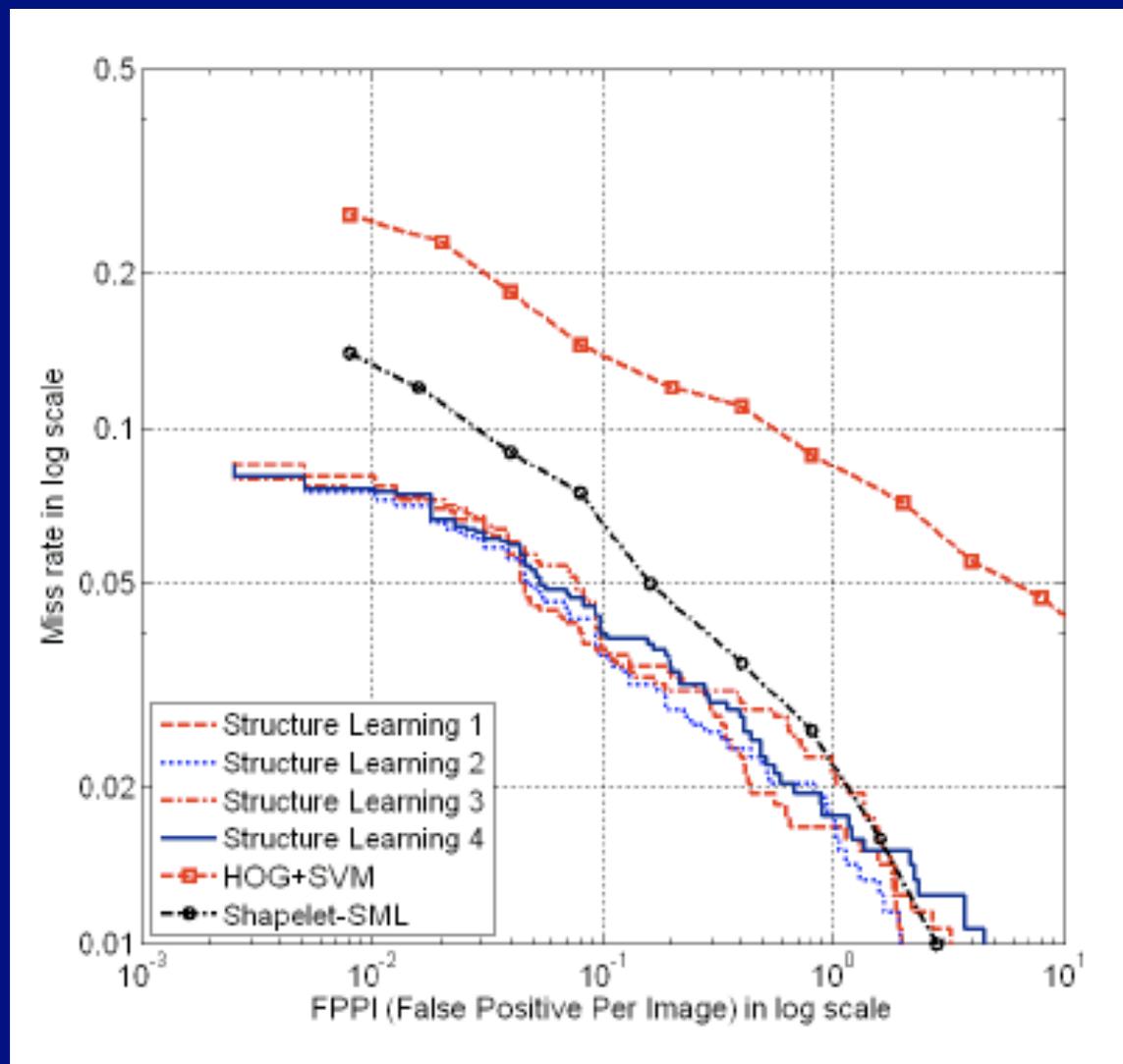
# Why are humans important?

- Surveillance
  - prosecution; intelligence gathering; crime prevention
  - HCI; architecture;
- Synthesis
  - games; movies;
- **Safety applications**
  - pedestrian detection
- People are interesting
  - movies; news

# Example: Pedestrian detection

- Detect pedestrians using object recognition technology (later)
  - Whole body detectors fail on funny configurations
  - Don't know correct criteria for identifying limb segments





Tran + Forsyth 07 vs Dalal+Triggs 05

# Why are humans important?

- Surveillance
  - prosecution; intelligence gathering; crime prevention
  - HCI; architecture;
- Synthesis
  - games; movies;
- Safety applications
  - pedestrian detection
- **People are interesting**
  - movies; news

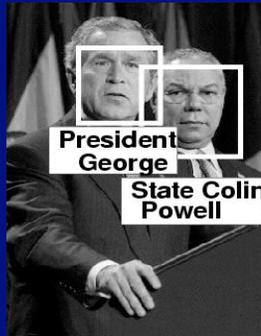
# News Faces

- 5e5 captioned news images
- Mainly people “in the wild”
- Correspondence problem
  - some images have many (resp. few) faces, few (resp. many) names (cf. Srihari 95)
- Process
  - Extract proper names
  - Detect faces (Vogelhuber Schmid 00) 44773 big face responses
  - Rectify faces 34623 properly rectified
  - Kernel PCA rectified faces
  - Estimate linear discriminants
  - Now have (face vector; name\_1, ..., name\_k) 27742 for k ≤ 4
- Apply a form of modified k-means



President George W. Bush makes a statement in the Rose Garden while Secretary of Defense Donald Rumsfeld looks on, July 23, 2003. Rumsfeld said the United States would release graphic photographs of the dead sons of Saddam Hussein to prove they were killed by American troops. Photo by Larry Downing/Reuters





US President George W. Bush (L) makes remarks while Secretary of State Colin Powell (R) listens before signing the US Leadership Against HIV /AIDS , Tuberculosis and Malaria Act of 2003 at the Department of State in Washington, DC. The five-year plan is designed to help prevent and treat AIDS, especially in more than a dozen African and Caribbean nations(AFP/ Luke Frazza)



German supermodel Claudia Schiffer gave birth to a baby boy by Caesarian section January 30, 2003, her spokeswoman said. The baby is the first child for both Schiffer, 32, and her husband, British film producer Matthew Vaughn, who was at her side for the birth. Schiffer is seen on the German television show 'Bet It...?!' ('Wetten Dass...?!') in Braunschweig, on January 26, 2002. (Alexandra Winkler/Reuters)

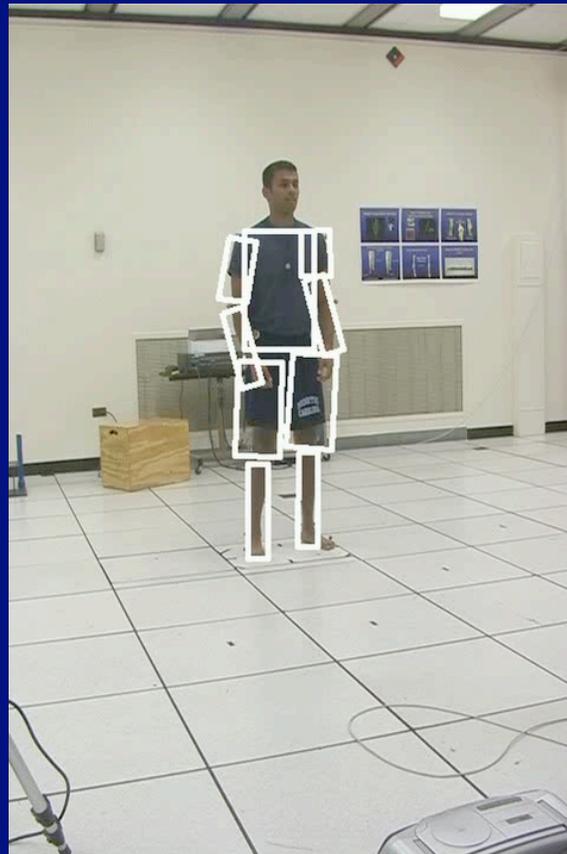


British director Sam Mendes and his partner actress Kate Winslet arrive at the London premiere of 'The Road to Perdition', September 18, 2002. The films stars Tom Hanks as a Chicago hit man who has a separate family life and co-stars Paul Newman and Jude Law. REUTERS/Dan Chung

# Core Problems

- It is not known what needs to be known
  - or, what should we extract from video to do what task?
- It is hard to find people
  - Appearance
  - Aspect
- It is hard to track people in detail
  - Small parts that move fast and unpredictably
- It is hard to describe what they are doing
  - Behaviour composes
    - sometimes in complex ways
  - A canonical vocabulary is not known

# Motion transduction



# Pictorial structures

- For models with the right form, one can test “everything”
  - model is a set of cylindrical segments linked into a tree structure
    - model should be thought of as a 2D template
      - segments are cylinders, so no aspect issue there
      - 3D segment kinematics implicitly encoded in 2D relations
      - easy to build in occlusion
  - putative image segments are quantized
  - => dynamic programming to search all matches
  - Known segment colour - Felzenszwalb-Huttenlocher 00
  - Learned models of colour, layout, texture - Ramanan Forsyth 03, 04

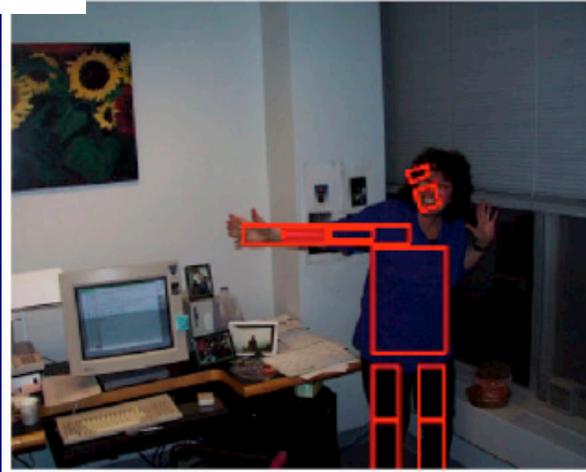
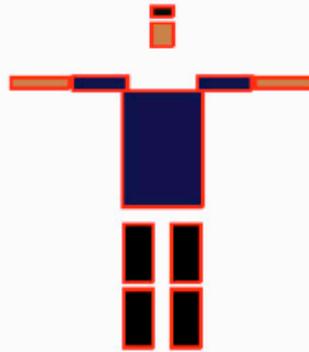
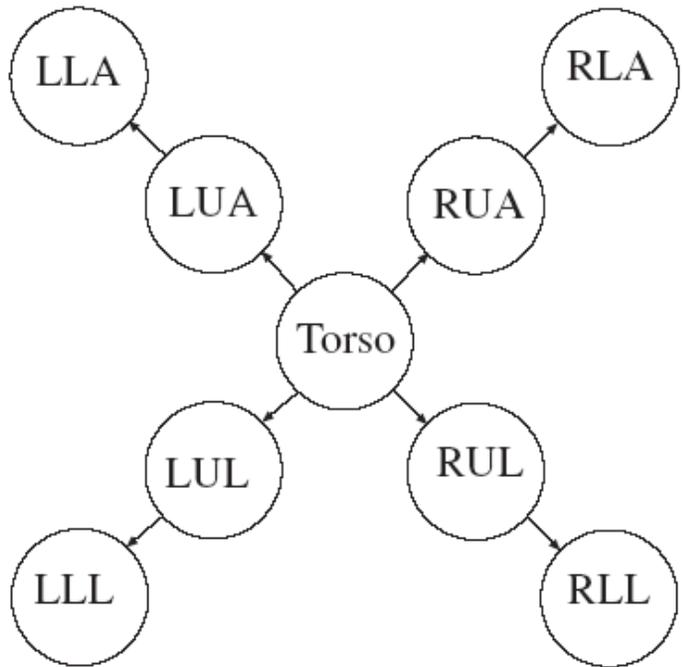


Figure from "Efficient Matching of Pictorial Structures,"  
 P. Felzenszwalb and D.P. Huttenlocher, Proc. Computer Vision and Pattern Recognition  
 2000, c 2000, IEEE as used in Forsyth+Ponce, pp 636, 640

# Human tracking options

- People as blobs (+appearance)
  - Grimson et al 98; Stauffer et al 00; Haritaoglu et al 98, 00; Okuma et al 04
- People as motion fields
  - Bregler 97; Boyd+Little 98
- People as blobs+motion fields
  - Efros et al 03
- Kinematics
  - Hogg 83; Rohr 93; Deutscher et al 00; Toyama+Blake 02; SidenbladhBlackFleet 00; JuBlackYacoob 96; Song Perona 00; etc

# Why is kinematic tracking hard?

- It's hard to detect people
  - until recently, all human trackers were manually started
- People move fast, and can move unpredictably
  - dynamics gives limited constraint on future configuration
  - appearance changes over time (shading, aspect, etc)
- Some body parts are small and tend to have poor contrast
  - particularly difficult to track
    - lower arms (small, fast, look like other things);
    - upper arms (poor contrast)



variation in pose & aspect



self-occlusion & clutter



variation in appearance

# Strategies

- Markov model of (appearance, configuration)
  - 3D Models
    - compare to image
      - variations in dynamical constraints, complexity of inference
        - Hogg 83; Rohr 93; Bregler+Malik 98; Sidenbladh Black Fleet 00; Deutscher Blake Reid 00
  - 2D model
    - Ju Black Yacoob 96; Cham + Rehg 99
- Not quite Markov, but
  - templates encode appearance, then assume markovian dynamics
    - Toyama+Blake 02
- Track by detection
  - Song+Perona (motion) 00; Ioffe+Forsyth (appearance) 01; Mori+Malik (appearance) 02

# Opportunistic detection

People take on a variety of poses, aspects, scales



self-occlusion

rare pose

motion blur

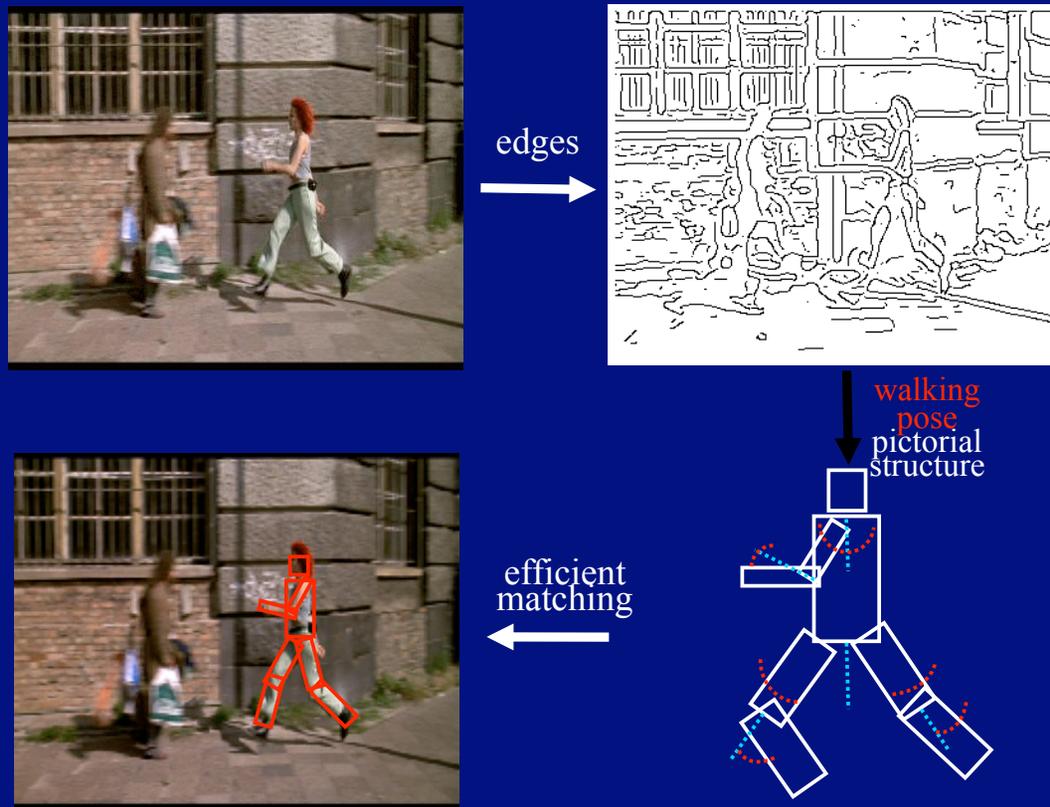


non-distinctive pose

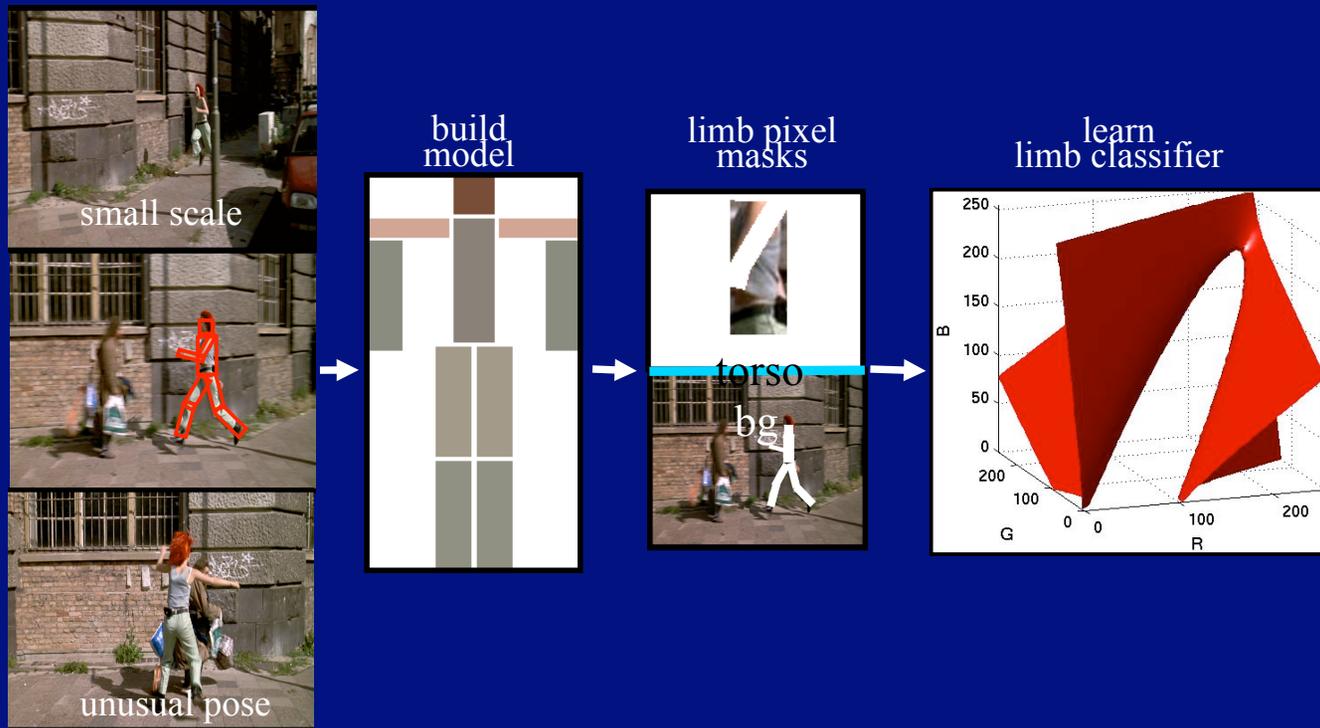
too small

just right  
detect this

# Stylized pose detector



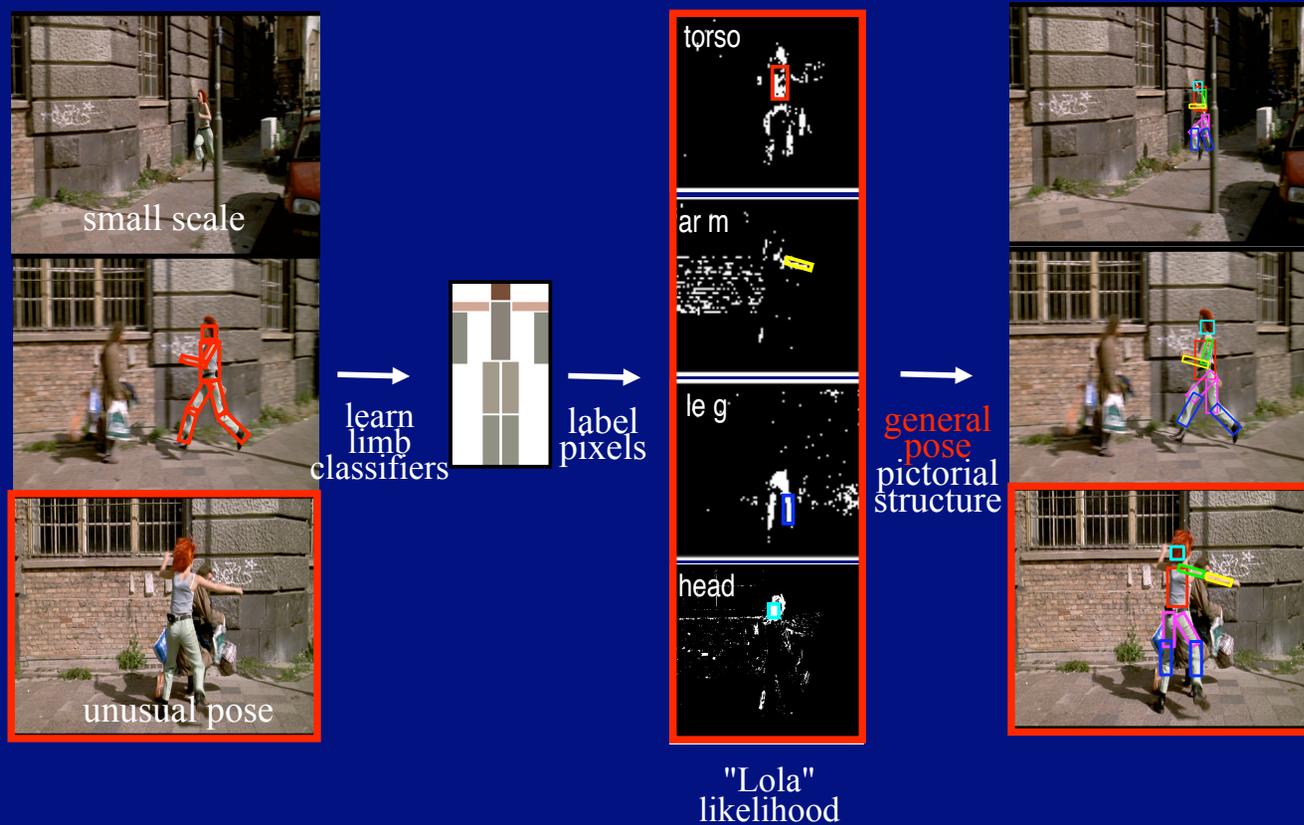
# Model building





Ramanan, Forsyth and Zisserman CVPR05

# Build and detect models



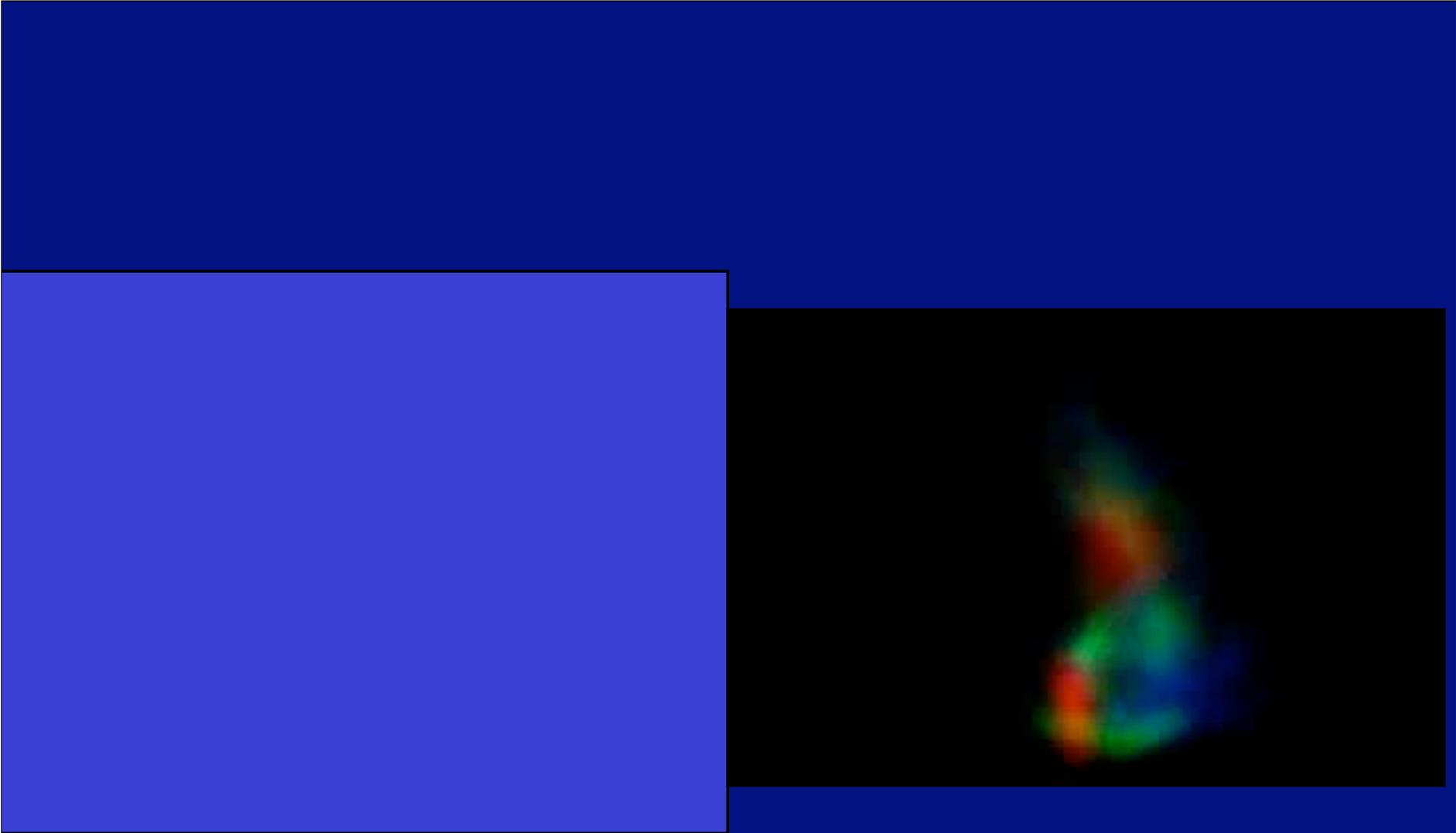


Ramanan, Forsyth and Zisserman CVPR05



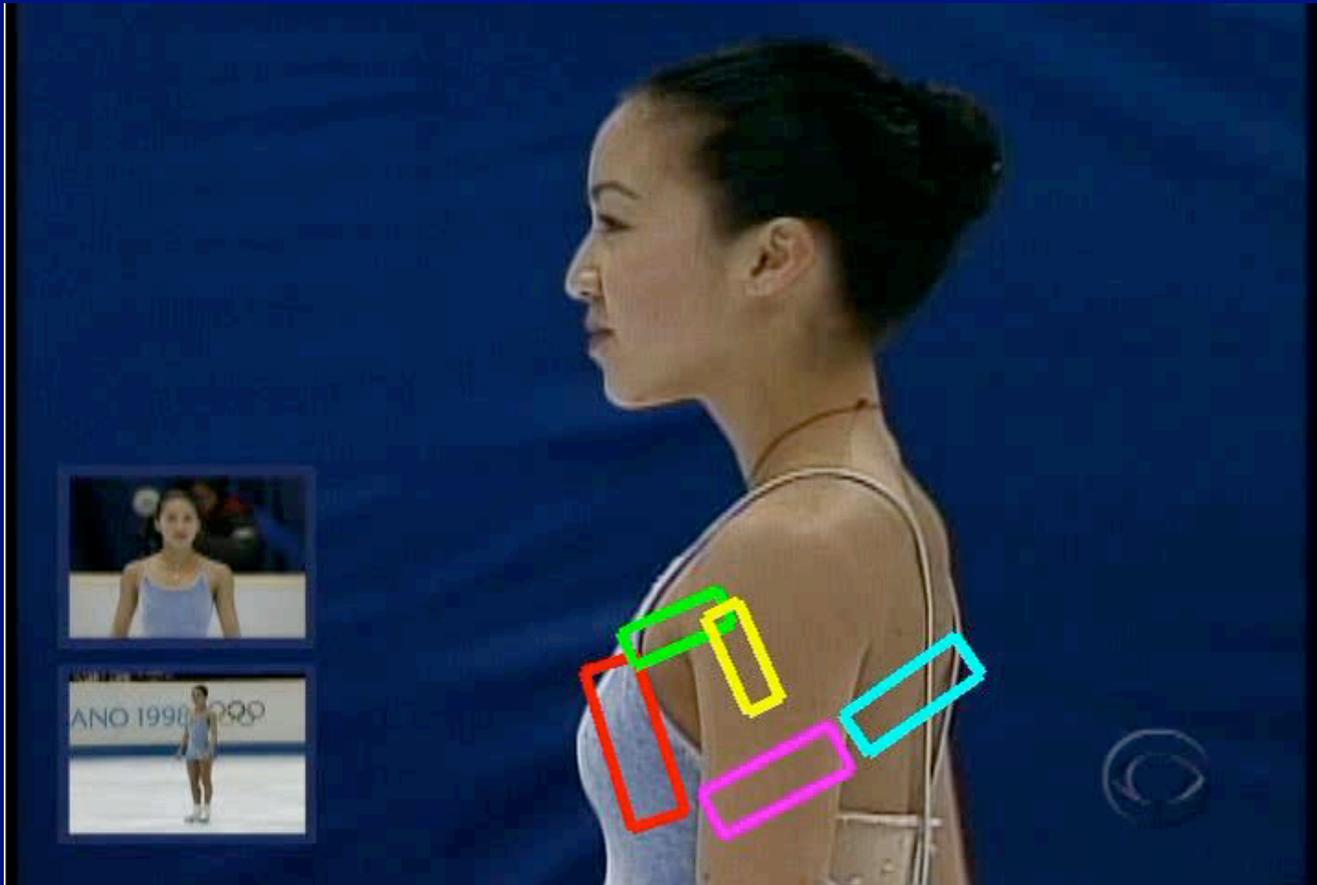
Ramanan, Forsyth and Zisserman CVPR05







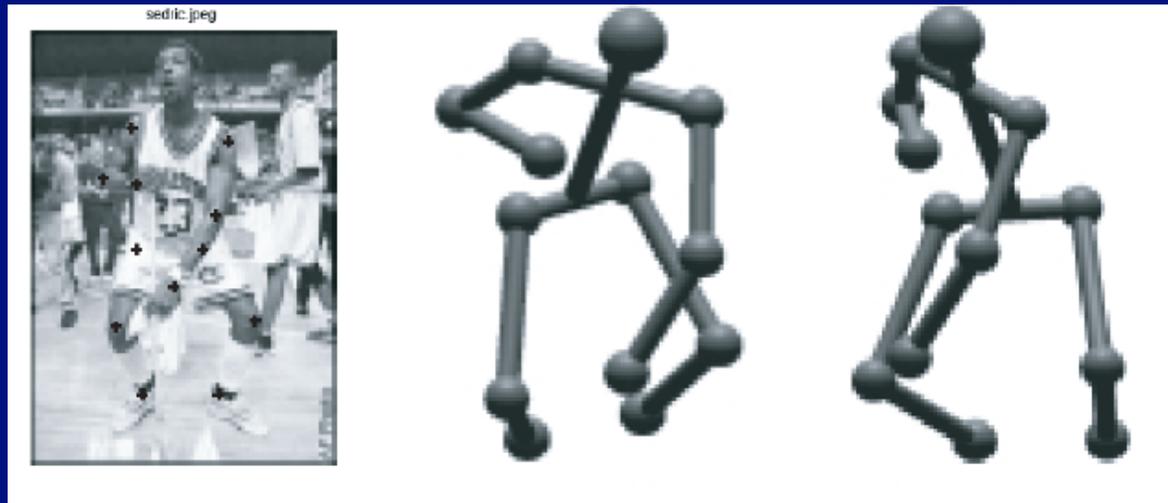
Ramanan, Forsyth and Zisserman CVPR05



Ramanan, Forsyth and Zisserman CVPR05

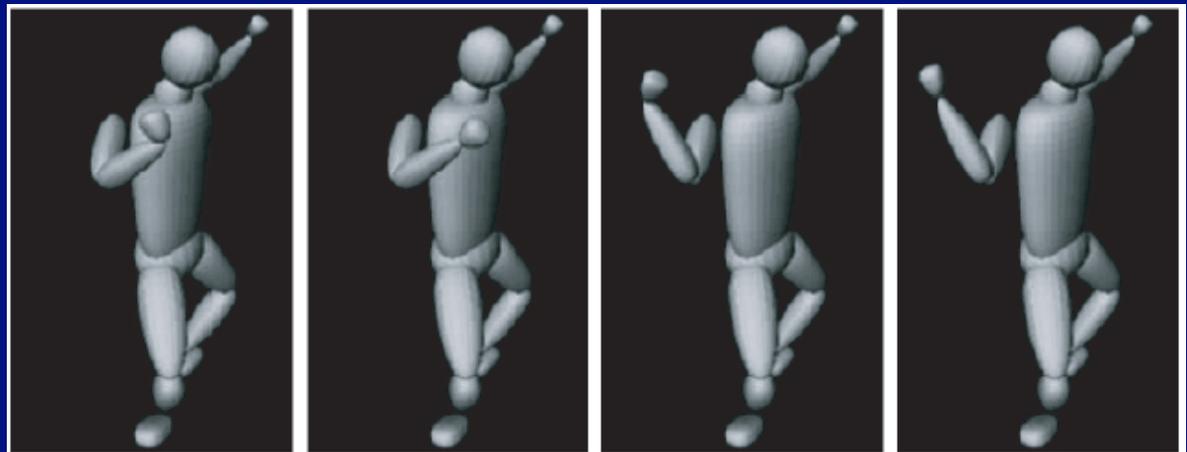
# Lifting

- Infer 3D configuration from image configuration
- Useful for
  - view independent activity recognition
  - user interfaces
  - video motion capture



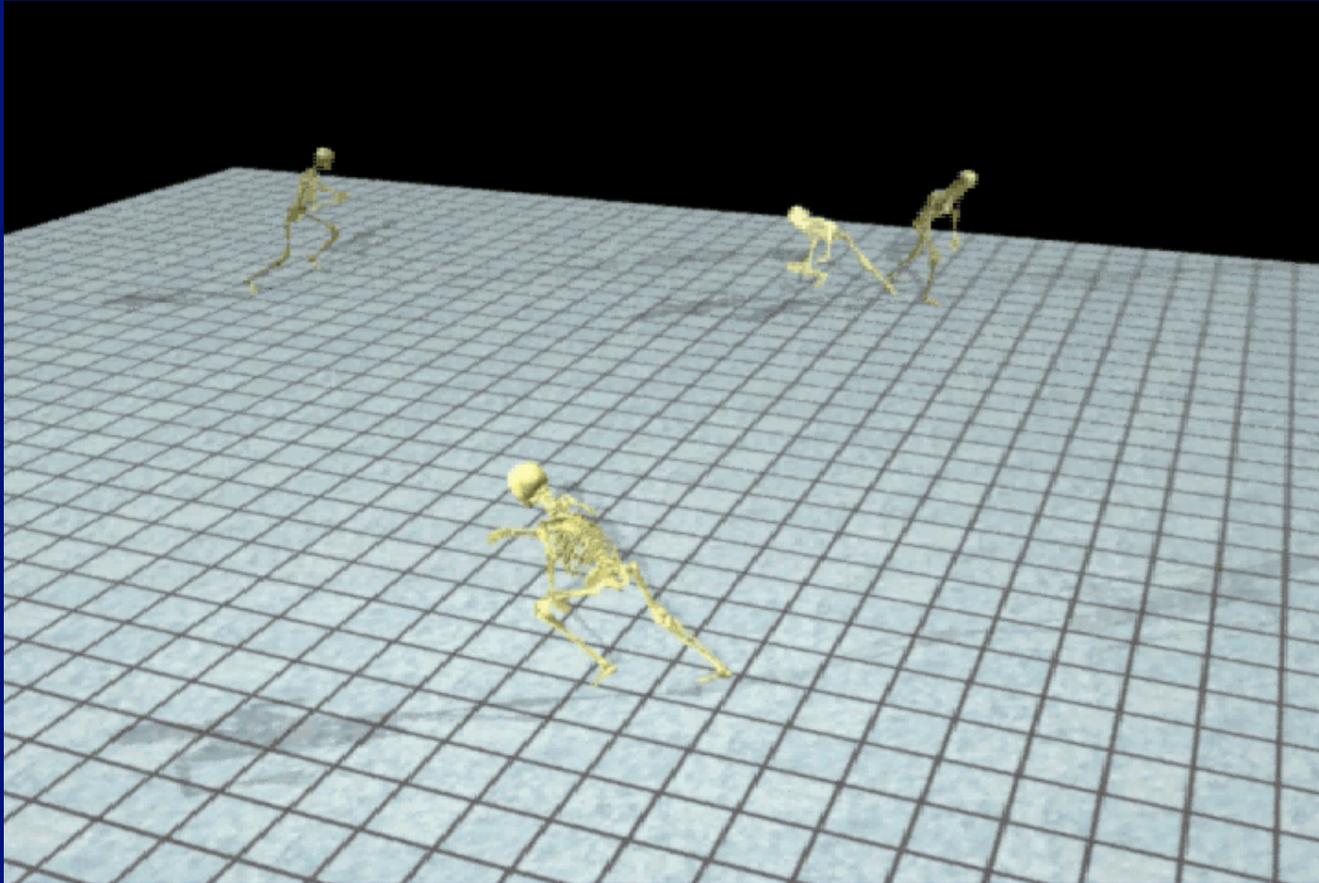
# Ambiguity

- Troubled question
  - lifts are ambiguous (Orthography; Sminchicescu+Triggs 03; etc)
  - but ambiguities
    - can be ignored
      - Taylor 00; Barron+Kakadiaris 00
    - can be dodged
      - Ramanan+Forsyth 03; Howe et al 00
- Summary+musings in Forsyth et al 06



Sminchicescu+Triggs, 03

# Animating people



# Points

- Some properties of motion, illustrated by animation
  - motion composes
    - across time
    - across the body
  - motion can be easy to annotate
    - but good from bad is hard
  - motion clusters well

# Motion synthesis

- Methods
  - By animator
  - By combining observations
    - old tradition of move trees; also (Kovar et al 02, Lee et al 02, Arikan +Forsyth 02, Arikan et al 03, Gleicher et al 03)
  - By physical models
    - old tradition; (Witkin+Kass, 88; Witkin+Popovic 99; Funge et al 88; Fang+Pollard 03, 04)
  - By biomechanical models
    - old tradition; (Liu+Popovic 02; Abe et al 04; Wu+Popovic 03; Liu +Popovic 02)
  - By statistical models
    - old tradition (e.g. Ramsey+Silverman 97); Li et al 02; Safanova et al 04; Mataric et al 99; Mataric 00; Jenkins+Mataric 04;

# Motion graph

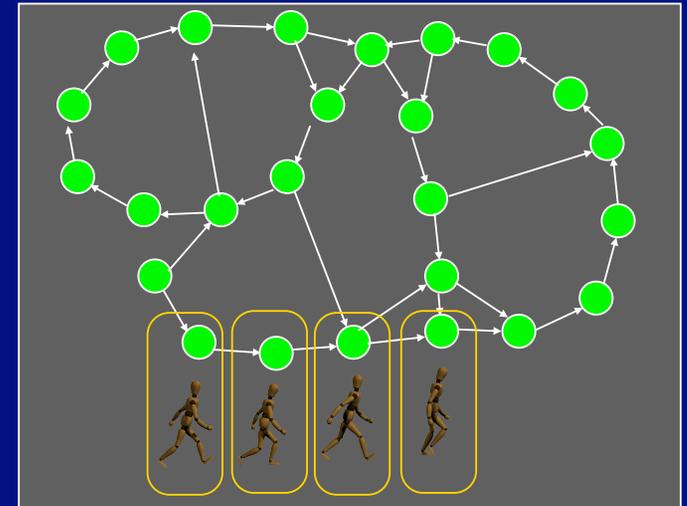
- Take measured frames of motion as nodes
  - from motion capture, given us by our friends
- Edge from frame to any that could succeed it
  - decide by dynamical similarity criterion
  - see also (Kovar et al 02; Lee et al 02)
- A path is a motion
- Search with constraints
  - like root position+orientation, etc.
  - Local (Kovar et al 02)
  - With some horizon
    - Lee et al 02; Ikemoto, Arikan+Forsyth 05
  - Whole path
    - Arikan+Forsyth 02; Arikan et al 03

Motion Graph:

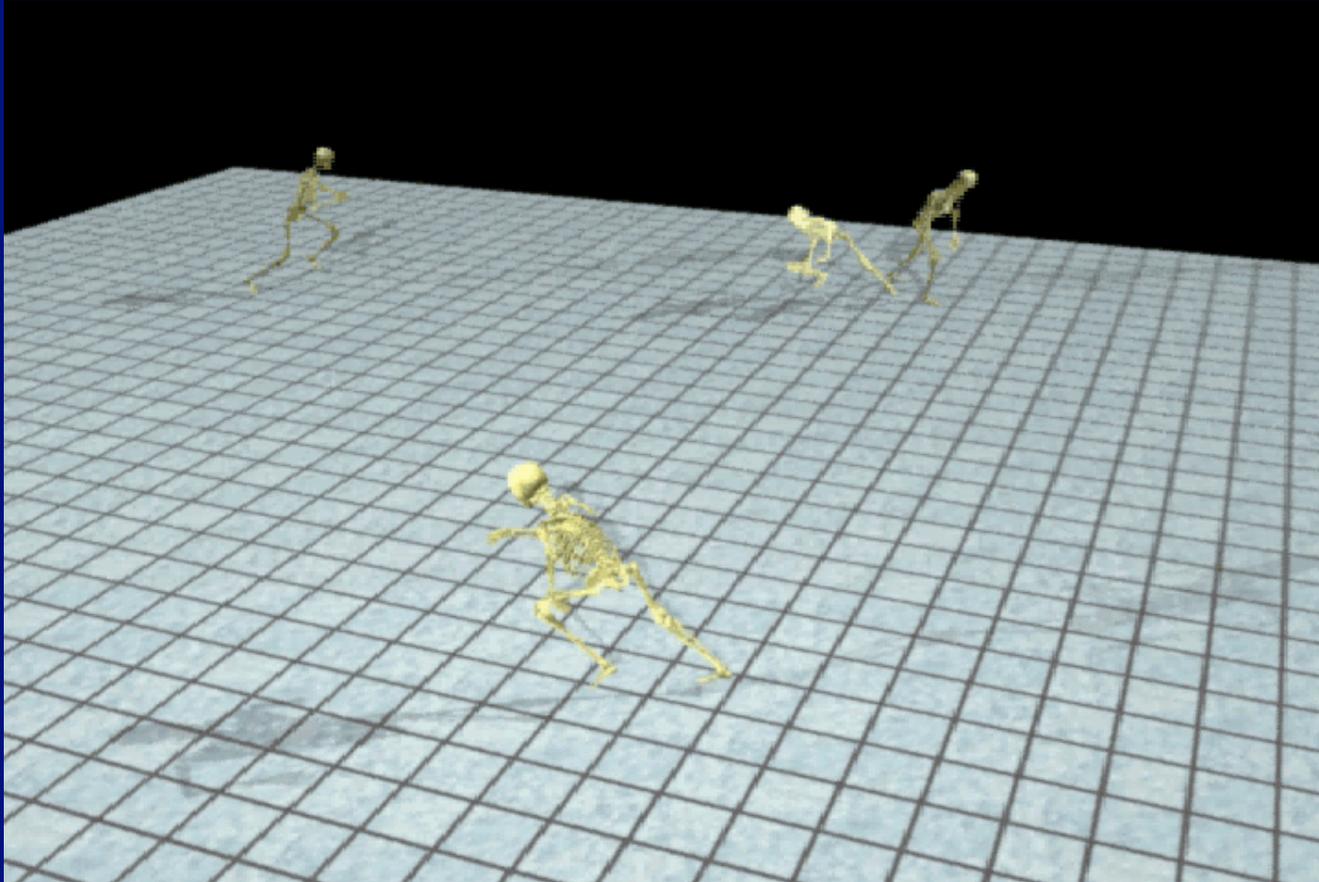
Nodes = Frames

Edges = Transition

A path = A motion

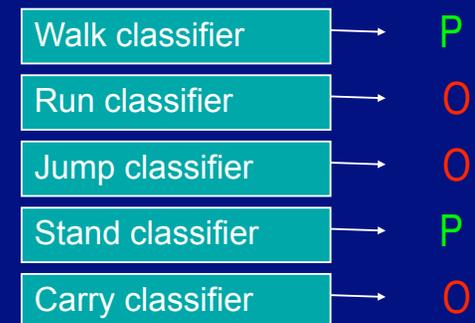






# Annotation - desirable features

- Composability
  - run and wave;
- Comprehensive but not canonical vocabulary
  - because we don't know a canonical vocabulary
- Speed and efficiency
  - because we don't know a canonical vocab.
- Can do this with one classifier per vocabulary item
  - use an SVM applied to joint angles
  - form of on-line learning with human in the loop
  - works startlingly well (in practice 13 bits)



	?	?	?	...	?	$n$ - frames
Walk	P	P	P		P	
Run	●	●	●		●	
Jump	●	●	●		●	
Wave	P	P	○		○	
Carry	●	●	●		●	

Motion demand

## Synthesis by dynamic programming



# Dynamic programming practicalities

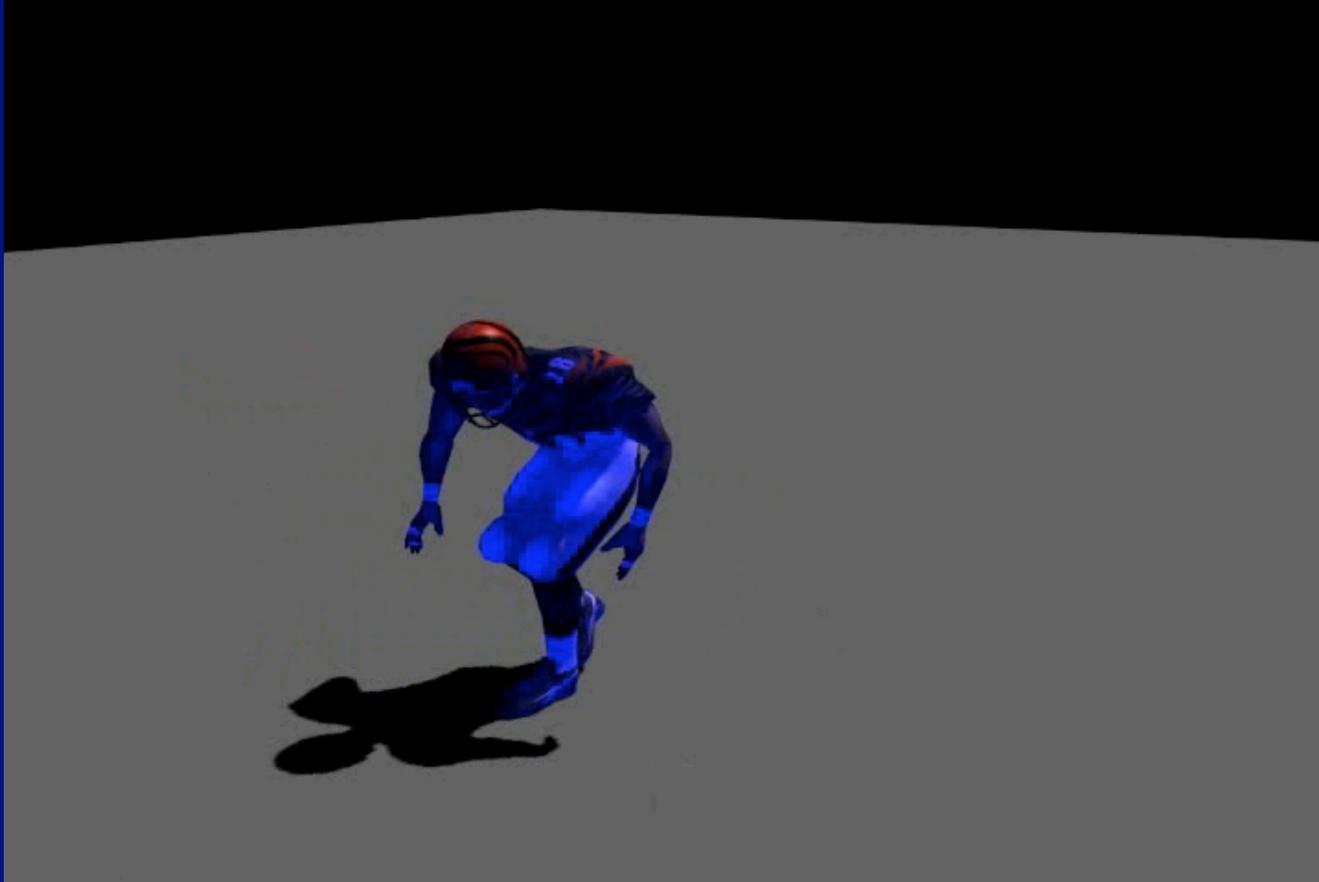
- Scale
  - Too many frames to synthesize
  - Too many frames in motion graph
- Obtain good summary path, refine
  - Form long blocks of motion, cluster
  - DP on stratified sample
    - split blocks on “best” path
    - find similar subblocks
      - DP on this lot
        - etc. to 1-frame blocks



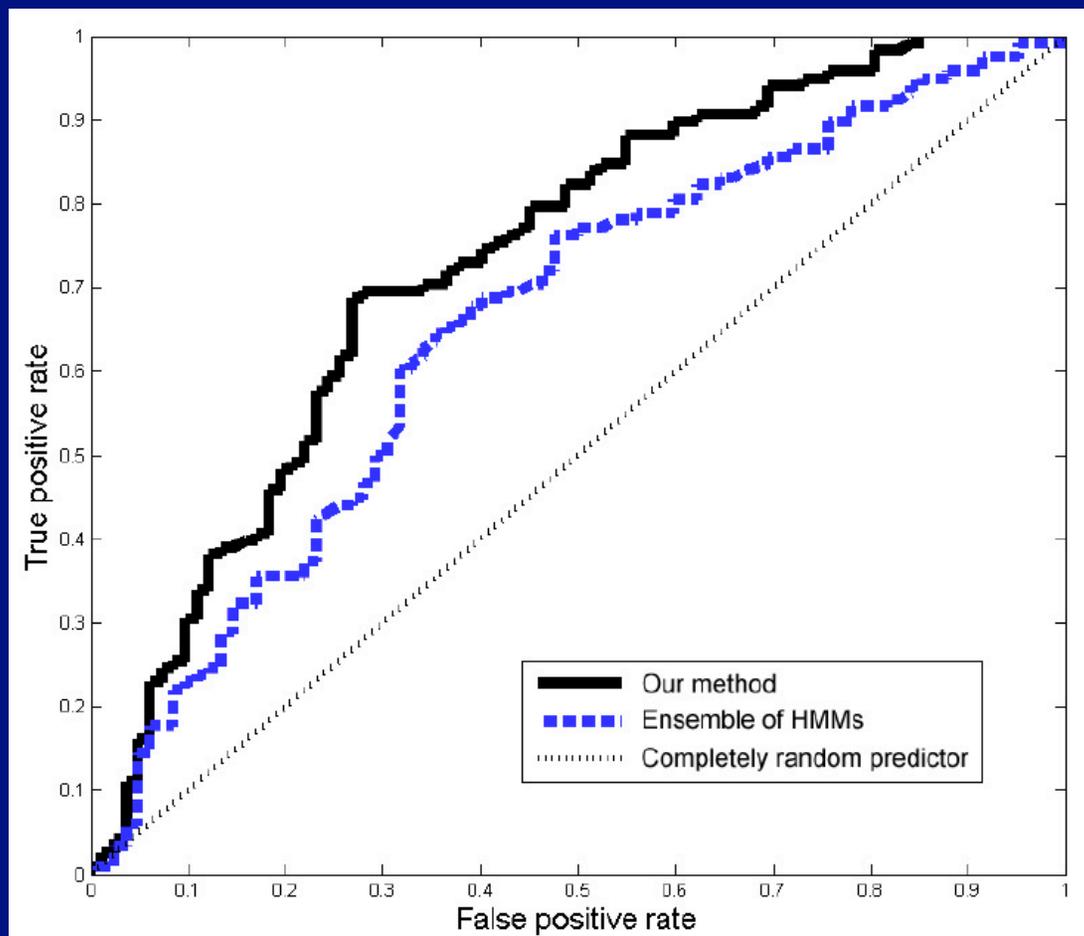
# Transplantation

- Motions clearly have a compositional character
  - Why not cut limbs off some motions and attach to others?
    - we get some bad motions
  - build a classifier to tell good from bad
    - avoid foot slide by leaving lower body alone



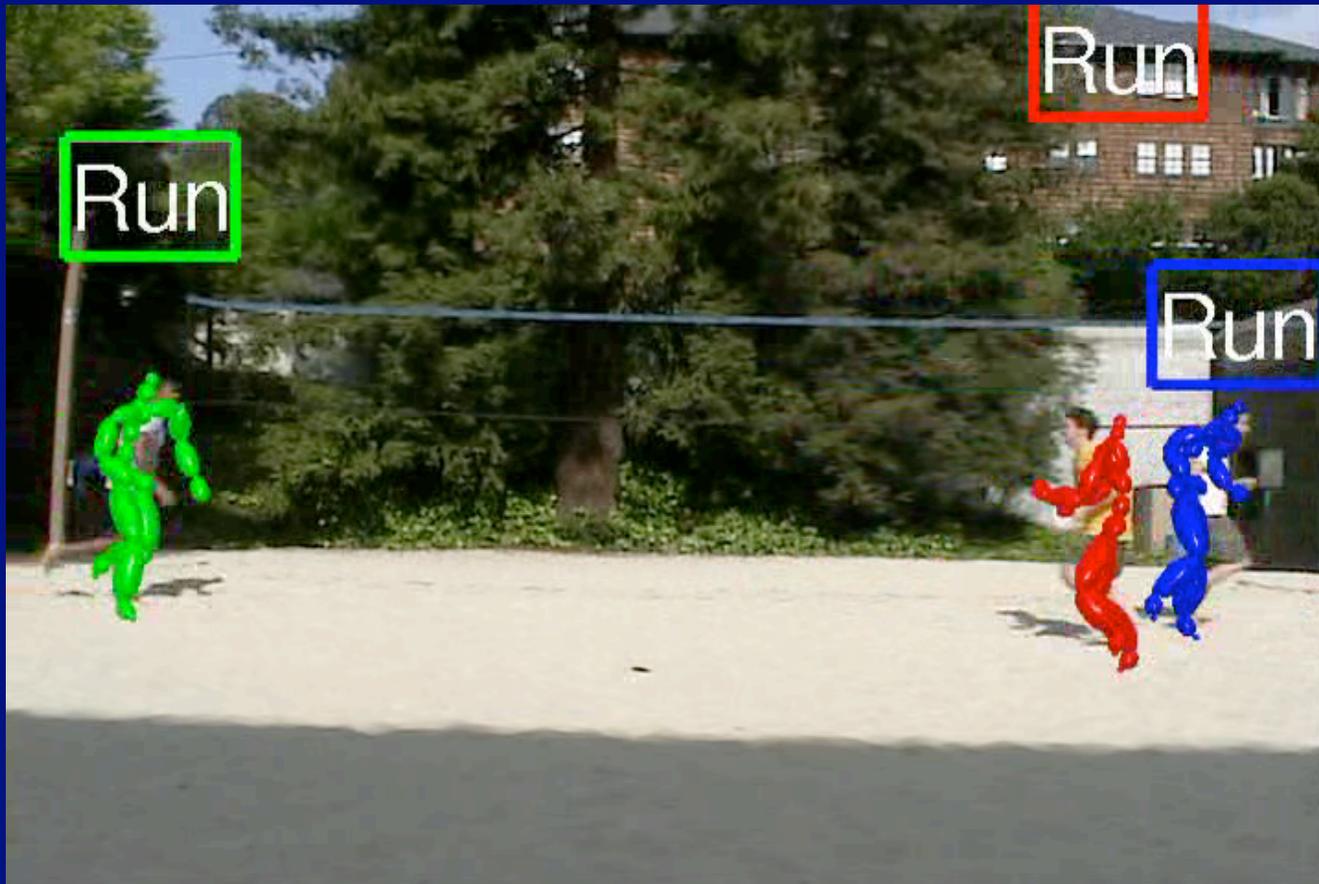


# It is hard to tell good from bad automatically



cf Ren et al 05 for HMM's

# Activity recognition



# Themes

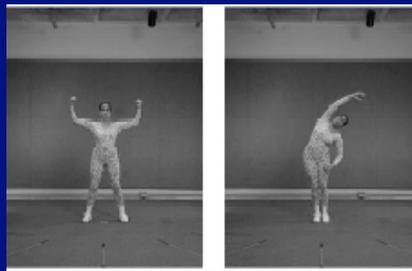
- Activity recognition has important special properties
  - No taxonomy - the structure of categories is hard, not well understood
  - Activity composes in complex ways
- Hence, most activities have no name
  - we're not really tagging videos "run" vs. "walk"
- Signal representations should
  - respect composition
  - be comparative

# Composition and Activity

- Composition is an important source of complexity
  - (flexibility for planning, control)
- We can join motions up in time to make new motions
  - The process is now quite well understood
  - Good quality can be obtained
  - Useful in animation
- We can join up parts of motion across the body
  - But it doesn't always work (and we don't know why, really)

# Naming activities

- Absence of a canonical vocabulary is a serious problem
  - strategies
    - adopt specialized domains (Bobick+Davis 01, Efros et al 03)
    - guess a vocabulary (Efros et al 03)
    - match motion to motion and avoid the issue (Efros et al 03)
    - use vocab useful for synthesis (Ramanan et al 03)



Bobick & Davis. PAMI01

# Activity recognition

- **By comparison to labelled data**
  - benefit from temporal smoothing
    - aka motion synthesis
- **By inference on a generative model**
  - so we can search for activities without having ever seen them
    - composition over body and space
- **By discriminative method**
  - transfer learning by feature construction deals with
    - aspect
    - shortage of training data

# Annotating observations by synthesis

