

# Activity, Appearance, Aspect and Attributes

D.A. Forsyth, UIUC

with

Okan Arikan (UT Austin), Nazli Ikizler (Boston U), Leslie Ikemoto (animate-me), Derek Hoiem (UIUC), Ali Farhadi (UIUC), Ian Endres (UIUC), Ryan White (Euclid media)

Obtain dataset

Build features

Mess around with classifiers, probability, etc

Produce representation

Computer vision

Obtain dataset

Build features

---

Light entertainment  
(the way we do it)

Mess around with classifiers, probability, etc

---

Computer vision

Produce representation

# Big questions

Computer vision

- What signal representation should we use for activity recognition?
    - Compare
      - Appearance (do not segment bits and pieces explicitly)
      - Kinematic (segment bits and pieces explicitly)
- 
- 

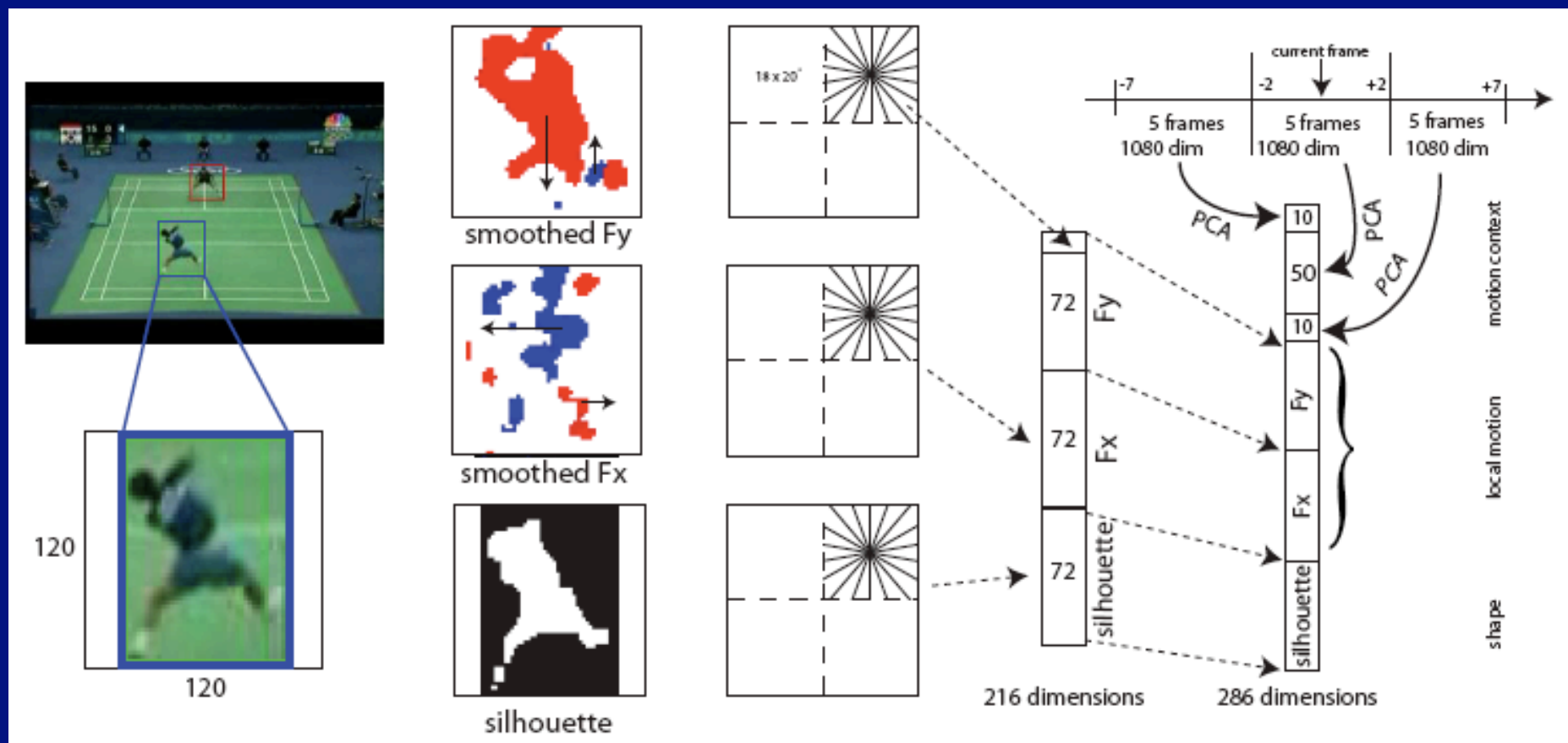
Computer vision

- Does a discriminative framework make any sense? For activity
  - Compare
    - Walk; run; etc
    - (rather vague)

# Appearance features

- Less nasty segmentation
  - (body from background, perhaps not even that)
    - Spatio-temporal volumes
      - (e.g. Davis+Bobick 97; Blank et al 05)
    - Motion trends/flow fields
      - (e.g. BobickDavis 96; Davis 01; Efros et al 03; Laptev+Perez 07; Laptev et al 08)
    - Spatio-temporal interest points
      - (e.g. Niebles et al 06; 08; Scovanner et al 07)
    - Various mixtures of these

# An Appearance feature



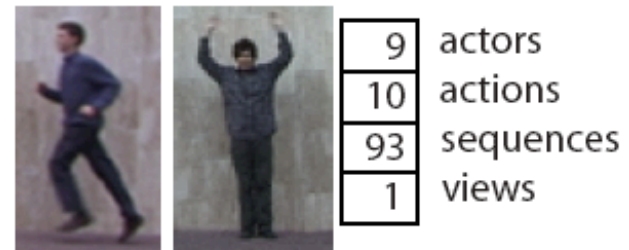
Tran and Sorokin 08, after Duygulu and Ikizler 07

# Datasets

IXMAS



Weizman



Our dataset



UMD



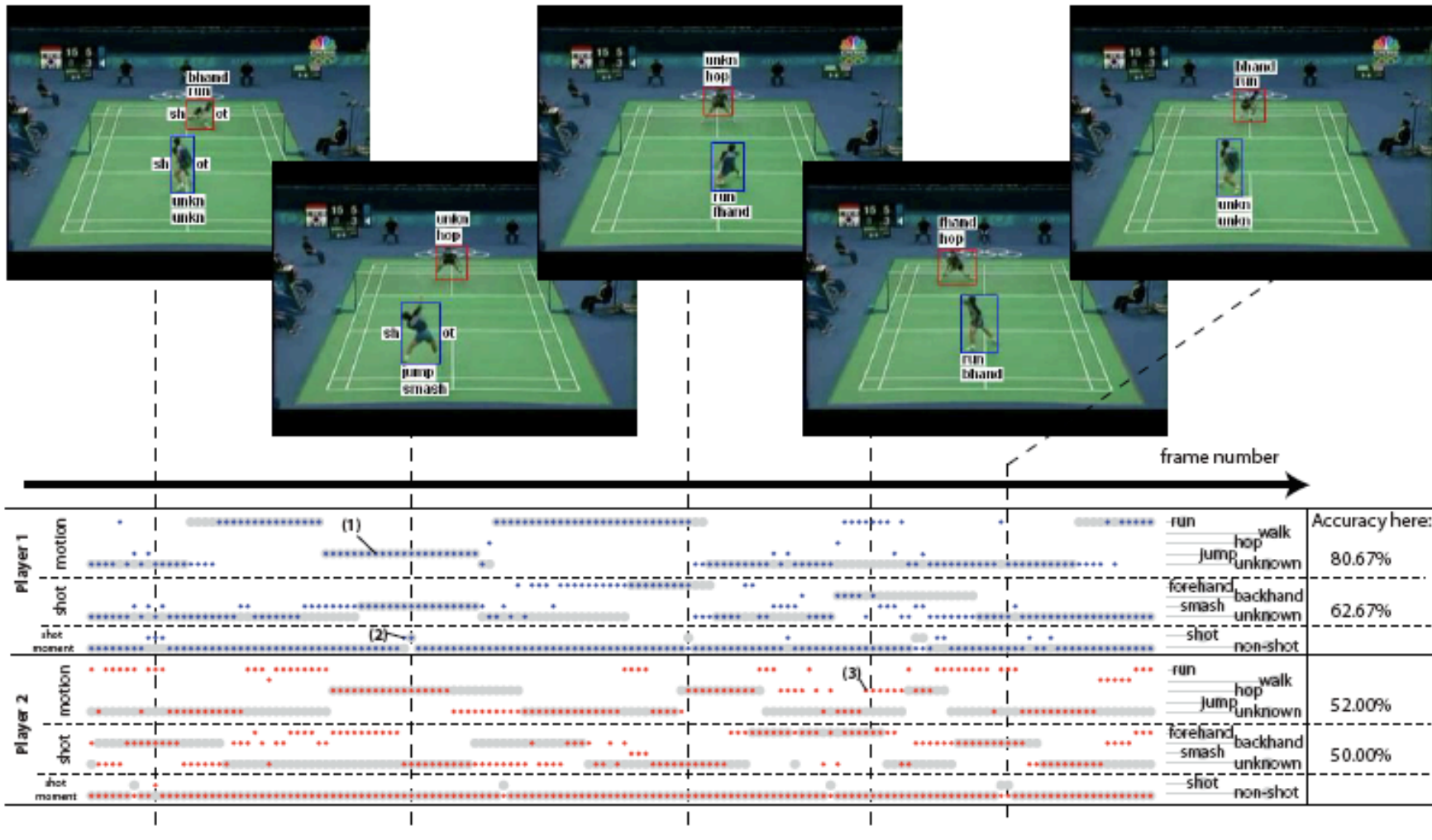
# Discriminative results

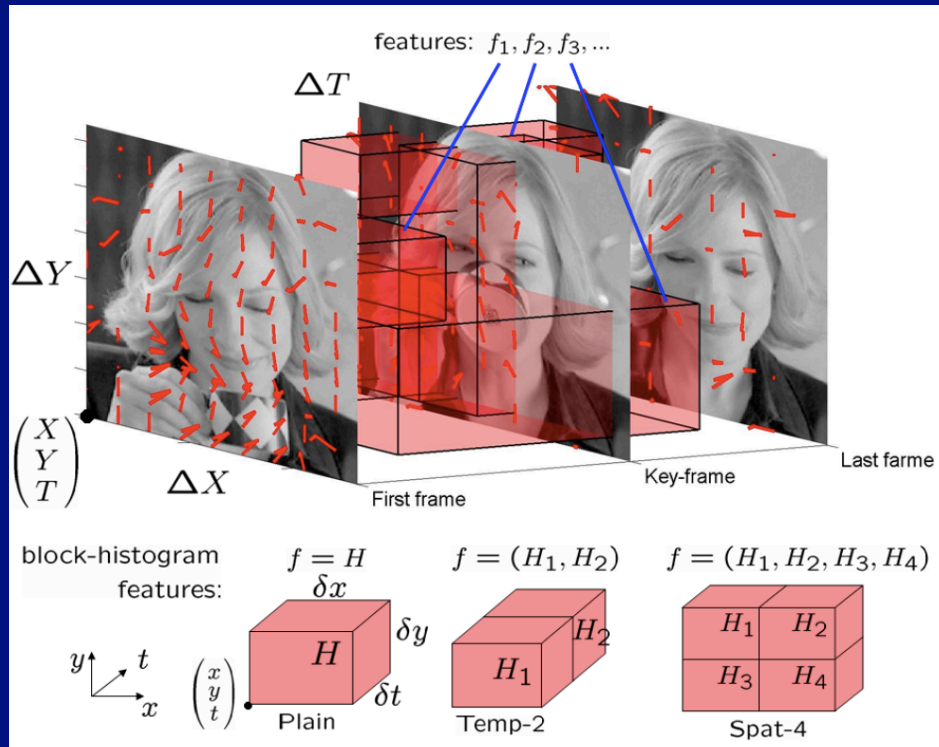
Dataset	Algorithm	Chance	Protocols								
			Discriminative task				Reject	Few examples			
			L1SO	L1AAO	L1AO	L1VO	UNa	FE-1	FE-2	FE-4	FE-8
Weizman	NB(k=300)	10.00	91.40	93.50	95.70	N/A	0.00	N/A	N/A	N/A	N/A
	1NN	10.00	95.70	95.70	96.77	N/A	0.00	53.00	73.00	89.00	96.00
	1NN-M	10.00	100.00	100.00	100.00	N/A	0.00	72.31	81.77	92.97	100.00
	1NN-R	9.09	83.87	84.95	84.95	N/A	84.95	17.96	42.04	68.92	84.95
	1NN-MR	9.09	<b>89.66</b>	<b>89.66</b>	<b>89.66</b>	N/A	<b>90.78</b>	N/A	N/A	N/A	N/A
Our	NB(k=600)	7.14	98.70	98.70	98.70	N/A	0.00	N/A	N/A	N/A	N/A
	1NN	7.14	98.87	97.74	98.12	N/A	0.00	58.70	76.20	90.10	95.00
	1NN-M	7.14	<b>99.06</b>	<b>97.74</b>	<b>98.31</b>	N/A	0.00	88.80	94.84	95.63	98.86
	1NN-R	6.67	95.86	81.40	82.10	N/A	81.20	27.40	37.90	51.00	65.00
	1NN-MR	6.67	<b>98.68</b>	<b>91.73</b>	<b>91.92</b>	N/A	<b>91.11</b>	N/A	N/A	N/A	N/A
IXMAS	NB(k=600)	7.69	80.00	78.00	79.90	N/A	0.00	N/A			
	1NN	7.69	81.00	75.80	80.22	N/A	0.00				
	1NN-R	7.14	<b>65.41</b>	<b>57.44</b>	<b>57.82</b>	N/A	<b>57.48</b>				
UMD	NB(k=300)	10.00	100.00	N/A	N/A	97.50	0.00	N/A			
	1NN	10.00	100.00	N/A	N/A	97.00	0.00				
	1NN-R	9.09	<b>100.00</b>	N/A	N/A	<b>88.00</b>	<b>88.00</b>				

Works well, depending on task; not rejecting improves things  
metric learning improves things



# Youtube video





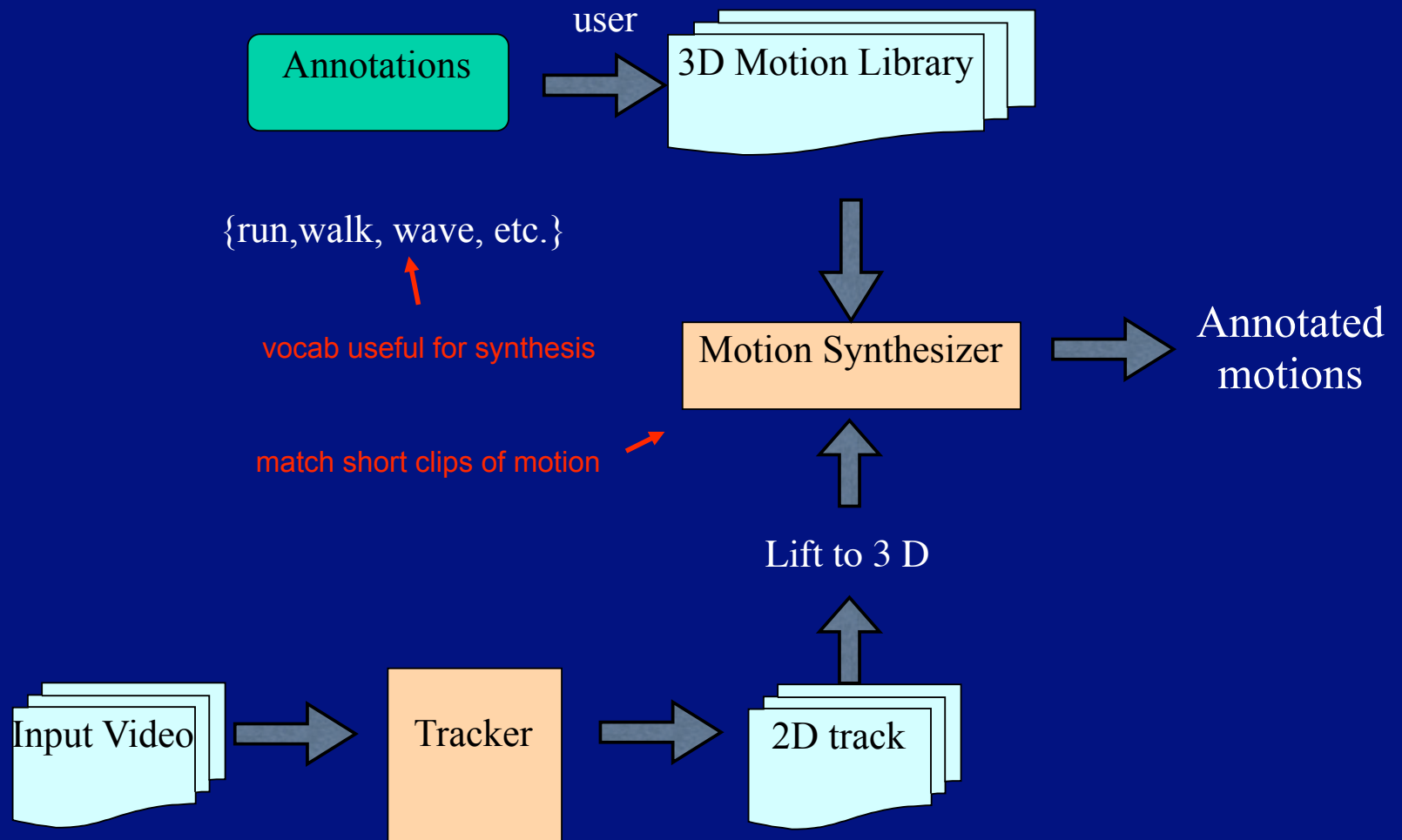
Laptev Perez 2007  
see also Laptev et al 08

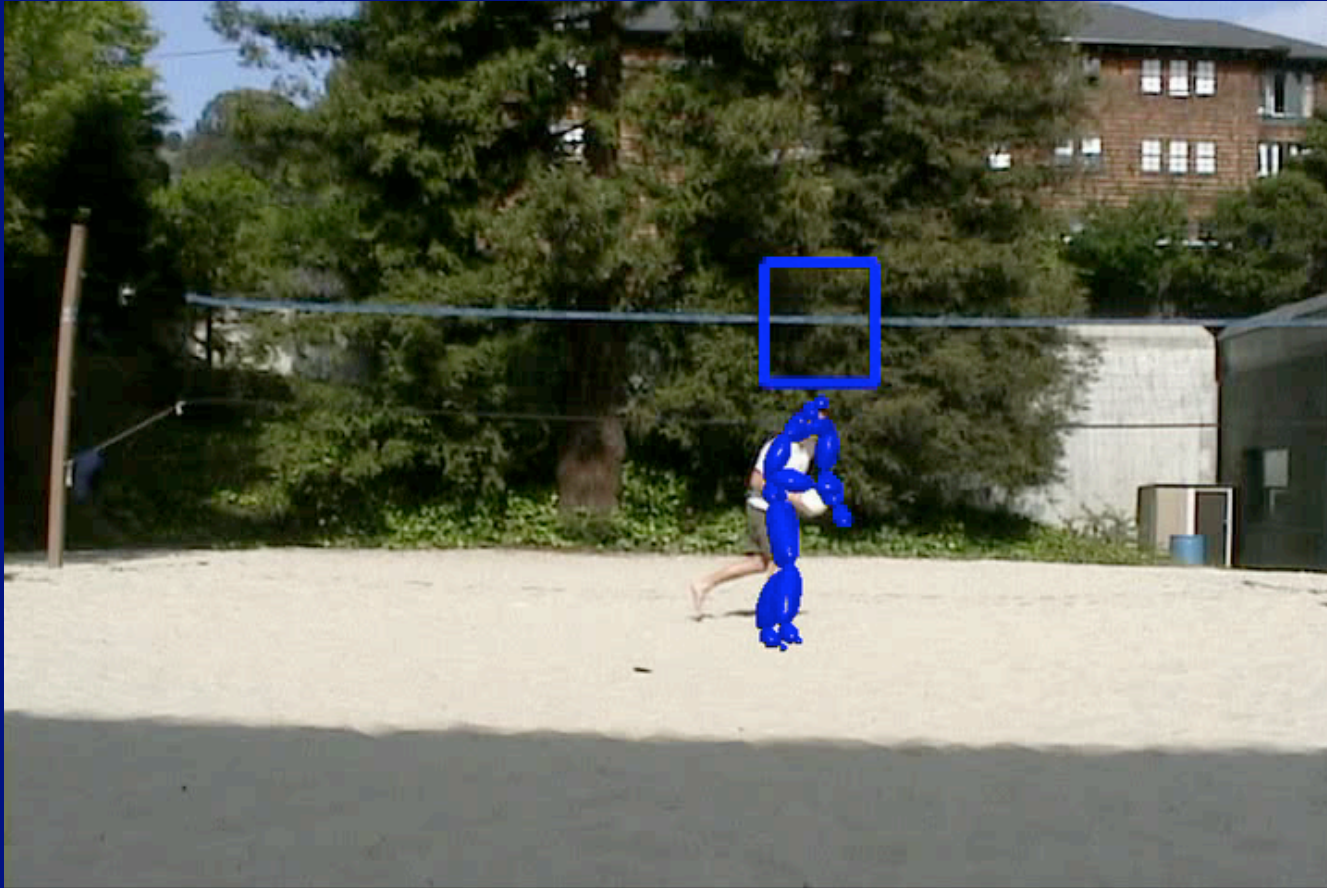


# Kinematic features

- Find body parts
  - with geometric/appearance model (deformable template)
  - cardboard people
    - (eg Ju et al 96; Sidenbladh et al 2000)
  - pictorial structures
    - (eg Felzenszwalb Huttenlocher 05)
  - kinematic tracks
    - (eg Ramanan et al 05)
  - repeated model-based segmentation
    - (eg Ferrari et al 08)

# Annotating observations by synthesis





Ramanan+Forsyth 03

# Criteria

- Base accuracy?
  - appearance wins hands down on current datasets
- Aspect
  - appearance can be fixed
- Do they solve the right problem?
  - advantage: kinematic



# IXMAS and Aspect

Camera 0



Camera 4



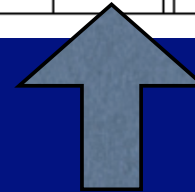
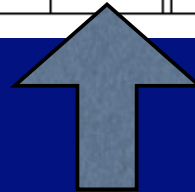
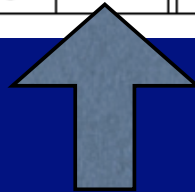
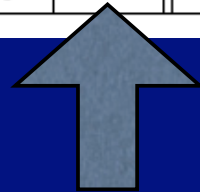
# The Effects of Aspect

	Camera 0		Camera 1		Camera 2		Camera 3		Camera 4	
FO	76		76		68		73		51	
	WT		WT		WT		WT		WT	
Camera 0	NA		35		16		8		10	
Camera 1	38		NA		15		8		11	
Camera 2	16		16		NA		6		11	
Camera 3	8		8		8		NA		8	
Camera 4	12		11		15		9		NA	



# Results

	Camera 0			Camera 1			Camera 2			Camera 3			Camera 4		
	QV	SS	CV	QV	SS	CV	QV	SS	CV	QV	SS	CV	QV	SS	CV
Camera 0	76	76	<b>84</b>	72	78	<b>79</b>	61	69	<b>79</b>	62	<b>70</b>	68	30	45	<b>76</b>
Camera 1	69	<b>77</b>	72	76	78	<b>85</b>	64	<b>74</b>	<b>74</b>	68	67	<b>70</b>	41	44	<b>66</b>
Camera 2	62	66	<b>71</b>	67	71	<b>82</b>	68	74	<b>87</b>	67	64	<b>76</b>	43	54	<b>72</b>
Camera 3	63	69	<b>75</b>	72	70	<b>75</b>	68	63	<b>79</b>	73	68	<b>87</b>	44	44	<b>76</b>
Camera 4	51	39	<b>80</b>	55	39	<b>73</b>	51	52	<b>73</b>	53	34	<b>79</b>	51	66	<b>80</b>



Farhadi Kamali 08

		test views					
		cam1	cam2	cam3	cam4	cam5	All
training views	cam1	76.4	77.6	69.4	70.3	44.8	67.2
	cam2	77.3	77.6	73.9	67.3	43.9	67.4
	cam3	66.1	70.6	73.6	63.6	53.6	65.0
	cam4	69.4	70.0	63.0	68.8	44.2	63.9
	cam5	39.1	38.8	51.8	34.2	66.1	45.2
	All	74.8	74.5	74.8	70.6	61.2	72.7

Junejo et al 08, different feature construction, same dataset

# The problem we have been solving

- Rack up a bunch of activity categories, and discriminate
  - how many categories are enough?
  - can one movement have two categories?
  - what are the categories?
    - the verb argument (probably) fails
      - if there are few movement, many goal verbs
    - introspection suggests too few words

AnswerPhone, GetOutCar, Handshake, Kiss, Hugperson, SitDown, SitUp, StandUp

Goal achieved by body movement

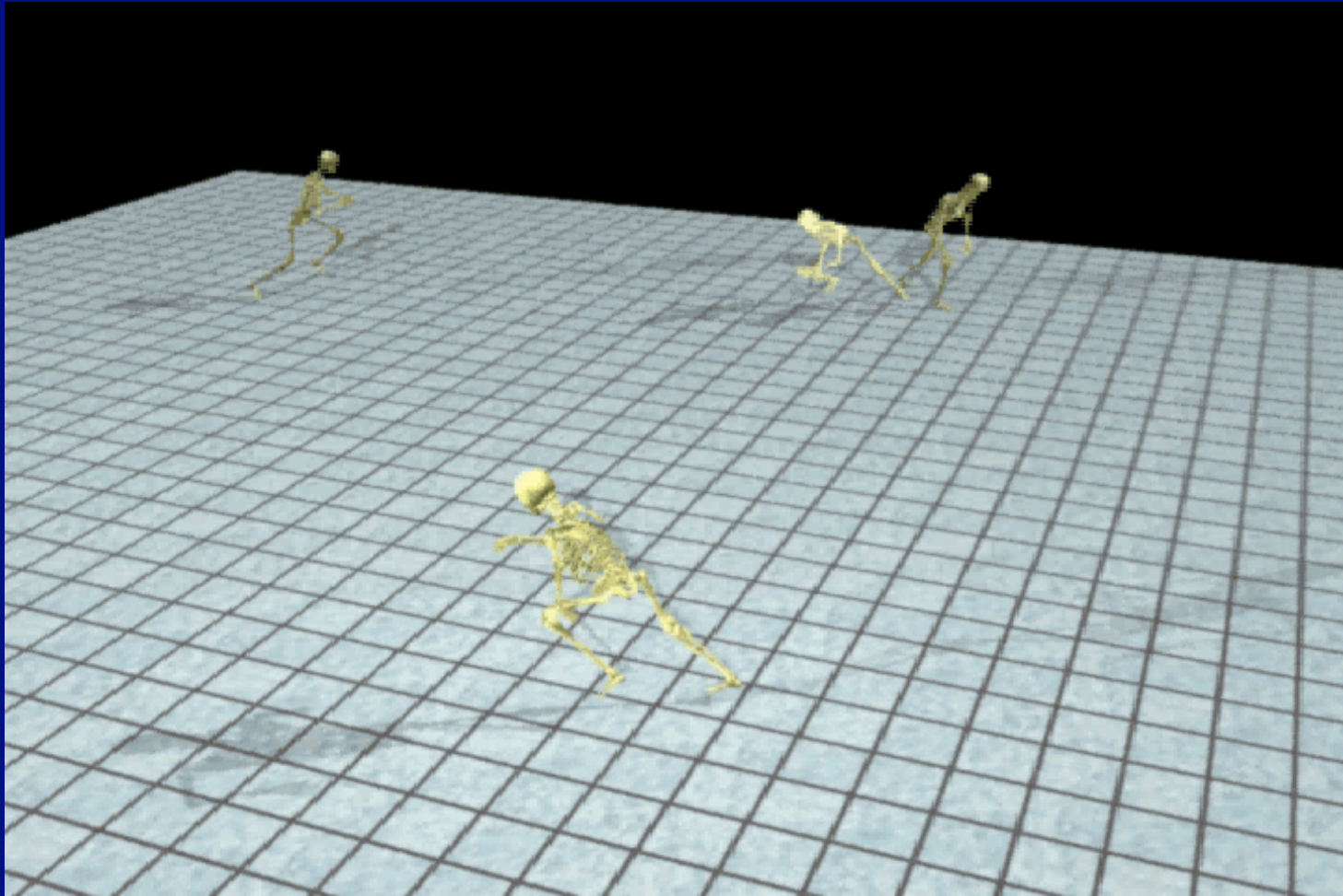
Body movement

# Components of the problem we should be trying to map

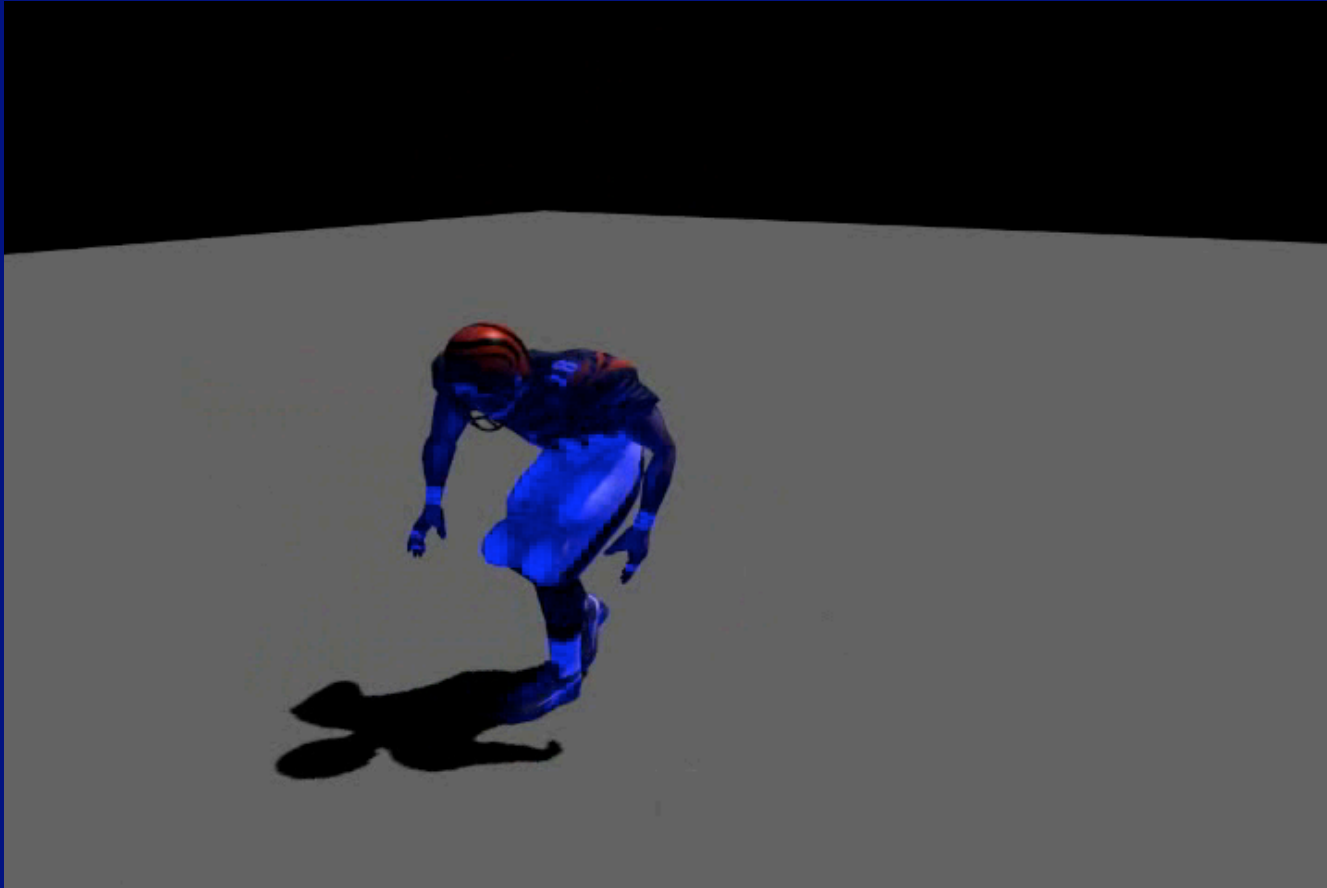
- Activity composes freely into complex structures
- Most human activities cause changes of state, meet goals
  - similar movements will meet different goals
  - different movements can meet the same goal
- We should probably be trying to “recognize” things
  - whose names we do not know
    - fluid, changing categories, affected by
      - nearby objects
      - observer, observation context
  - for which we have seen no examples

# 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)

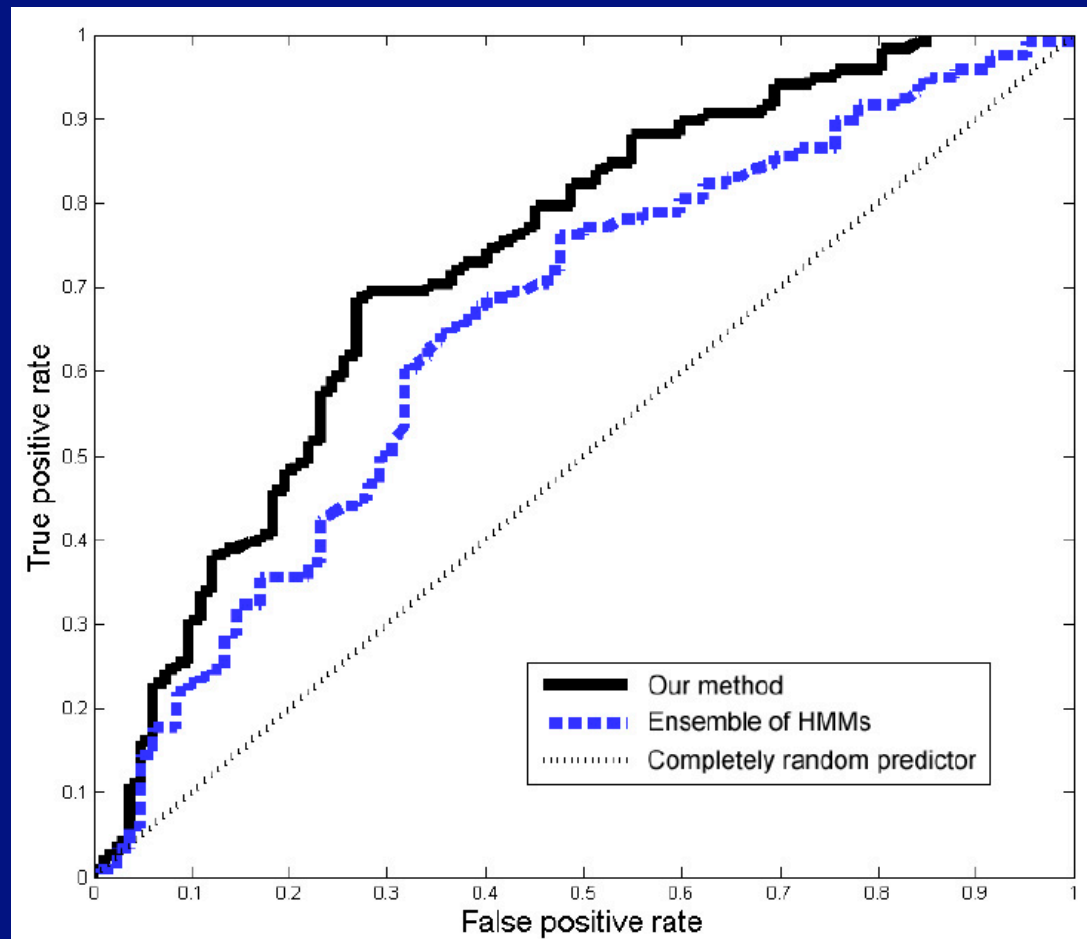






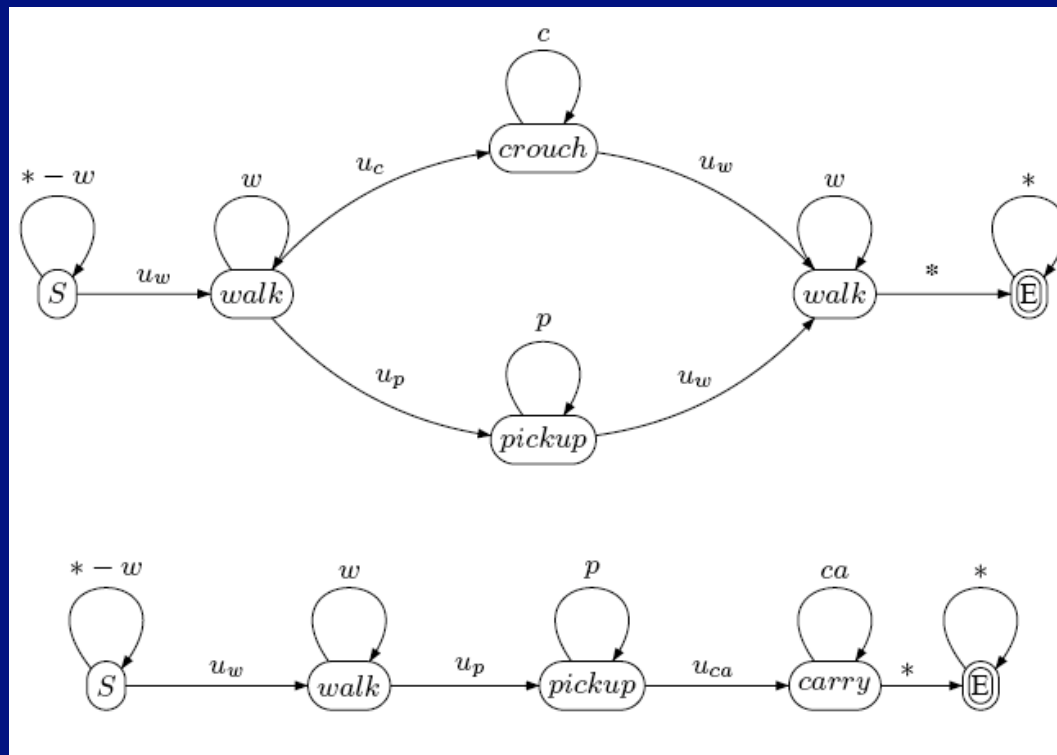


# Hard to tell good from bad



# “Recognizing” composites

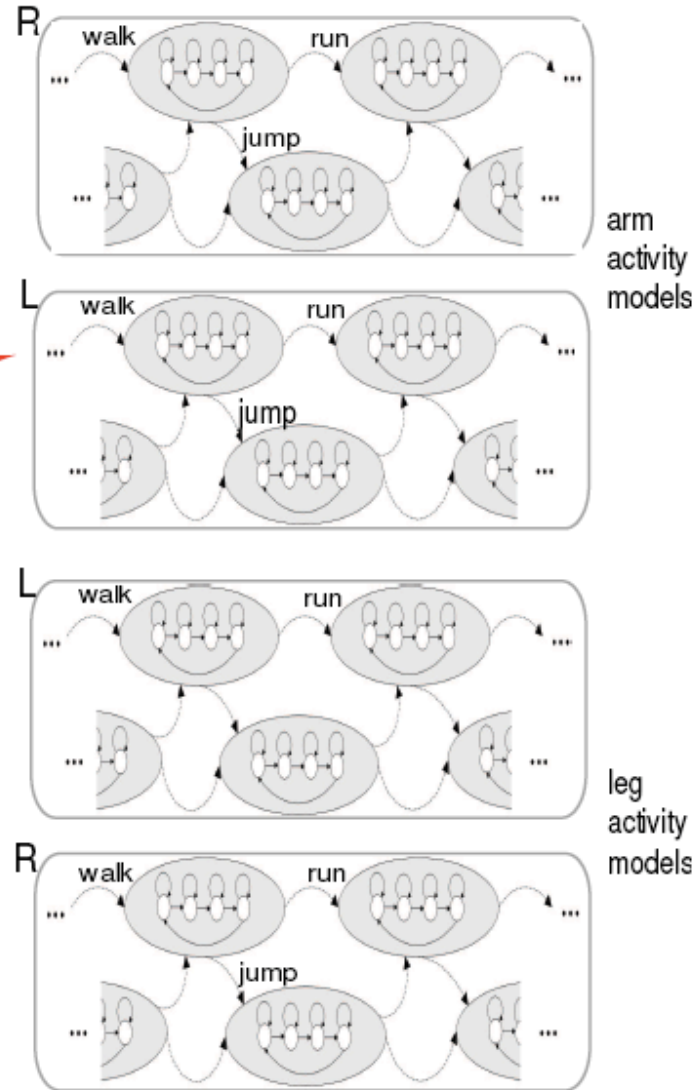
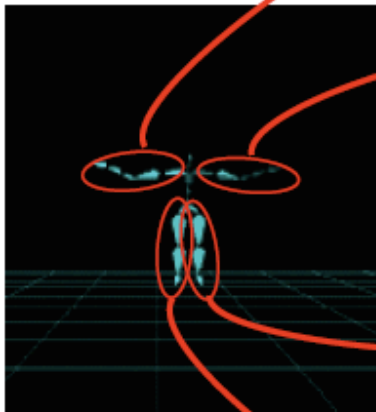
- Rank sequences by  $P(\text{FSAl data, model})$ 
  - e.g.  $P(\text{leg-walk-arm-walk-then-leg-walk-arm-reach | data, model})$
  - DP variant will do this easily



# Building a composite model

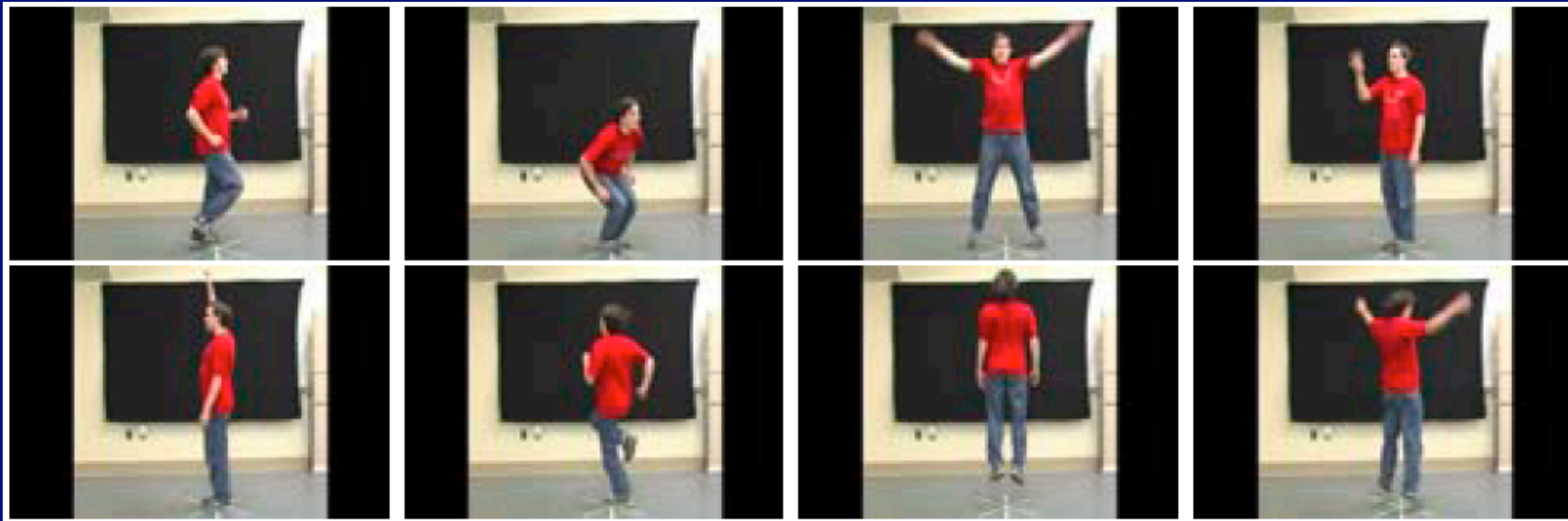
- Build a set of basic labels
  - guess them: walk, run, stand, reach, crouch, etc.
- Activity model:
  - Product of finite state automata for arms, legs built from MoCap
  - Arms, legs each have local short timescale activity models for labels
  - Link these models into a large model, using animation-legal transitions

# Composition



# Emission

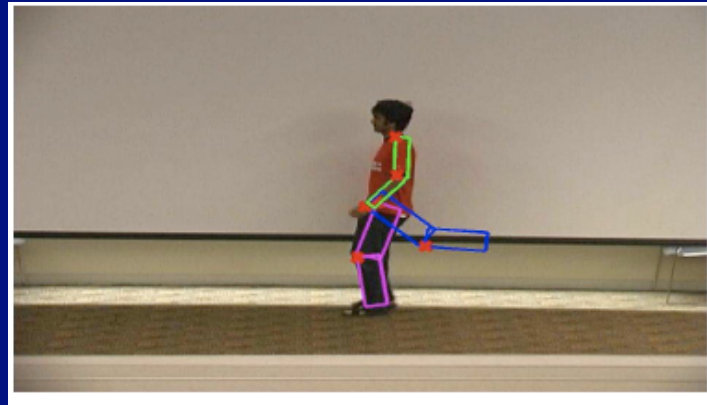
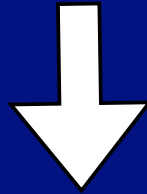
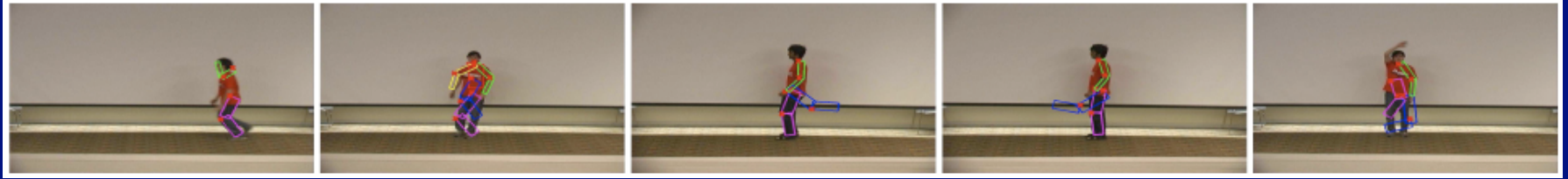
- Transduction
  - Track the body, as above
  - Lift “snippets” of each quarter
    - vector quantized
  - impose root consistency
- Emission
  - emit cluster center from state according to table
  - table learned by EM, known dynamical model



Ikizler Forsyth 07,08

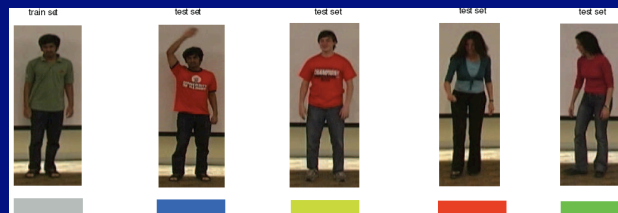
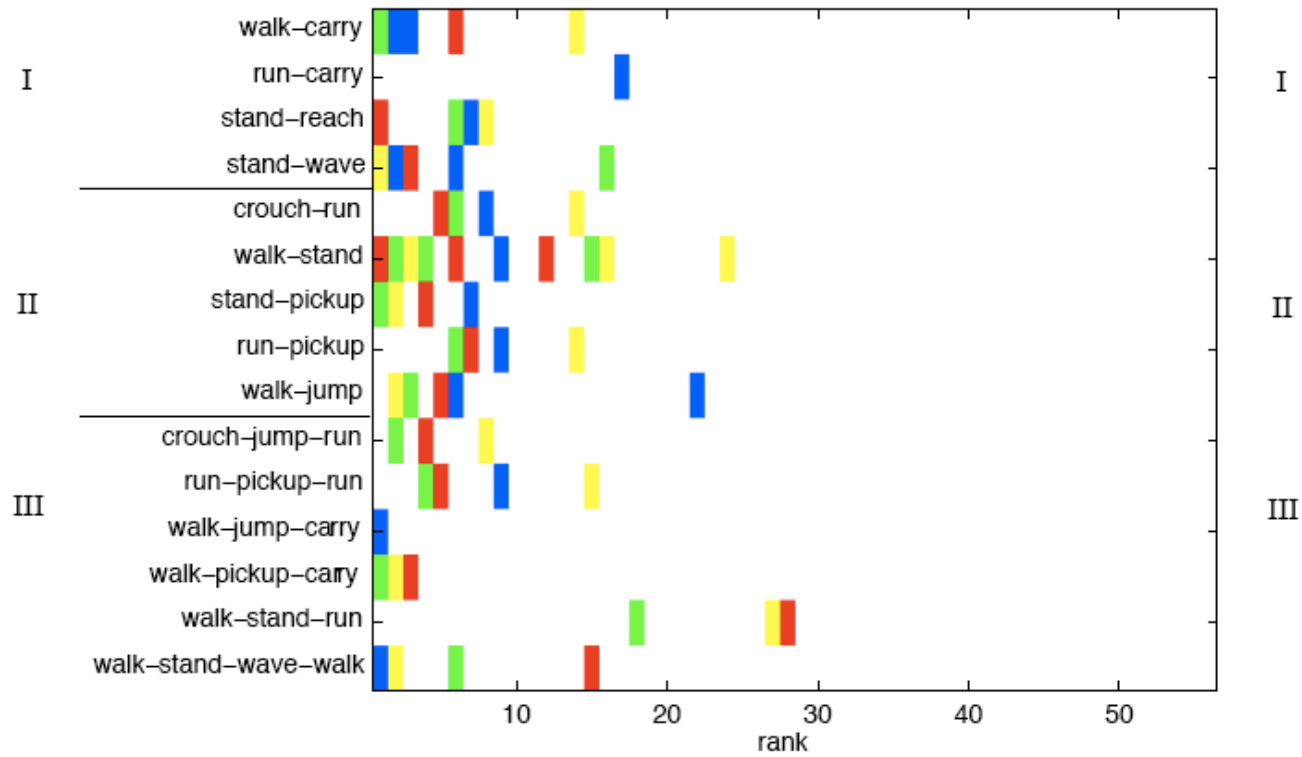
Context	# videos	Context	# videos
crouch-run	2	run-backwards-wave	2
jump-jack	2	run-jump-reach	5
run-carry	2	run-pickup-run	5
run-jump	2	walk-jump-carry	2
run-wave	2	walk-jump-walk	2
stand-pickup	5	walk-pickup-walk	2
stand-reach	5	walk-stand-wave-walk	5
stand-wave	2	crouch-jump-run	3
walk-carry	2	walk-crouch-walk	3
walk-run	3	walk-pickup-carry	3
run-stand-run	3	walk-jump-reach-walk	3
run-backwards	2	walk-stand-run	3
walk-stand-walk	3		

Figure 1. Contexts, # videos, # actions, # transitions, #1, #2, #3, #4, #5, #6, #7, #8, #9, #10, #11, #12, #13, #14, #15, #16, #17, #18, #19, #20, #21, #22, #23, #24, #25, #26, #27, #28, #29, #30, #31, #32, #33, #34, #35, #36, #37, #38, #39, #40, #41, #42, #43, #44, #45, #46, #47, #48, #49, #50, #51, #52, #53, #54, #55, #56, #57, #58, #59, #60, #61, #62, #63, #64, #65, #66, #67, #68, #69, #70, #71, #72, #73, #74, #75, #76, #77, #78, #79, #80, #81, #82, #83, #84, #85, #86, #87, #88, #89, #90, #91, #92, #93, #94, #95, #96, #97, #98, #99, #100



Ikizler Forsyth 07,08

### Our Method



Ikizler Forsyth 07,08

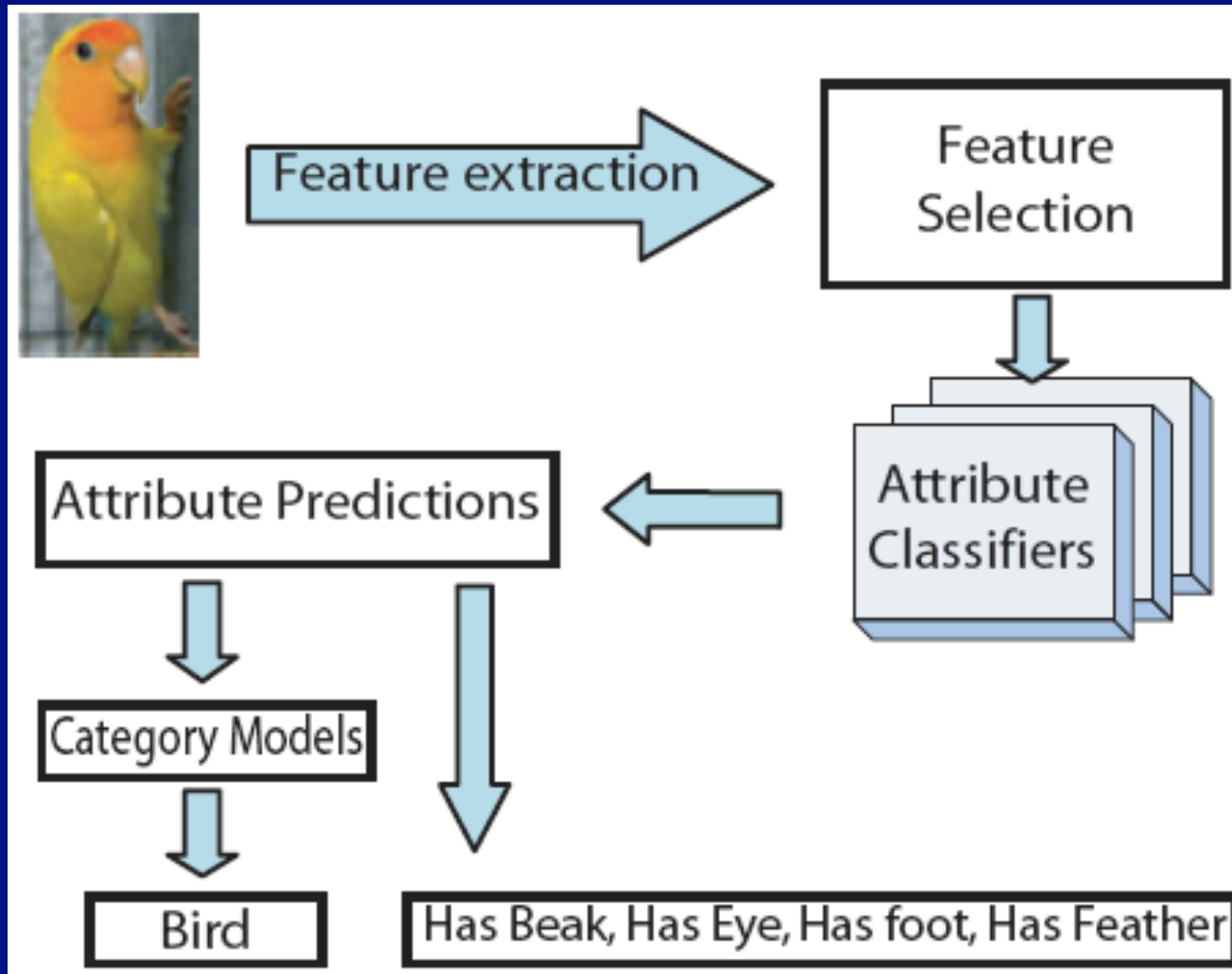


# How do you describe something whose name is unknown?

- Attributes
  - Properties shared by many object categories
  - Material (like)
    - glass, wood, furry, red, etc.
  - Part (like)
    - has wheel, has head, has tail, etc.
  - Shape (like)
    - is 2D Boxy, is cylindrical, etc
- What do we need to say about activity?
  - should we name activity, or reason about goals, intentions?
  - what about the objects nearby?

Farhadi et al 09;  
cf Blaschke 09;  
Ferrari Zisserman 07;

# General architecture





'is 3D Boxy'

'is Vert Cylinder'

'has Window' ~~'has Screen'~~

'has Row Wind' ~~'has Headlight'~~



'has Hand'

'has Arm'

~~'has Screen'~~

'has Plastic' ~~'has Saddle'~~

'is Shiny'



'has Head'

'has Hair'

'has Face'

~~'has Saddle'~~

'has Skin' ~~'has Wood'~~



'has Head'

'has Torso'

'has Arm'

'has Leg'

~~'has Wood'~~



'has Head'

'has Ear'

'has Snout'

'has Nose'

'has Mouth'



'has Head'

'has Ear'

'has Snout'

'has Mouth'

'has Leg'



~~'has Furniture Back'~~

~~'has Horn'~~

~~'s Screen'~~

'has Plastic'

'is Shiny'



'is 3D Boxy'

'has Wheel'

'has Window'

'is Round'

'has Torso'



'has Tail'

'has Snout'

'has Leg'

~~'has Text'~~

~~'has Plastic'~~



'has Head'

'has Ear'

'has Snout'

'has Leg'

'has Cloth'



'is Horizontal Cylinder'

~~'has Beak'~~

~~'has Wing'~~

~~'has Side mirror'~~

'has Metal'



'has Head'

'has Snout'

'has Horn'

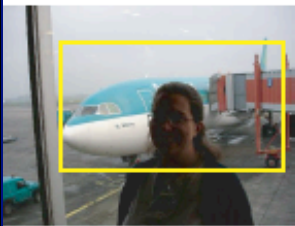
'has Torso'

~~'has Arm'~~

# How is an object different from typical?

- Pragmatics suggests this is how adjectives are chosen
  - If we are sure it's a cat, and we know that
    - an attribute is different from normal
    - the detector is usually reliable
  - we should report the missing/extra attribute

# Missing attributes



Aeroplane  
No "wing"



Car  
No "window"



Boat  
No "sail"



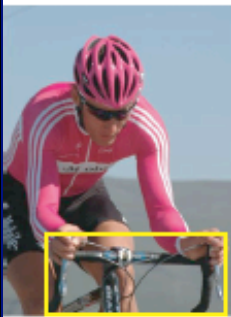
Aeroplane  
No "jet engine"



Motorbike  
No "side mirror"



Car  
No "door"



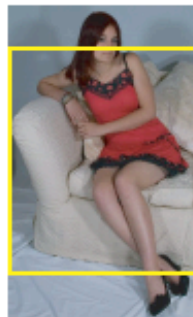
Bicycle  
No "wheel"



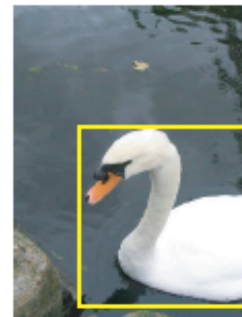
Sheep  
No "wool"



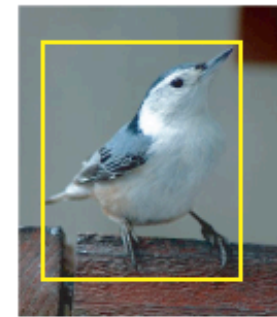
Train  
No "window"



Sofa  
No "wood"



Bird  
No "tail"



Bird  
No "leg"



Bus  
No "door"



# Extra attributes



Bird  
"Leaf"



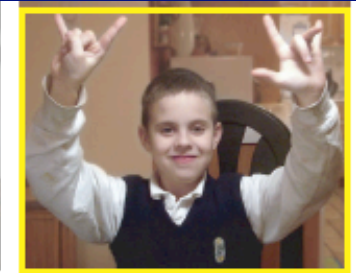
Bus  
"face"



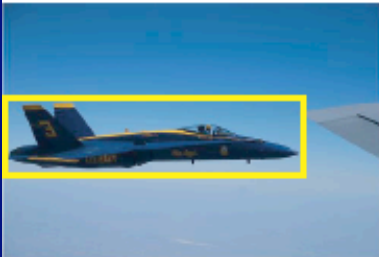
Motorbike  
"cloth"



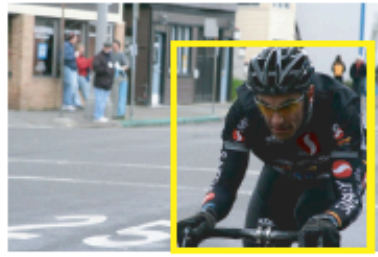
DiningTable  
"skin"



People  
"Furn.back"



Aeroplane  
"beak"



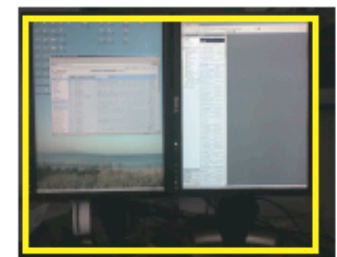
People  
"label"



Sofa  
"wheel"



Bike  
"Horn"



Monitor  
"window"

# Conclusions

- Absent taxonomy/composition is a major nuisance
  - if it were not for this question, appearance methods would win hands down
- What do we need to say about activity?
  - should we name activity, or reason about goals, intentions?
  - what about the objects nearby?
- Object recognition ~~is in a fool's paradise~~ has to deal with similar issues
  - unknown names, etc.