

Logistics

Class webpage

Send me an e-mail with at least 5
papers by this Friday

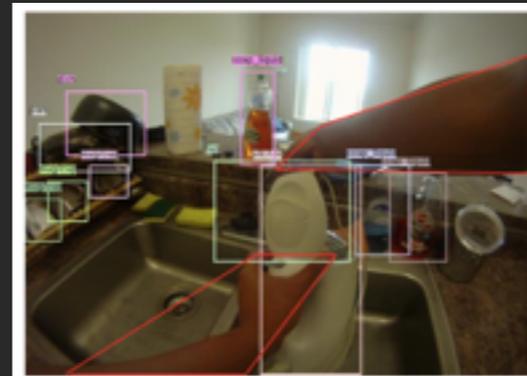
Next two weeks

This class: background on temporal analysis

Next class: background on deep models

Agenda

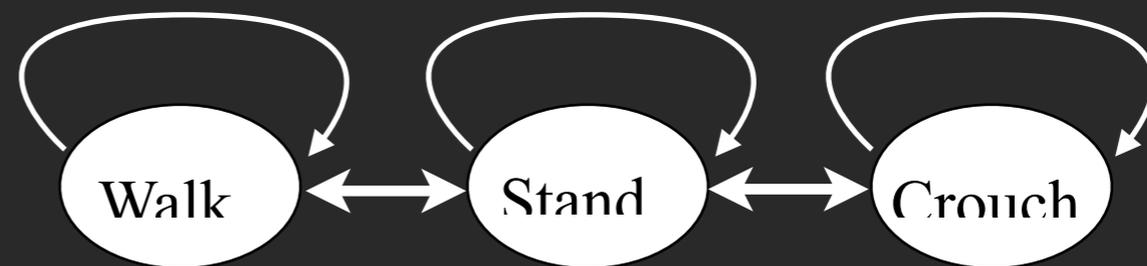
Data/benchmark analysis



Spatiotemporal features



Spatiotemporal models



[Caveat - I'll talk about many of this issues wrt my own work]

Image benchmarks



Torralba, et al. PAMI 2008.



Xiao, et al. CVPR 2010.

Like it or not, crucial for advances in the field

Large-scale annotated video datasets are more rare - why?

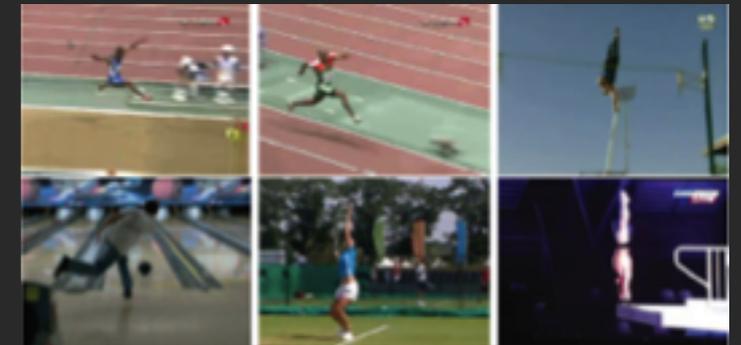
Action recognition benchmarks



UCF 101 Sports, 2013



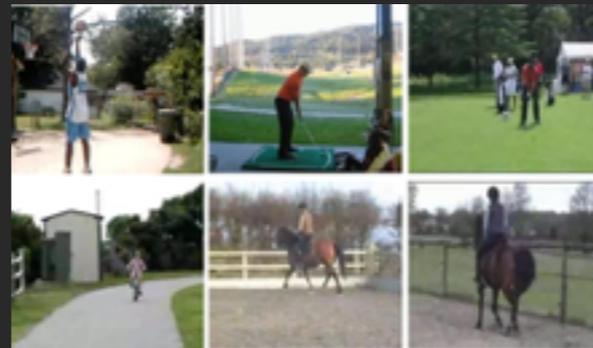
KTH, ICPR'04



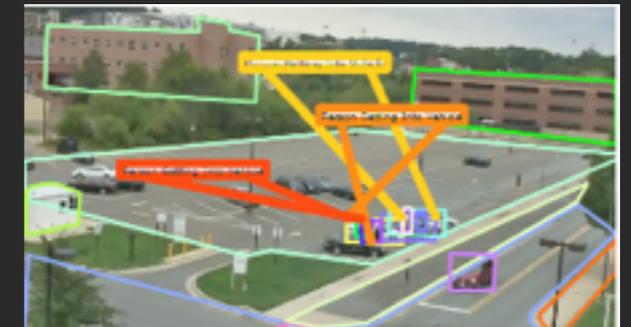
Olympics sport, BMVC'10



Hollywood, CVPR'09



UCF Youtube, CVPR'08



VIRAT, CVPR'11

- 1) Video is cumbersome to label (difficult to define natural categories outside sports)
- 2) Collecting interesting but natural video is surprisingly hard
- 3) Most current work focuses on K-way classification (similar to image recognition 10 years ago)

KTH



Classification performance around 100%
“Outdated”

TRECVID



“Woodworking” action

board-trick, feeding animal, fishing, wedding, woodworking, birthday, changing vehicle tire, flash mob, vehicle unstuck, grooming an animal, sandwich making, parade, parkour, repairing appliance and sewing,...

State-of-the-art is around 5-10% accuracy

Challenge 1: how we do know what to label?

Look for cues in language (how do people describe images/videos?)

RICK

Why weren't you honest with me? **Why did** you keep your marriage a secret?

Rick sits down with Lisa.

LISA

Oh, it wasn't my secret, Richard. Victor wanted it that way. Not even our closest friends knew about our marriage.



This is a lot of technology.
Somebody's screensaver of a pumpkin.
Black laptop is connected to black Dell monitor.
Old school Computer monitor with many stickers on it.
A refrigerator full of food.

Mining movie scripts
Everingham et al. BMVC
Laptev et al 08.

Ask people on turk for descriptions
Farhadi et al ECCV10

Challenge 1: how we do know what to label?

Look for cues in medical/labor literature on “activities of daily living” (ADLs)

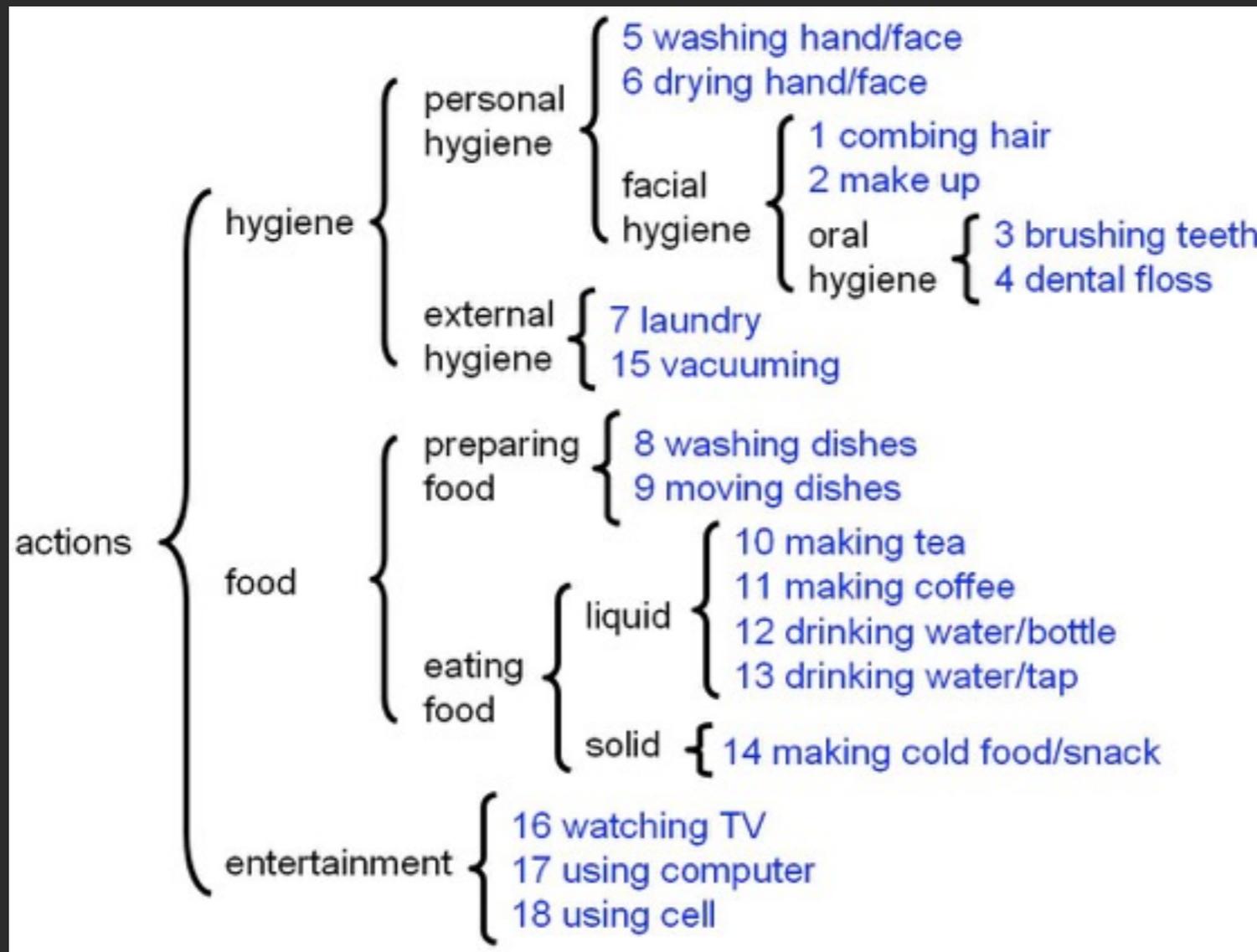
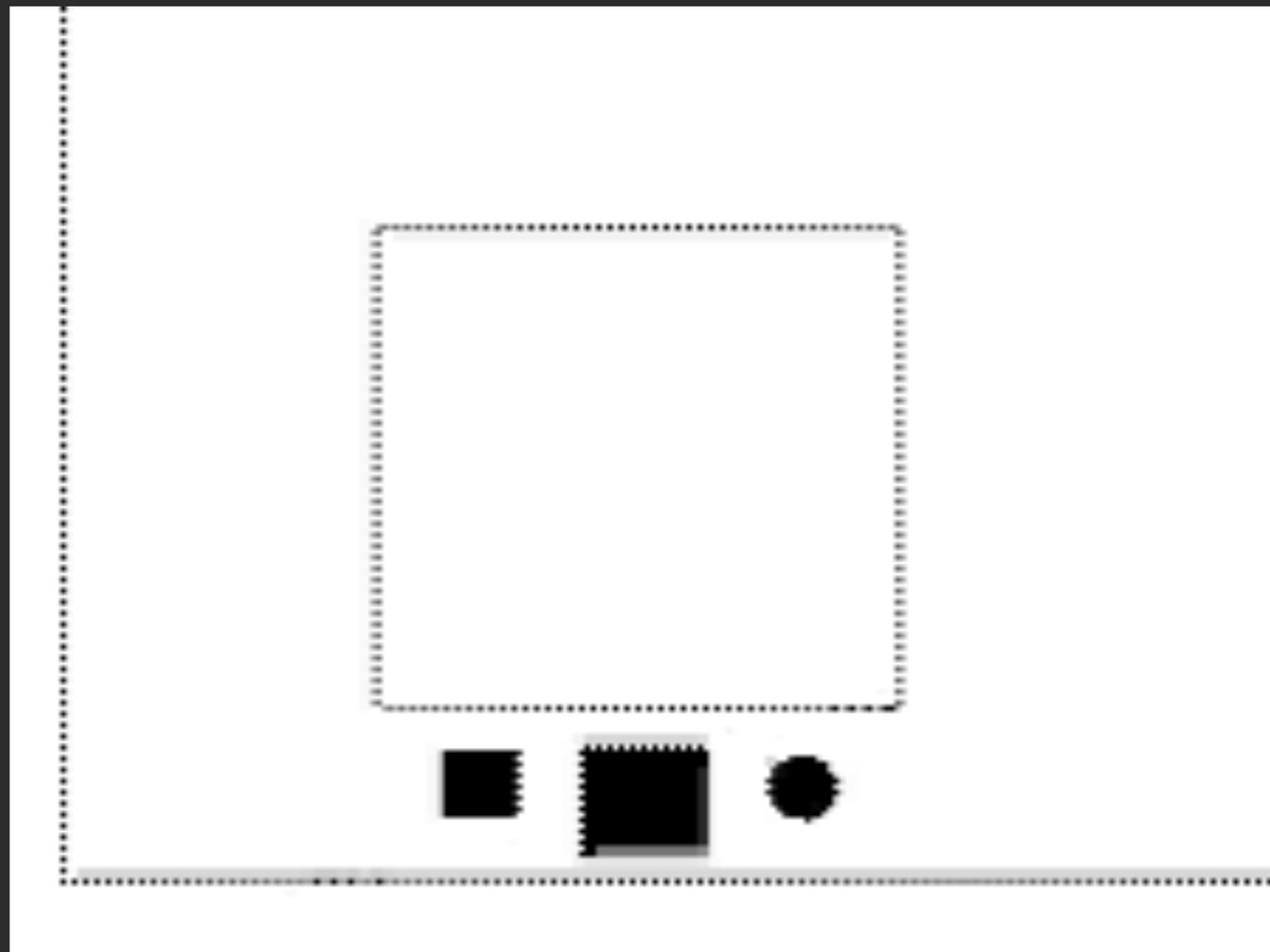


Table 1: AMAT Compound Tasks and Task Components

- I. Cut "Meat"
 1. Pick up knife and fork*
 2. Cut "meat" (Play-Doh)*†
 3. Fork to mouth
- II. Foam "Sandwich"
 4. Pick up foam "sandwich"
 5. "Sandwich" to mouth
- III. Eat With Spoon
 6. Pick up spoon
 7. Pick up dried kidney bean with spoon
 8. Spoon to mouth
- IV. Drink From Mug
 9. Grasp mug handle
 10. Mug to mouth
- V. Comb Hair
 11. Pick up comb
 12. Comb hair†
- VI. Open Jar
 13. Grasp jar top*
 14. Screw jar top open*
- VII. Tie Shoelace
 15. Tie shoelace*†
- VIII. Use Telephone
 16. Phone received to ear
 17. Press phone number

Challenge 1: how we do know what to label?

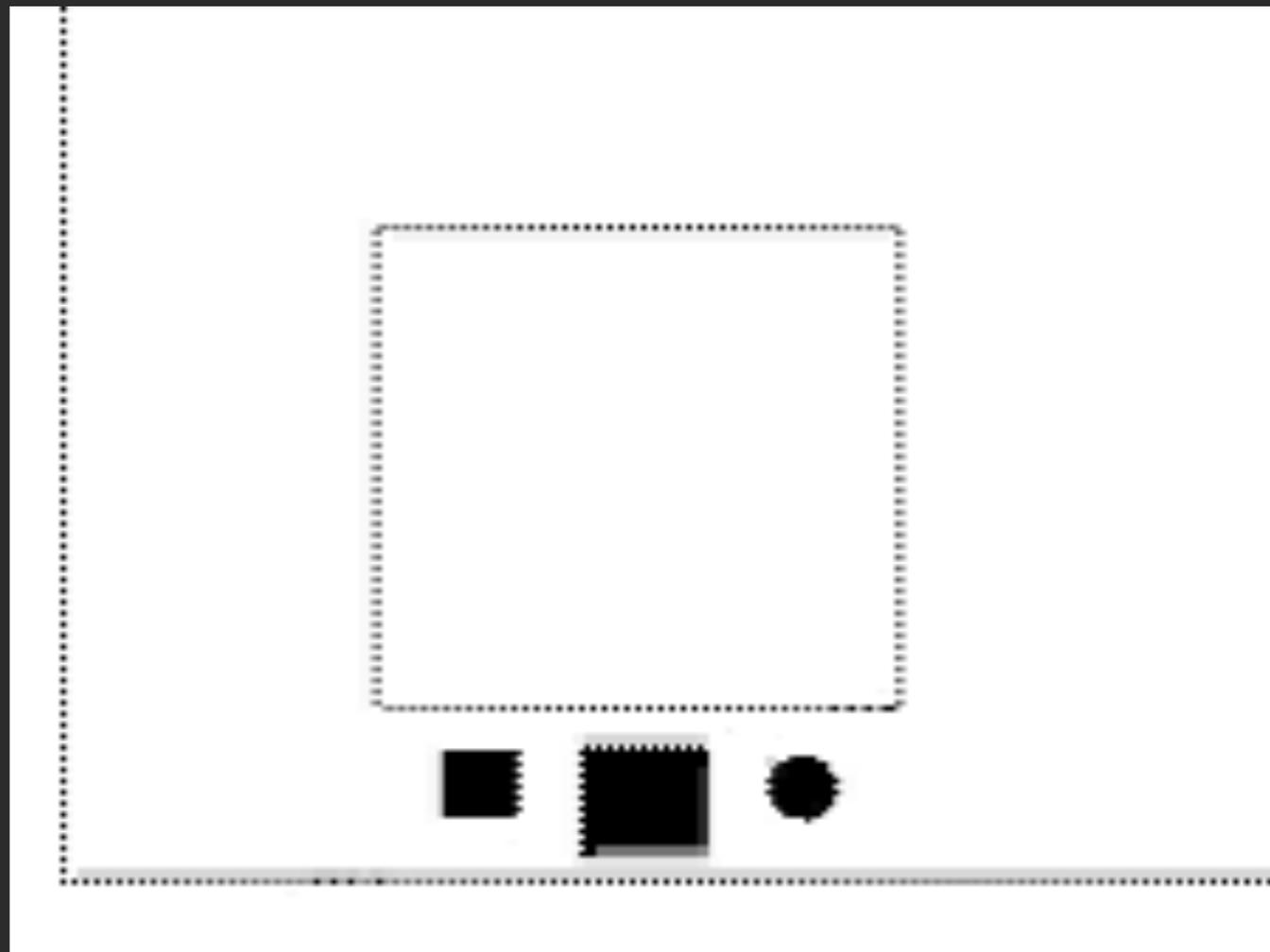
Actions vs goal-directed behaviors



Chase vs follow

Challenge 1: how we do know what to label?

Actions vs goal-directed behaviors



Chase vs follow

What is the relevant perceptual output here?

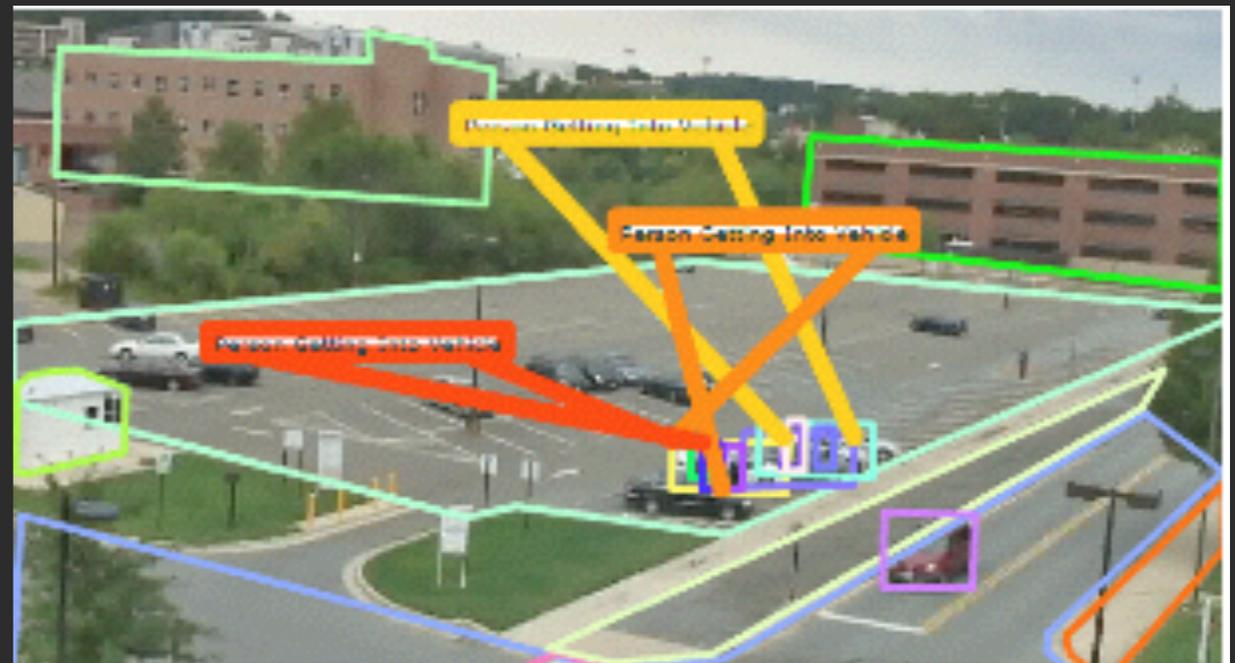


(Video thanks to David Forsyth)

Challenge 2: how we do obtain interesting data?



Script it, using actors



Use real but “boring” data

Challenge 2: how we do obtain interesting data?



Egocentric/wearable cameras

“Functional” ADLs
Easy to capture variety-rich data

Challenge 3: how we do produce detailed annotations?

Crowdsourcing labeling

Annotate every object, even stationary and obstructed objects, for the entire video. [Instructions](#) [+ New Object](#)

Car 12 Outside of view frame Occluded or obstructed Parked Driving Reversing

Person 11 Outside of view frame Occluded or obstructed Walking Running Standing

Car 10 Outside of view frame Occluded or obstructed Parked Driving Reversing

Car 9 Outside of view frame

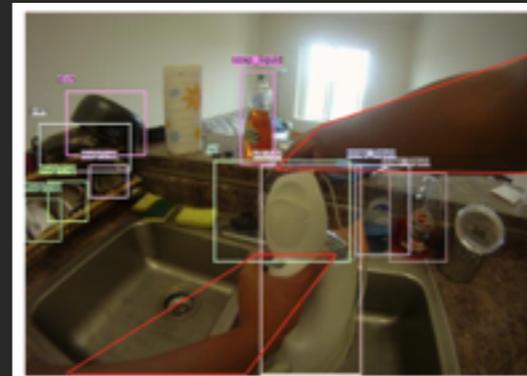
◀ Rewind ▶ **Play**

Disable Resize? Hide Boxes? Hide Labels? Slower Slow Normal Fast [Save Work](#)

Lessons: Interface design matters
Use experts, not the crowd
Active annotation helps

Roadmap

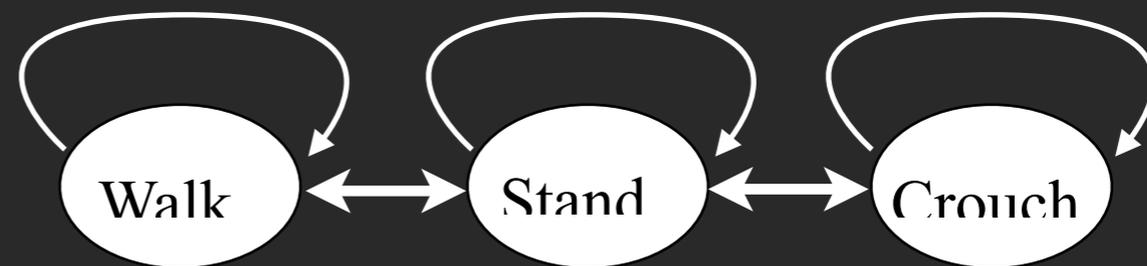
Data/benchmark analysis



Spatiotemporal features



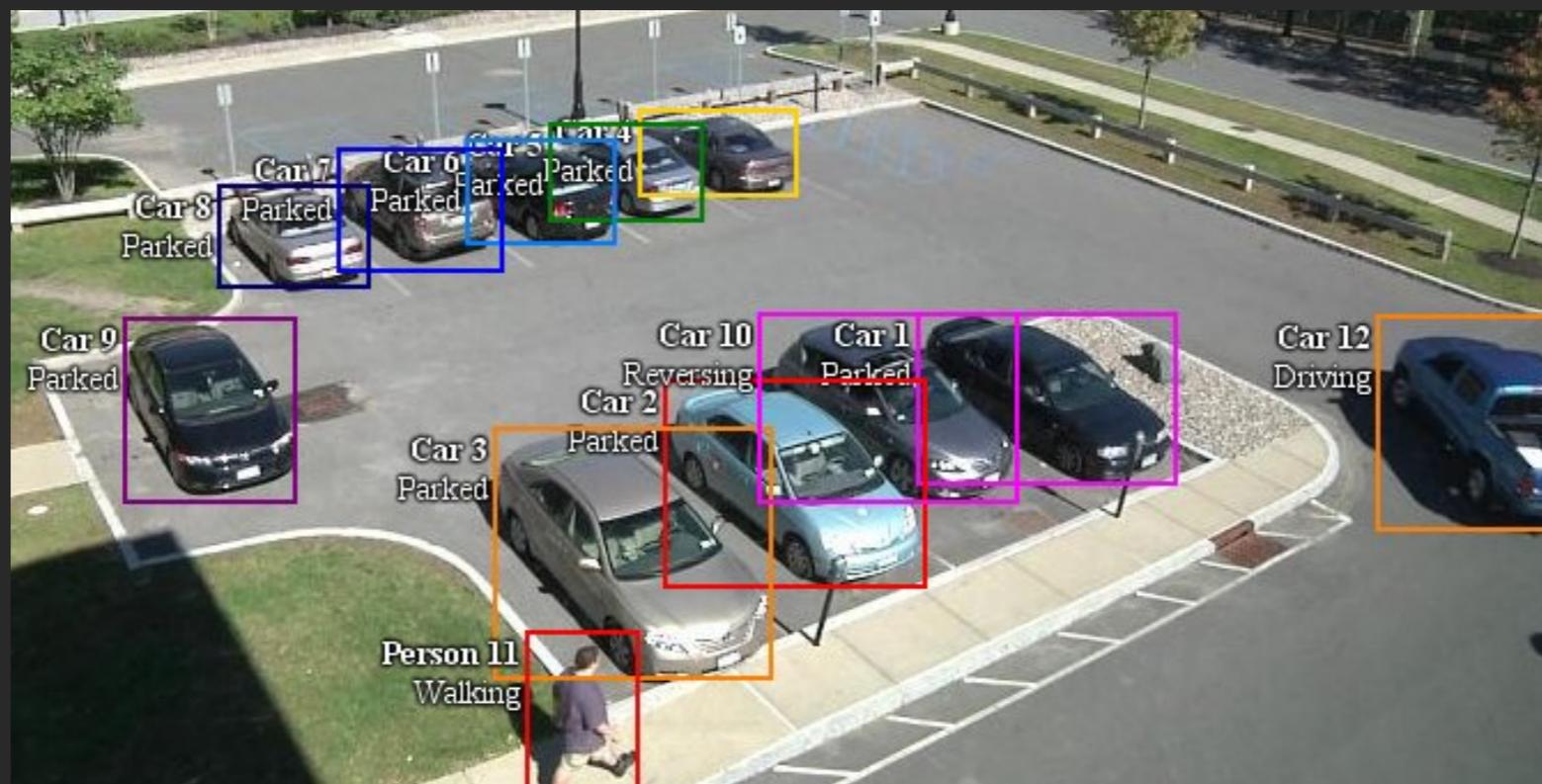
Spatiotemporal models



Spacetime features

Simple approach: just use spatial features

Surprisingly (and annoyingly) effective



Build a bank of static-image detectors (of poses, objects,)

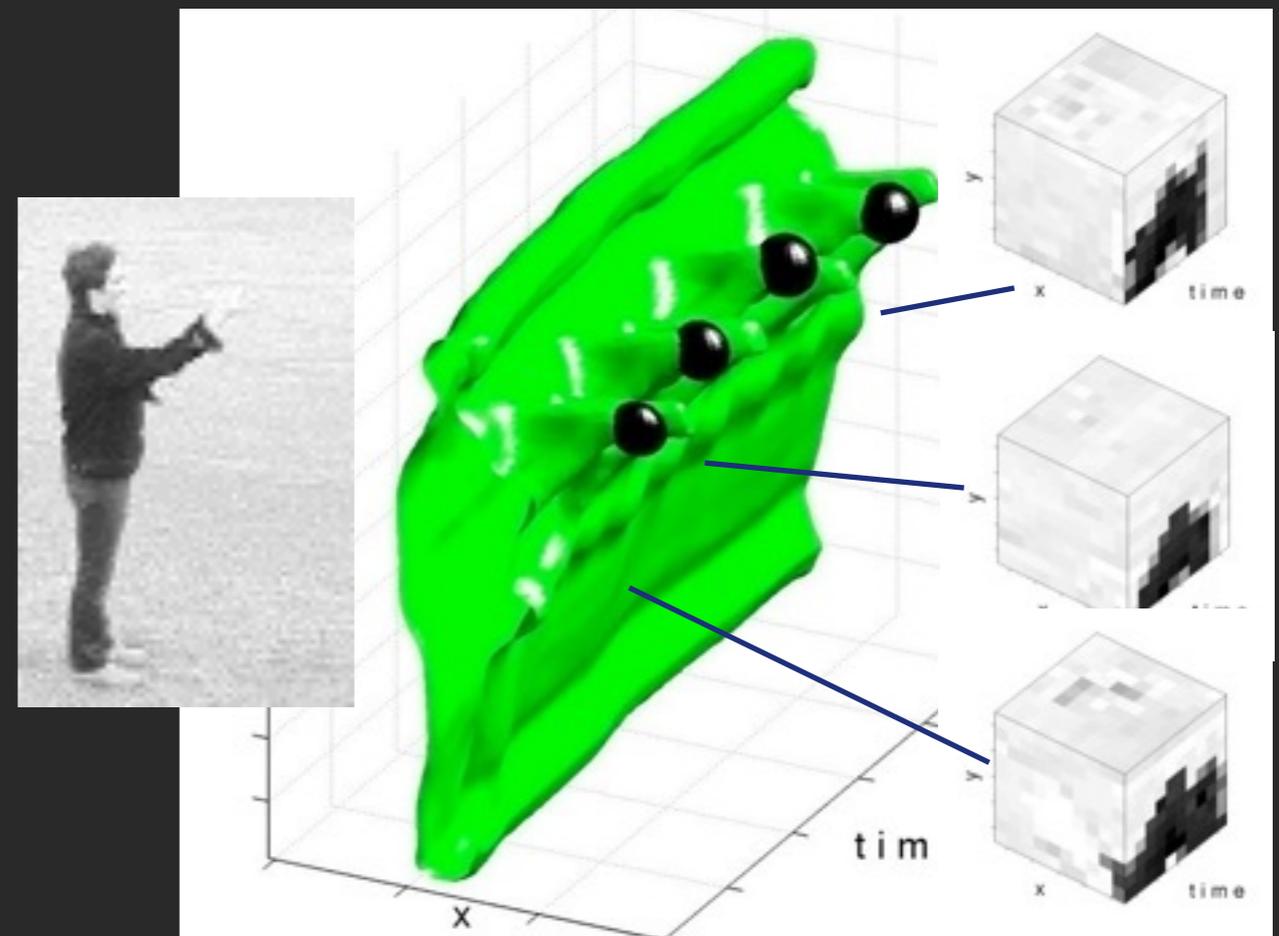
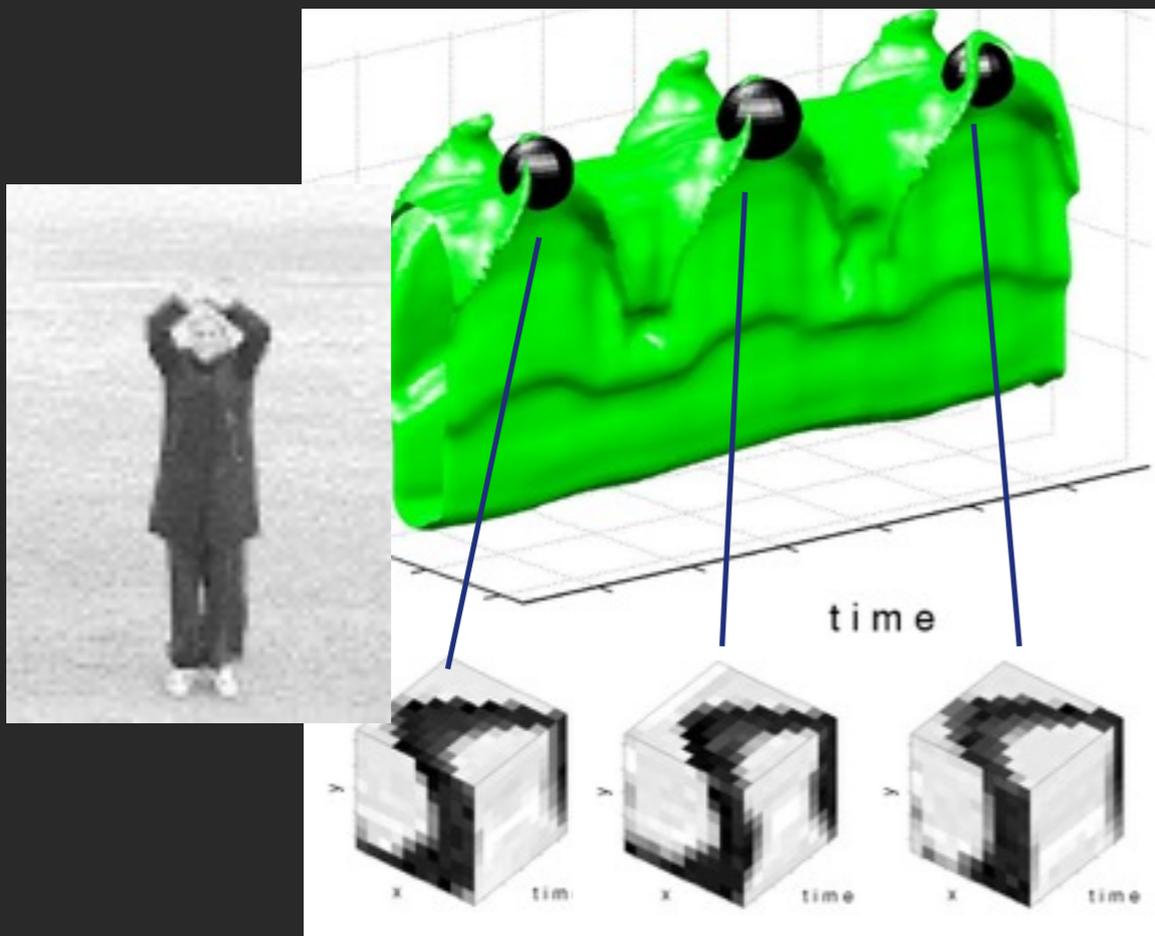
(cf. recent internships from past summer)

Aside: fine-grained activities



Exploiting motion

Spatiotemporal interest points (STIPs)



[Laptev 2005]

(Pre-deep) XYT descriptor evaluation



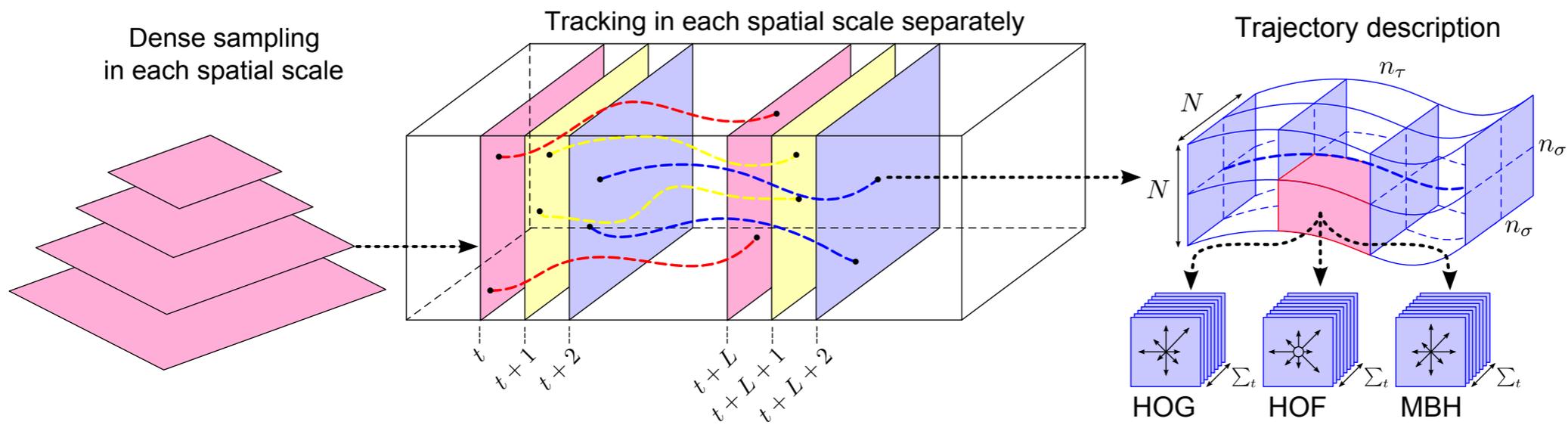
Detection (AP)

Descriptors

	Harris3D	Cuboids	Hessian	Dense
HOG3D	43.7%	45.7%	41.3%	45.3%
HOG/HOF	45.2%	46.2%	46.0%	47.4%
HOG	32.8%	39.4%	36.2%	39.4%
HOF	43.3%	42.9%	43.0%	45.5%
Cuboids	-	45.0%	-	-
E-SURF	-	-	38.2%	-

[Wang, Ullah, Kläser, Laptev, Schmid, 2009]

Compute features along dense trajectories



Action Recognition by Dense Trajectories

Heng Wang, Alexander Kläser, Cordelia Schmid, Liu Cheng-Lin

Capturing the “right” temporal motion



Image motion confounds camera translation, object translation, and nonrigid deformations

Capturing the “right” temporal motion



Image motion confounds camera translation, object translation, and nonrigid deformations



Stabilized camera



Stabilized object



Stabilized camera + object

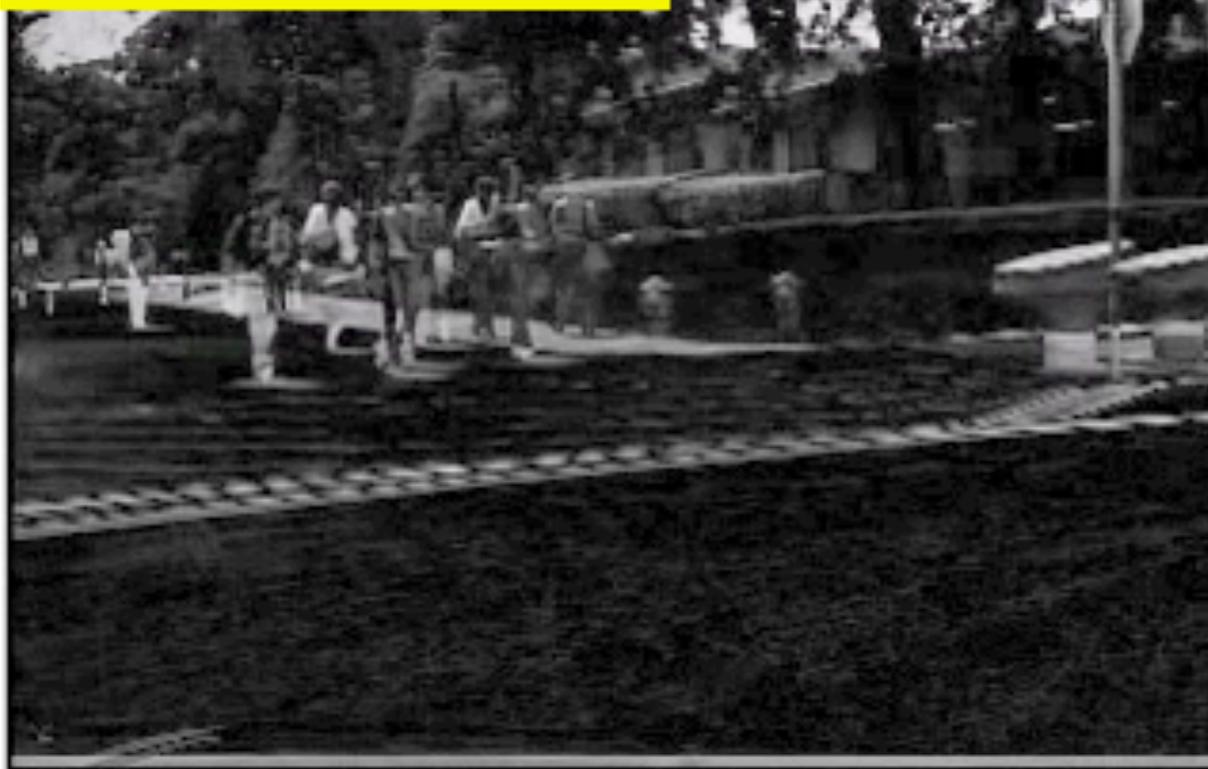
Unstabilized video



Stabilized video



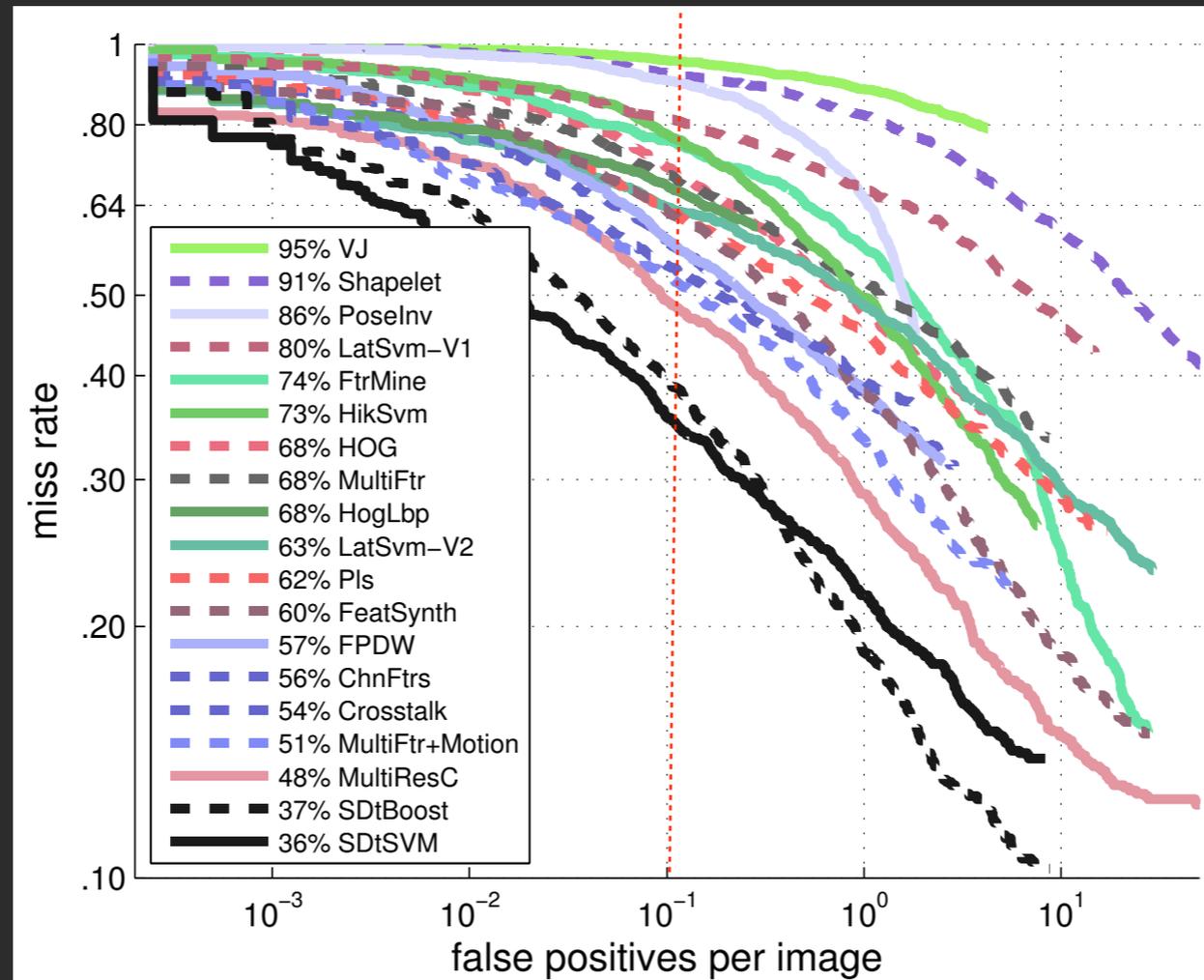
Unstabilized DoF



Stabilized DoF



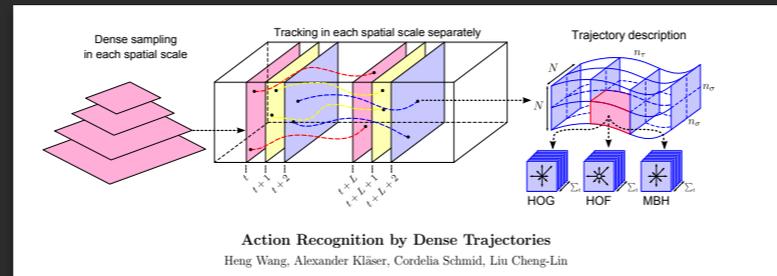
Motion features for detection in videos



Caltech Pedestrian Benchmark; reduce miss rate from 48% to 36%

Park et al CVPR13

Improved dense trajectories



factor out camera motion

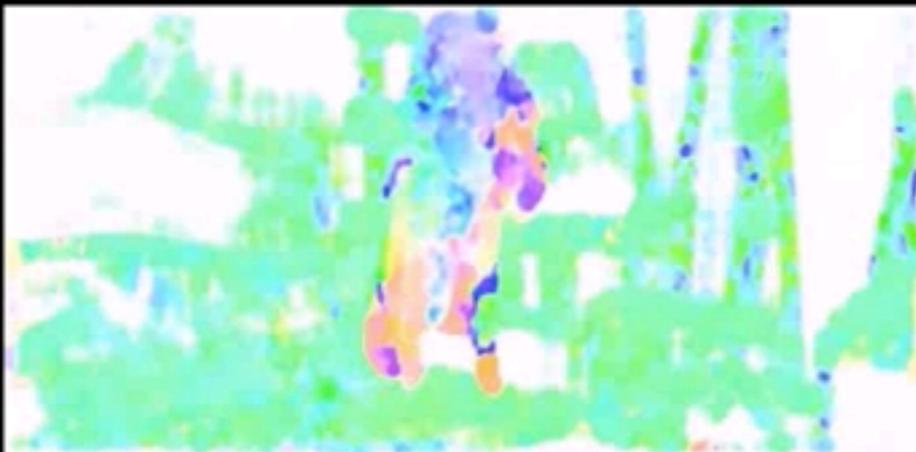
Image



RmTrack



Flow

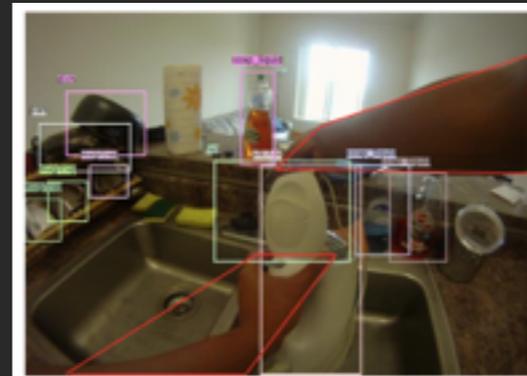


WarpFlow



Roadmap

Data/benchmark analysis

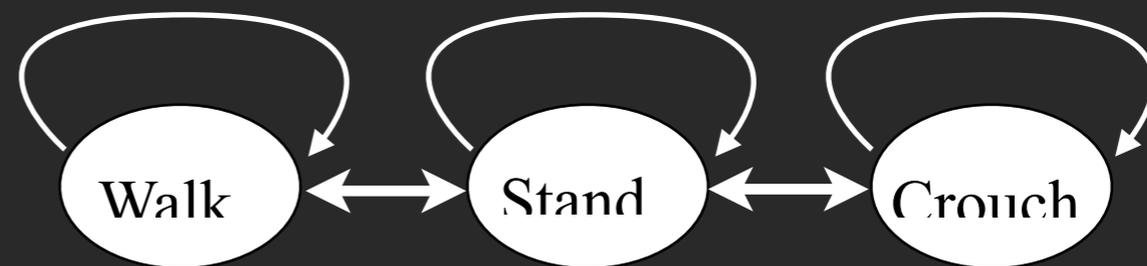


Spatiotemporal features



Motion Features
Tracking

Spatiotemporal models



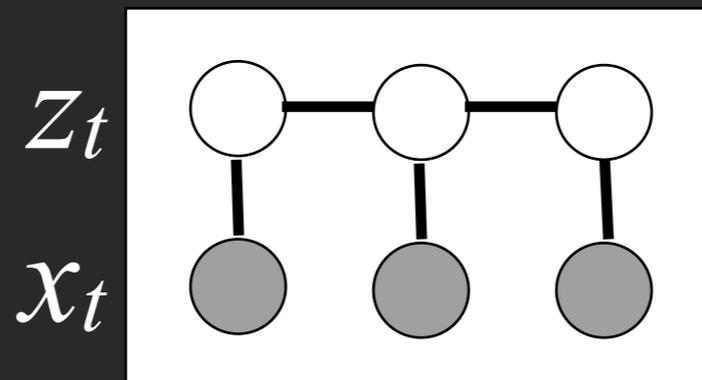
Why do we need to track?



Spacetime window maybe “shearing”

Tracking

Immense literature



$$P(x_{1:T}|z_{1:T}) = \prod_t p(z_t|z_{t-1})p(x_t|z_t)$$

temporal model

appearance model

Historically, last term has been focus of tracking community

Given z_{t-1} , predict z_t with $P(z_t|z_{t-1})$

e.g., particle filtering, Isard & Blake

Extreme form of problem: multi-object tracking



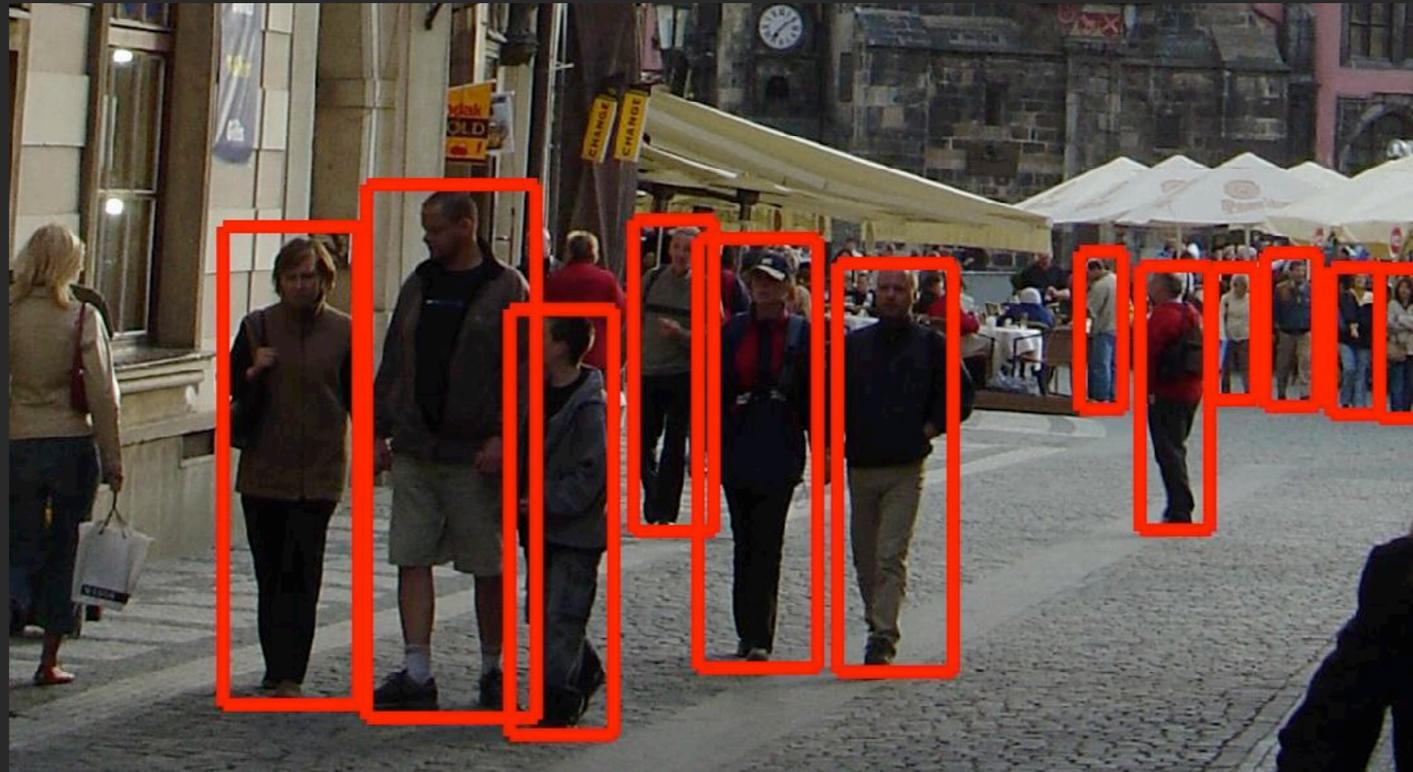
Estimate number of tracks and their extent

Do not assume manual initialization
Estimate birth and death of each track



Analogous to “multi-instance segmentation in spacetime”

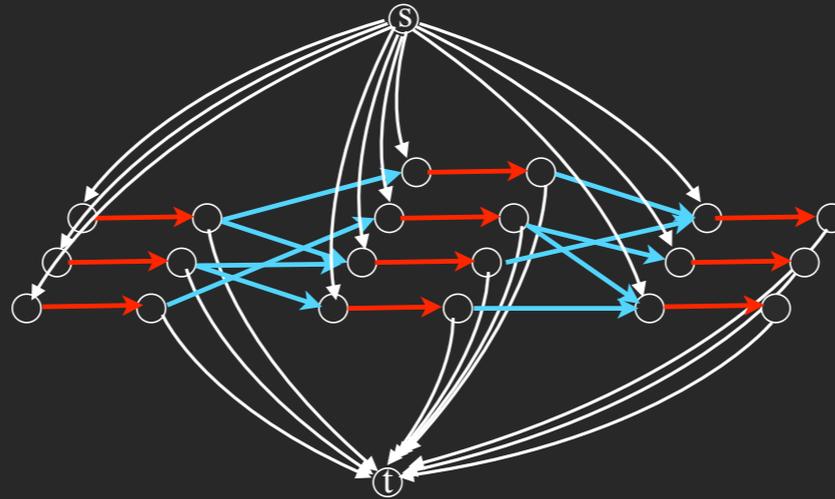
Tracking by detection



Detect candidates

Link detections over time into tracks

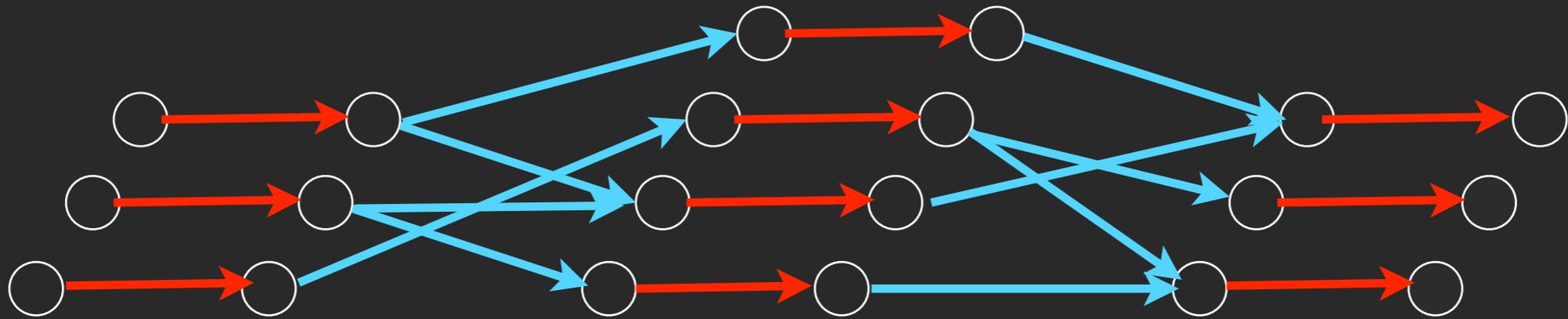
Multi-object tracking as integer/linear programming



View as combinatorial problem of what detections to turn on/off

Jiang et al CVPR07
Zhang et al CVPR08
Berclaz et al PAMI2011
Andriyenko and Schindler ECCV10
Pirsiavash, Ramanan, Fowlkes CVPR11
Butt and Collins CVPR13

Trellis Graph

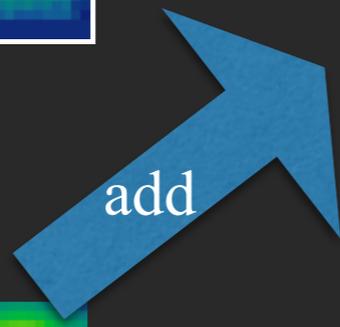
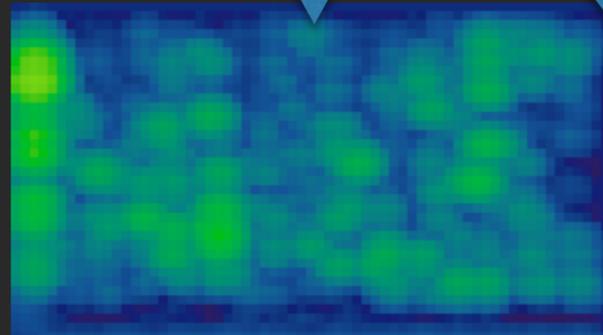
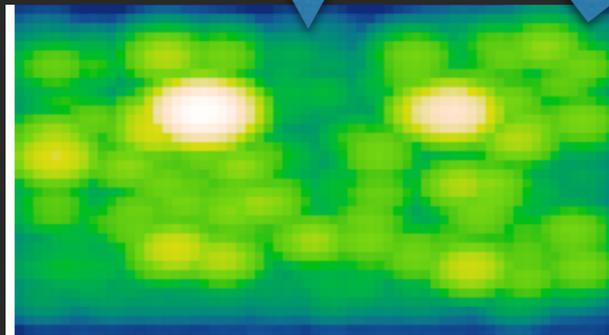
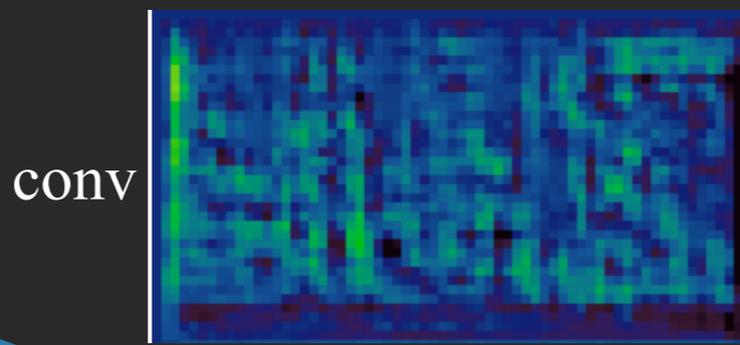
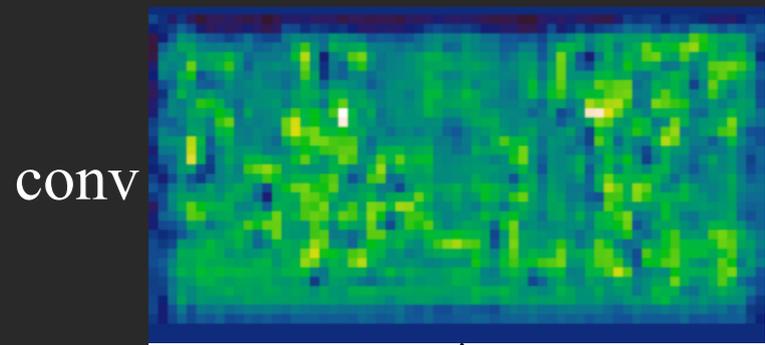
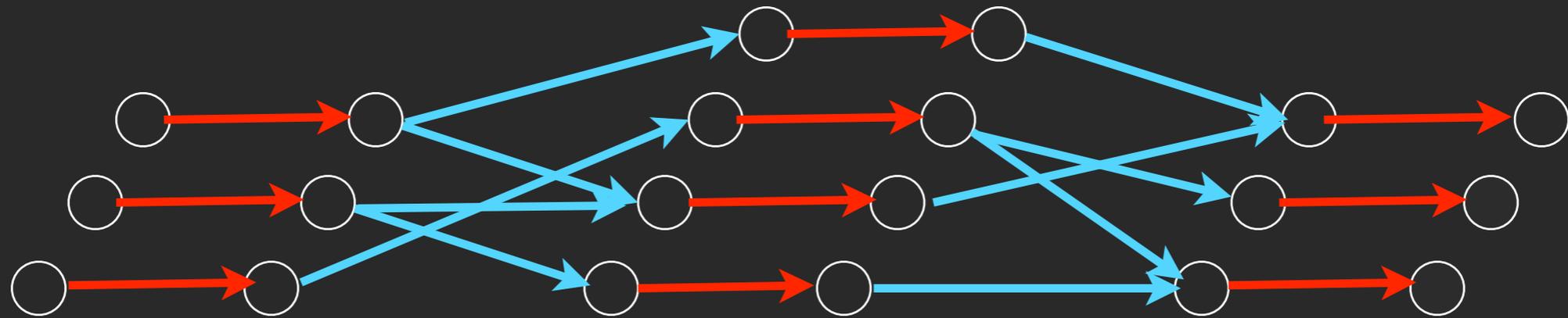


Local cost of window
Pairwise cost of transition

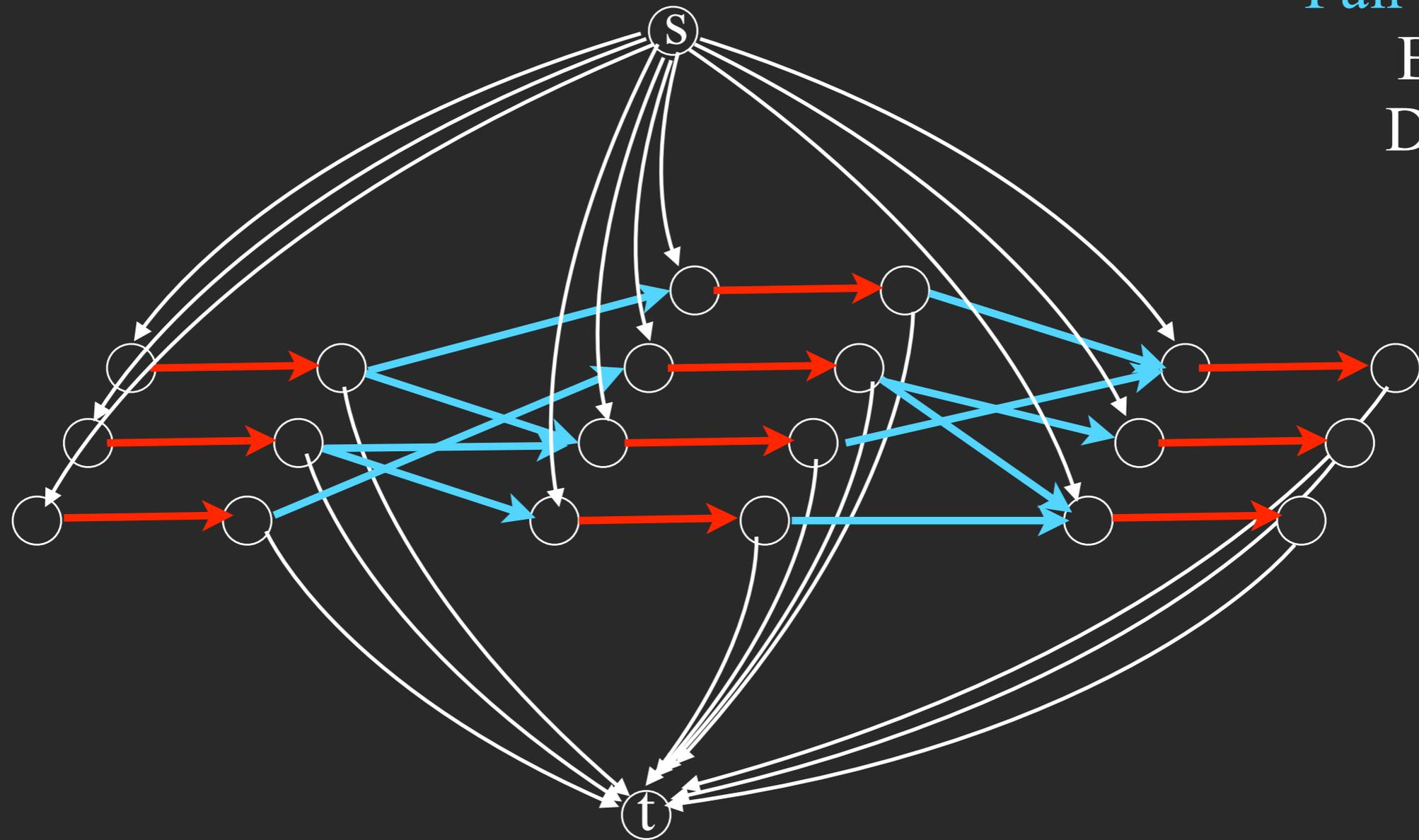
Use dynamic programming (DP) to find single track
(e.g., Viterbi algorithm)

Trellis Graph

Aside: looks a lot like max-pooling and linear convolution



Trellis Graph



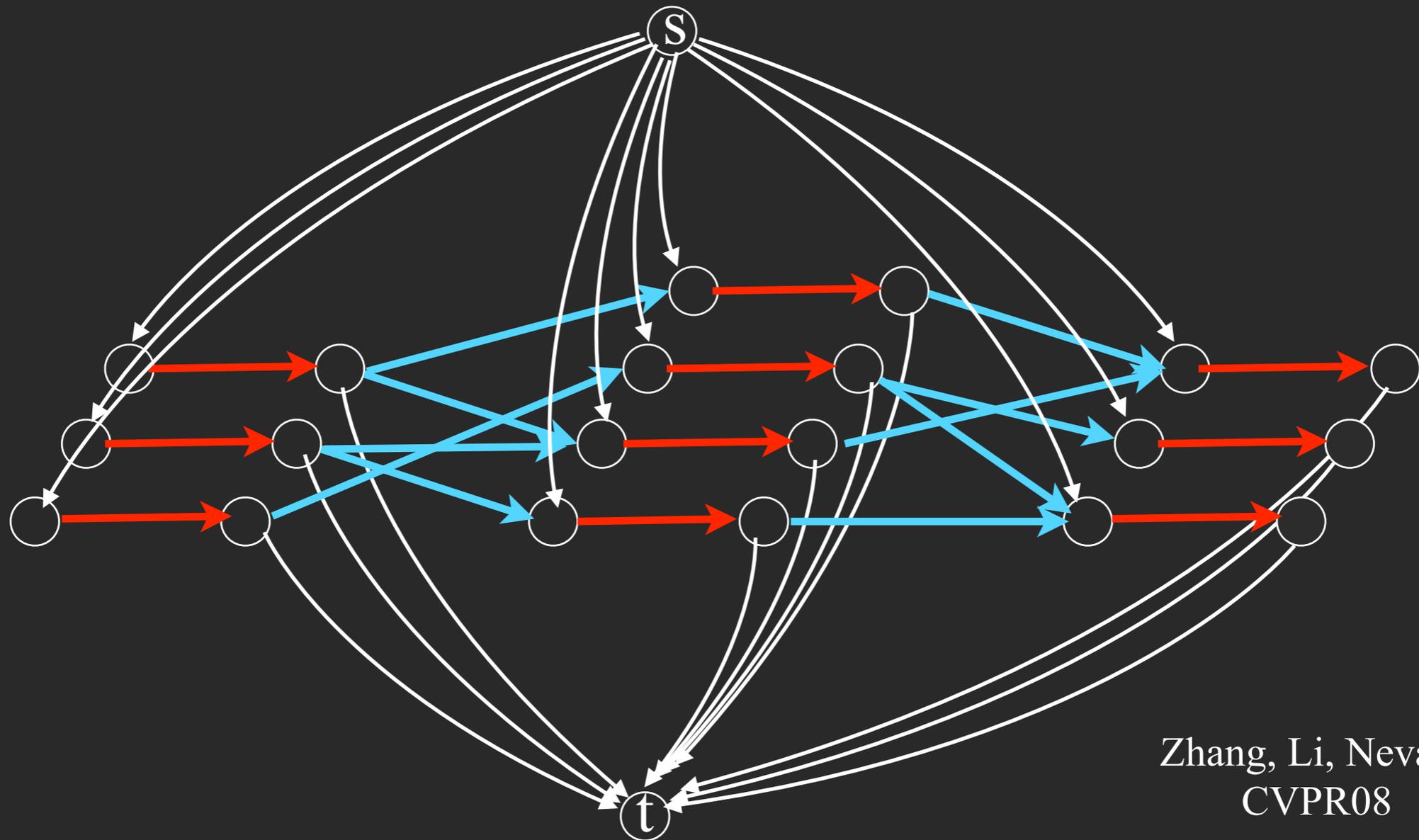
Local cost
Pairwise cost
Birth cost
Death cost

Shortest path from S to T = best variable-length track

Still can use DP

Min-cost flow problem

(generalization of min-cut / max-flow)



Zhang, Li, Nevatia
CVPR08

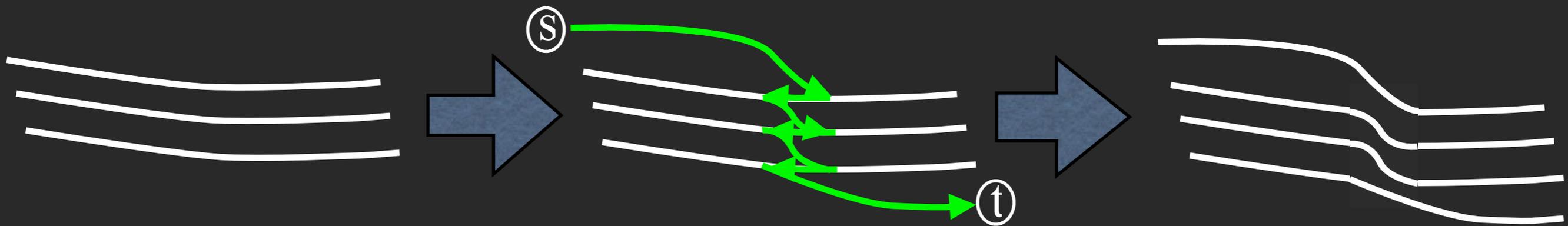
Cost of a K -unit flow = sum of flow along each edge * cost

- 1) Capacity along each edge is 1
- 2) Sum of flow into a node = sum of flow out
(ensures non-overlapping tracks)

Exact solution for $K > 1$

Problem: once we instantiate a track, we cannot edit it

Solution: compute shortest path on **residual graph** augmented with reserve edges



New tracks can “suck flow” out of existing tracks

Keep repeating until next instantiated track increases cost

Okay... so what about tracking articulations?



Which one is correct?

What **should** a single-image pose estimation alg. output?

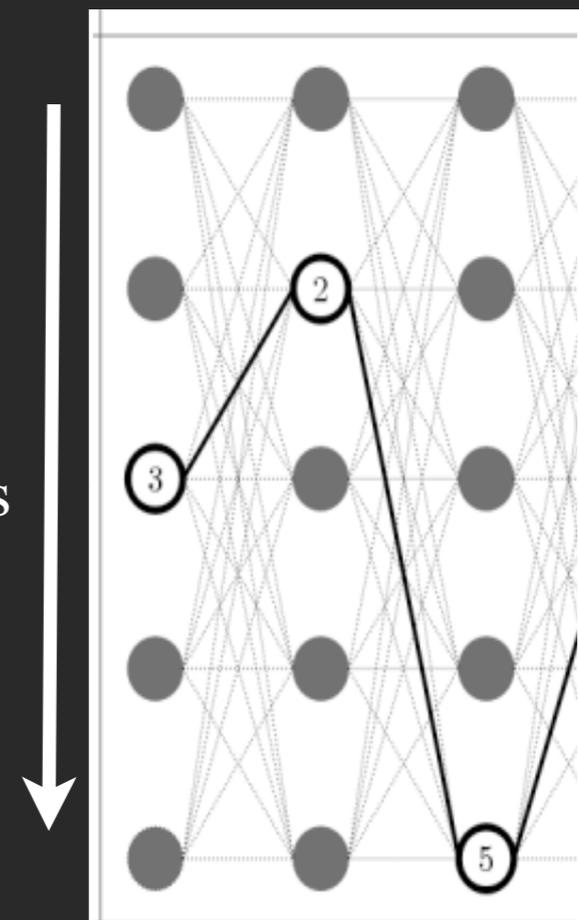
N-best decoding

Generate N high-scoring candidates with simple (tree) model, and evaluate with complex model

Popular in speech, but why not vision?

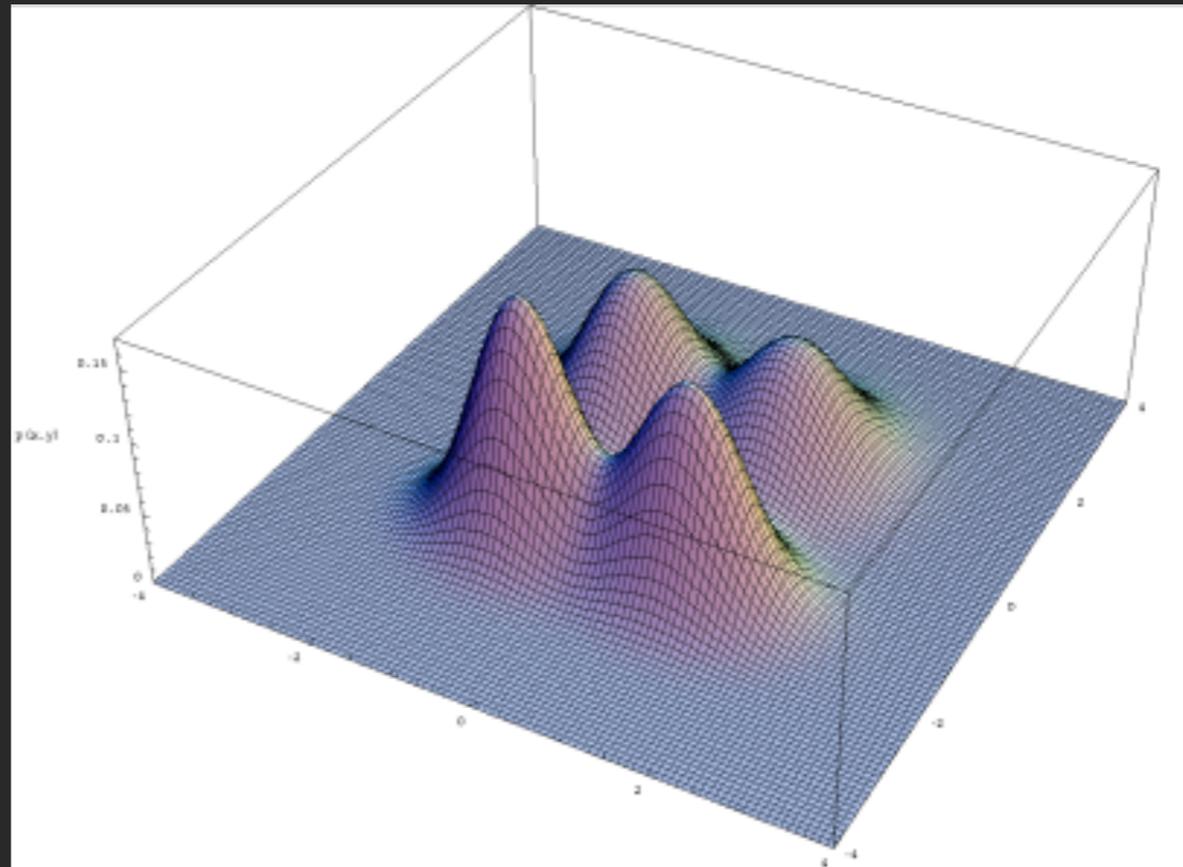


Pixel locations



head torso leg

N-best maximal decoding



N-best with “NMS” or “mode-finding”

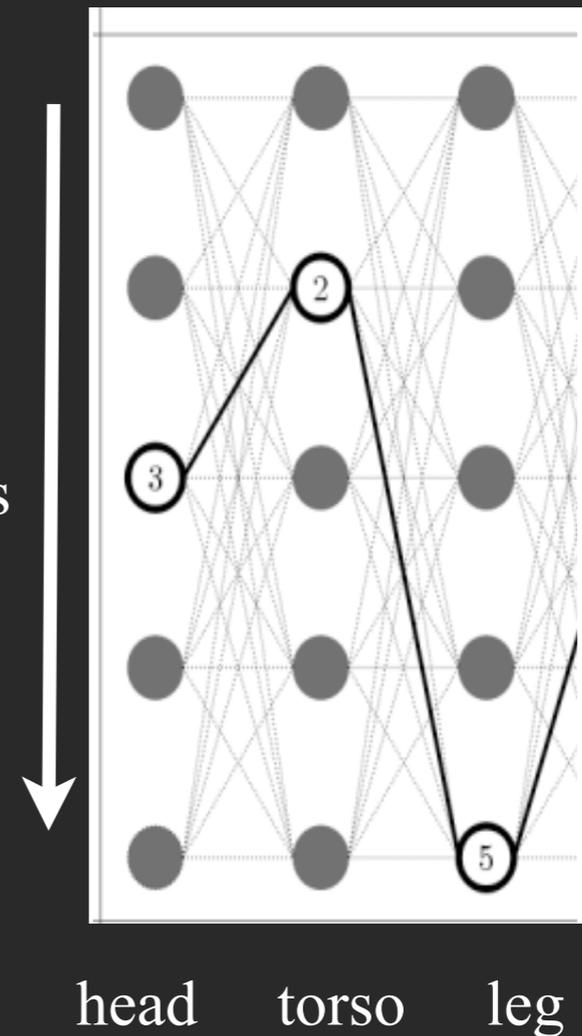
Park and Ramanan, ICCV11

Yadollahpour et al. ECCV12

N-best maximal decoding



Pixel
locations



Intuition: backtrack from all part “max-marginals”, not just root

(can we done without any noticeable increase in computation)

N-best maximal decoding

Park & Ramanan, "N-best decoders for part models" ICCV 2011



Find N-best "modes" rather than N-best poses

Philosophy: Delay hard decisions as much as possible

~~Candidate interest points~~

~~Candidate parts~~

Candidate poses

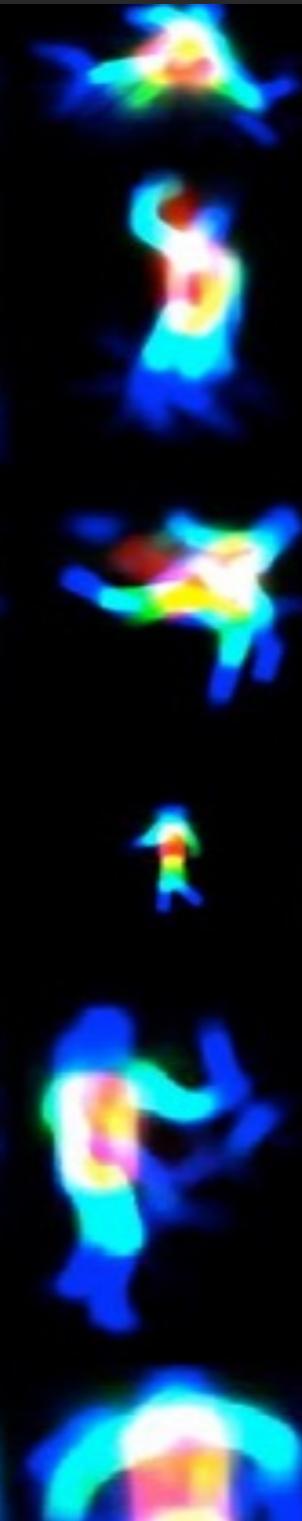
Maximal poses from a single frame



Correct one picked out by temporal context (tracker)

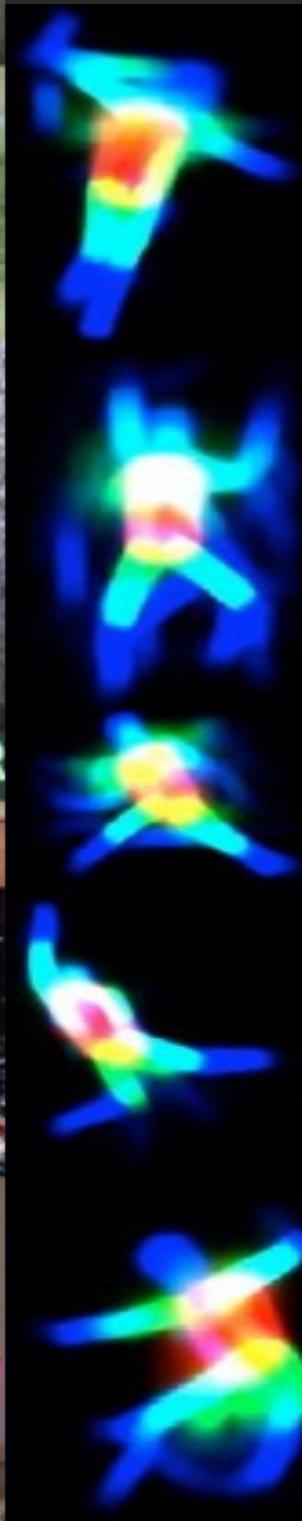
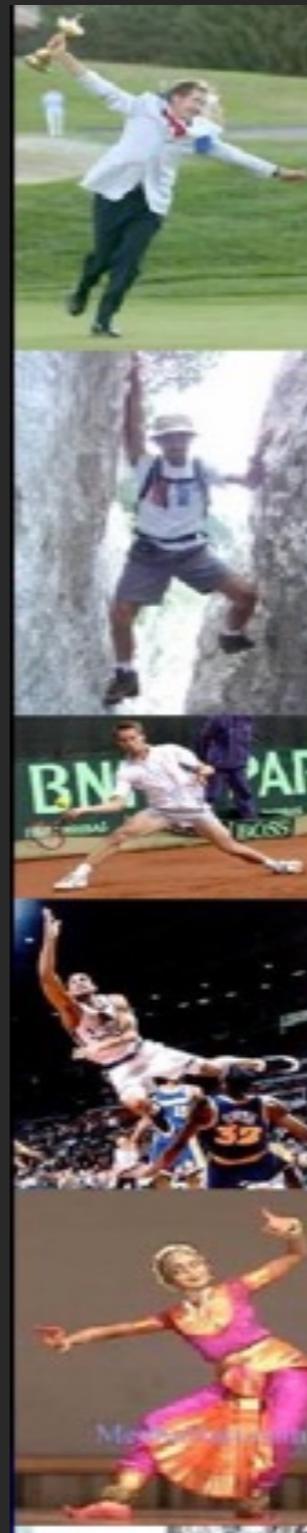
Aside: other ways of representing uncertainty

$$P(z|x) \propto e^{S(x,z)}$$



Log-linear
conditional models

Ramanan NIPS 06



Tracking by articulated detection



Problem: linking up these detections won't work

Recall: Why is finding people difficult?



variation in appearance



variation in pose and viewpoint



occlusion & clutter

Classic “nuisance factors” in image recognition

Recall: Why is finding people difficult?



~~variation in appearance~~



variation in pose and viewpoint

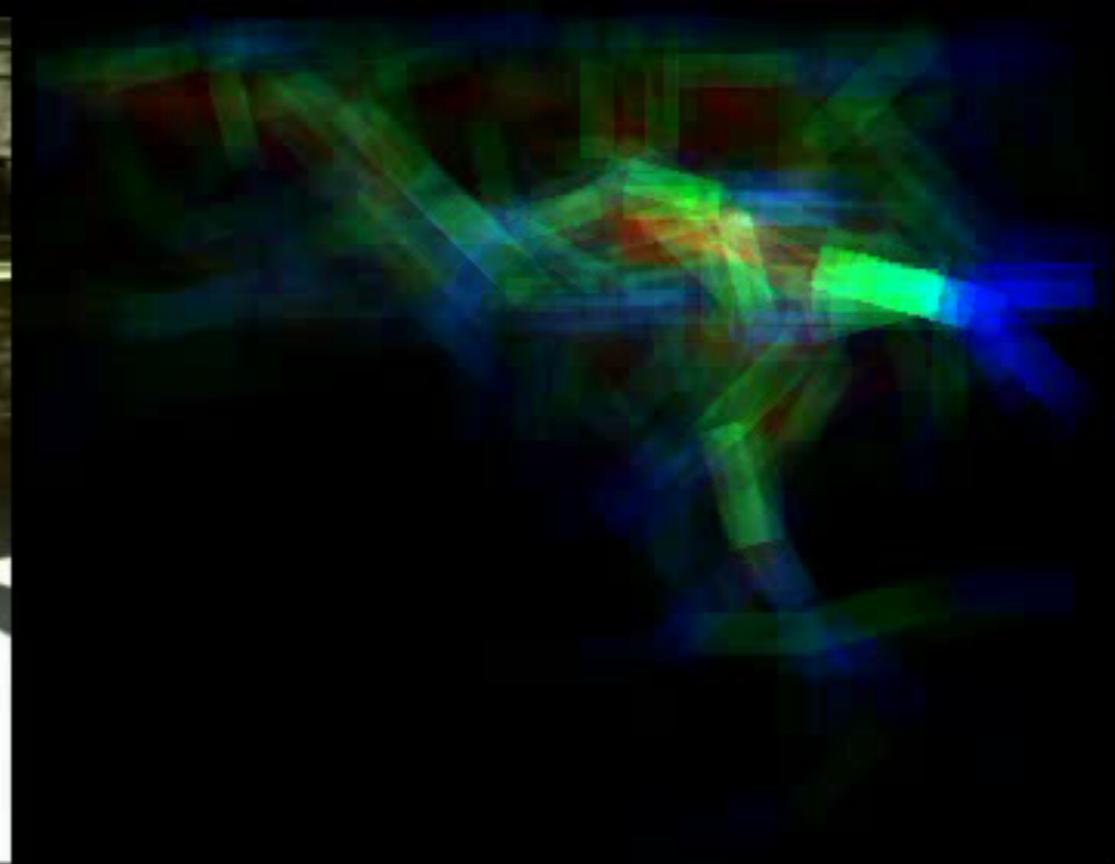
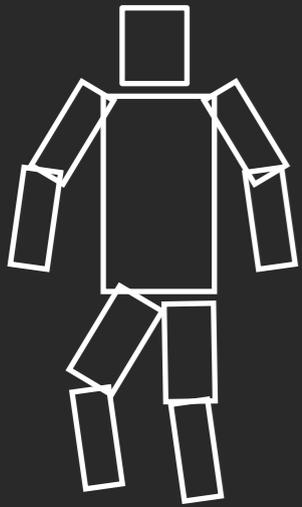


occlusion & clutter

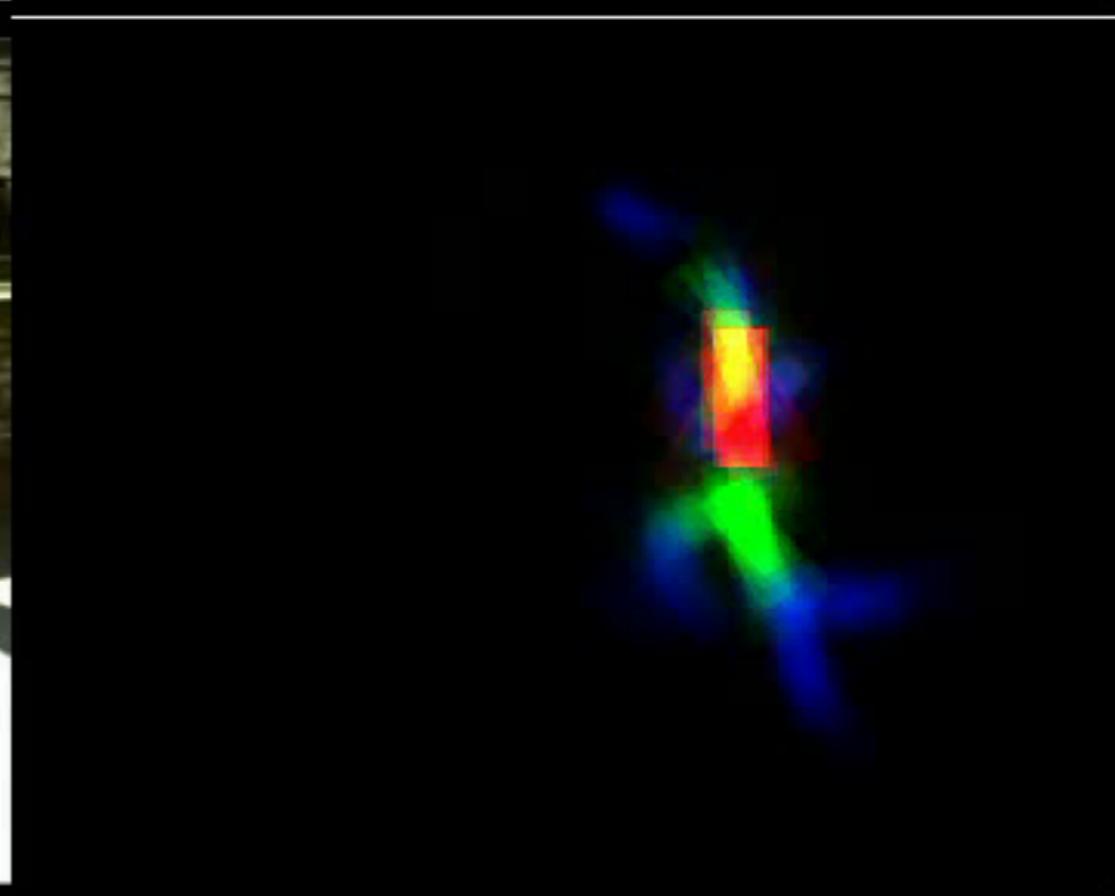
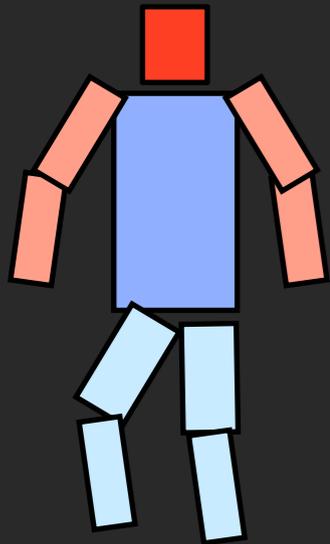
Classic “nuisance factors” in image recognition

Tracking by repeated detection

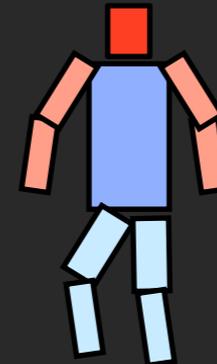
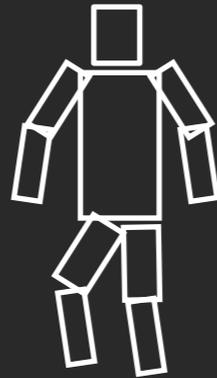
Generic
Person
Template



'Lola'
Template



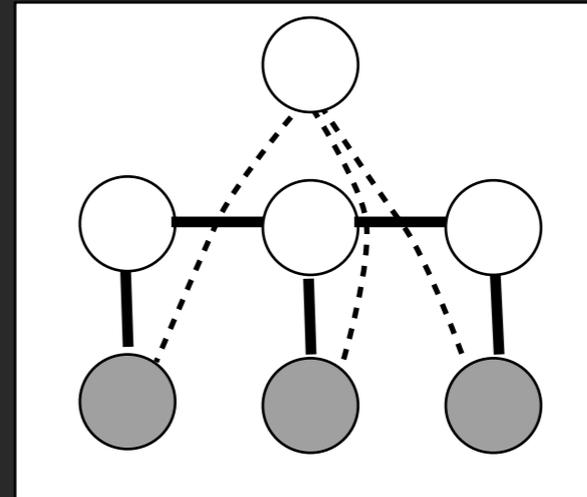
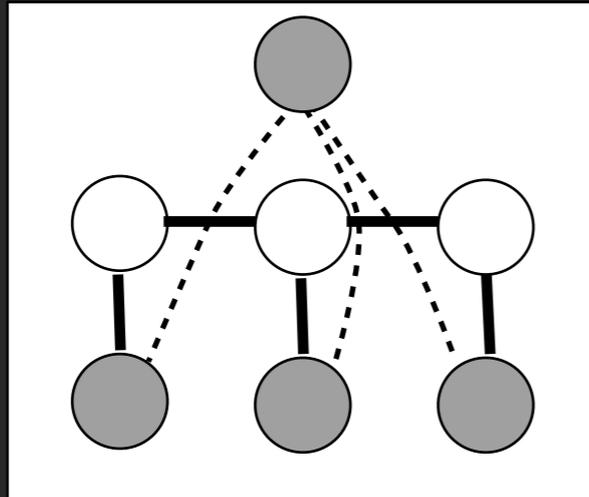
Tracking as model-building



Model-based tracking

Latent-variable tracking

ORourke & Badler 80
Hogg 83
Rehg & Kanade 95
Ioffe & Forsyth 01
Toyama & Blake 01
Sigal et al. 04



Ramanan et al.
PAMI 07

A generic object template must be **invariant**

We want to build a model of the object as we track it

Track through occlusions



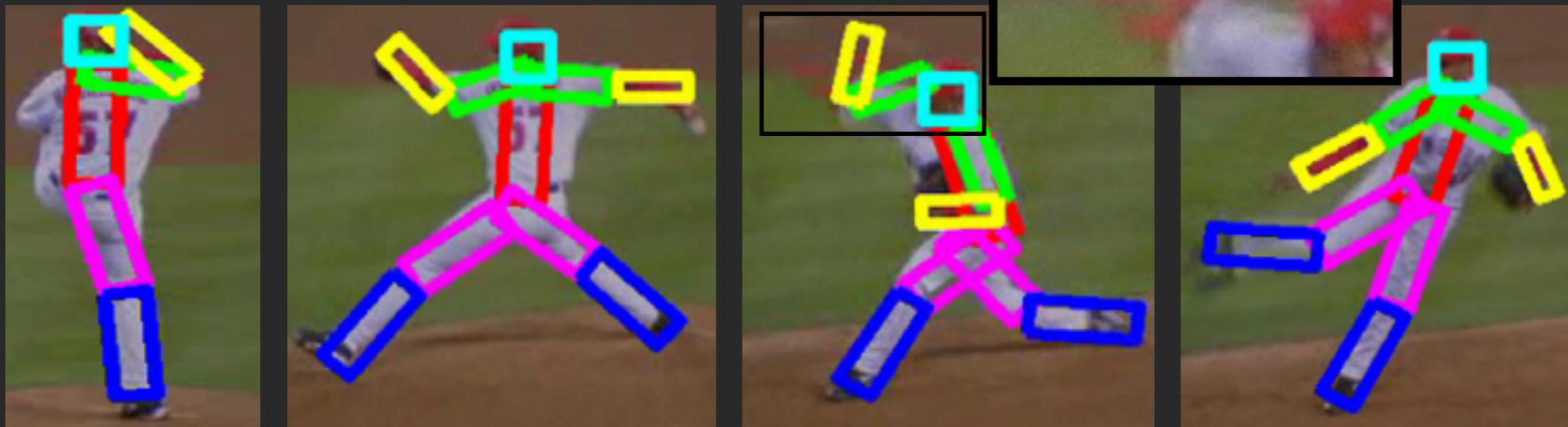
Ramanan, Forsyth, and Zisserman PAMI 07



Discriminative clothing models



2002 World Series

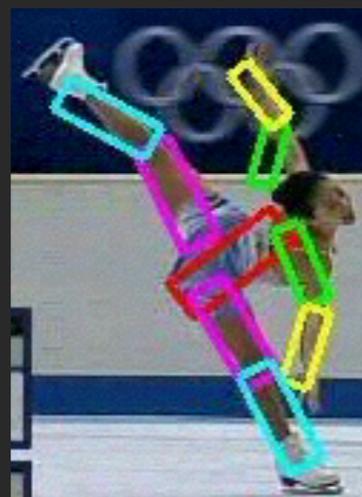


motion blur & interlacing

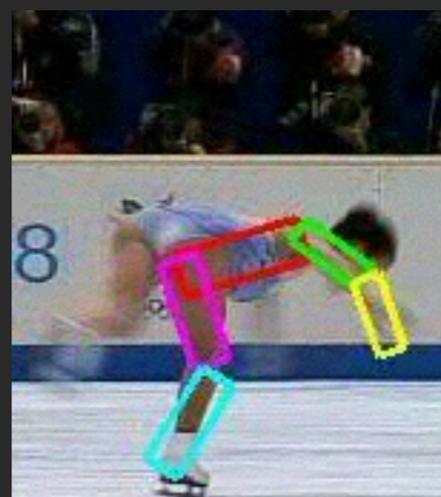
Track long footage (10,000 frames)



Michelle Kwan, 1998 Olympics



extreme pose



motion blur



fast movement



Olympic woes

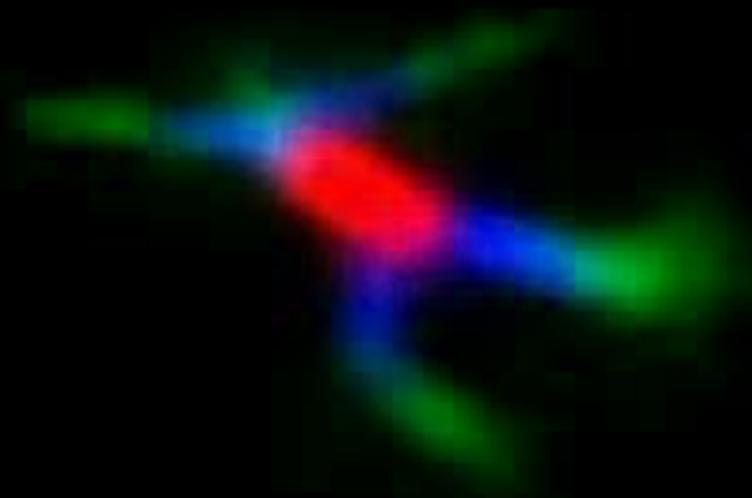
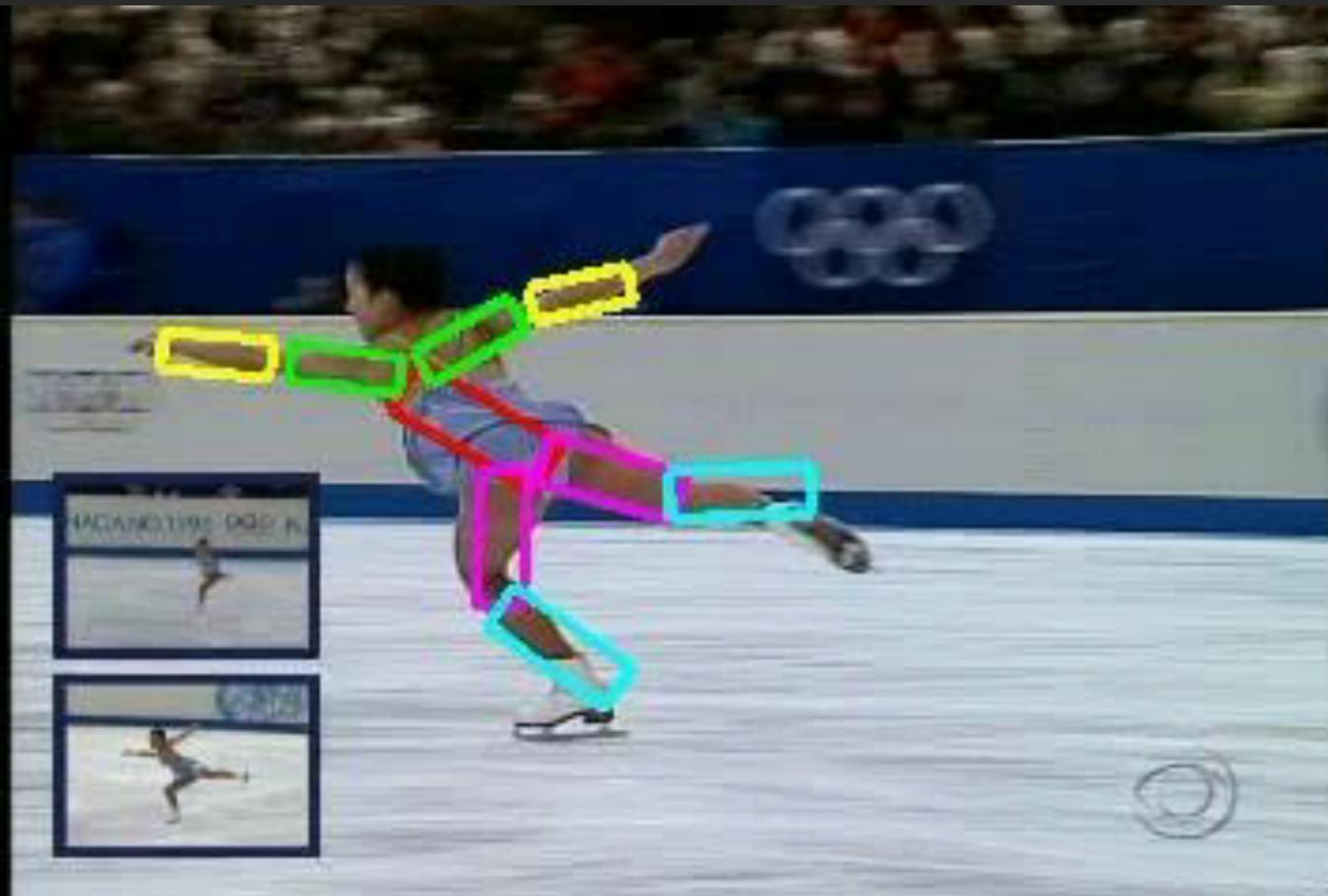
silver, not gold →



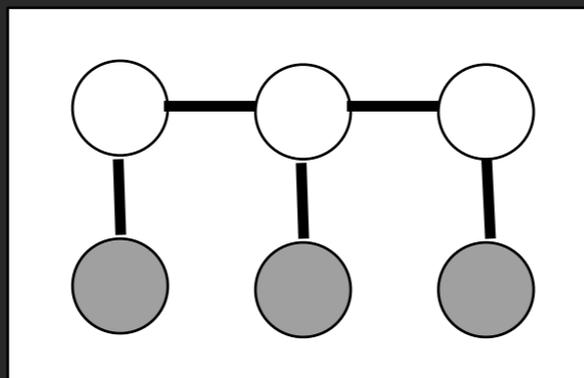
Kwan led after the short program. In the long program, skating to Lyra Angelica by the British composer William Awyn, the 17-year-old turned in a clean, if cautious, effort. Kwan didn't make a major error -- with only one slight wobble on a triple jump -- earning her a solid row of 5.9s on presentation from the judges. As flowers rained upon the ice from her fans, the gold medal, it seemed, was hers. Still, her conservative routine earned five 5.7s for technical merit, and the door was opened, however slight, for Lipinski.

http://espn.go.com/classic/biography/s/Kwan_Michelle.html

The culprit

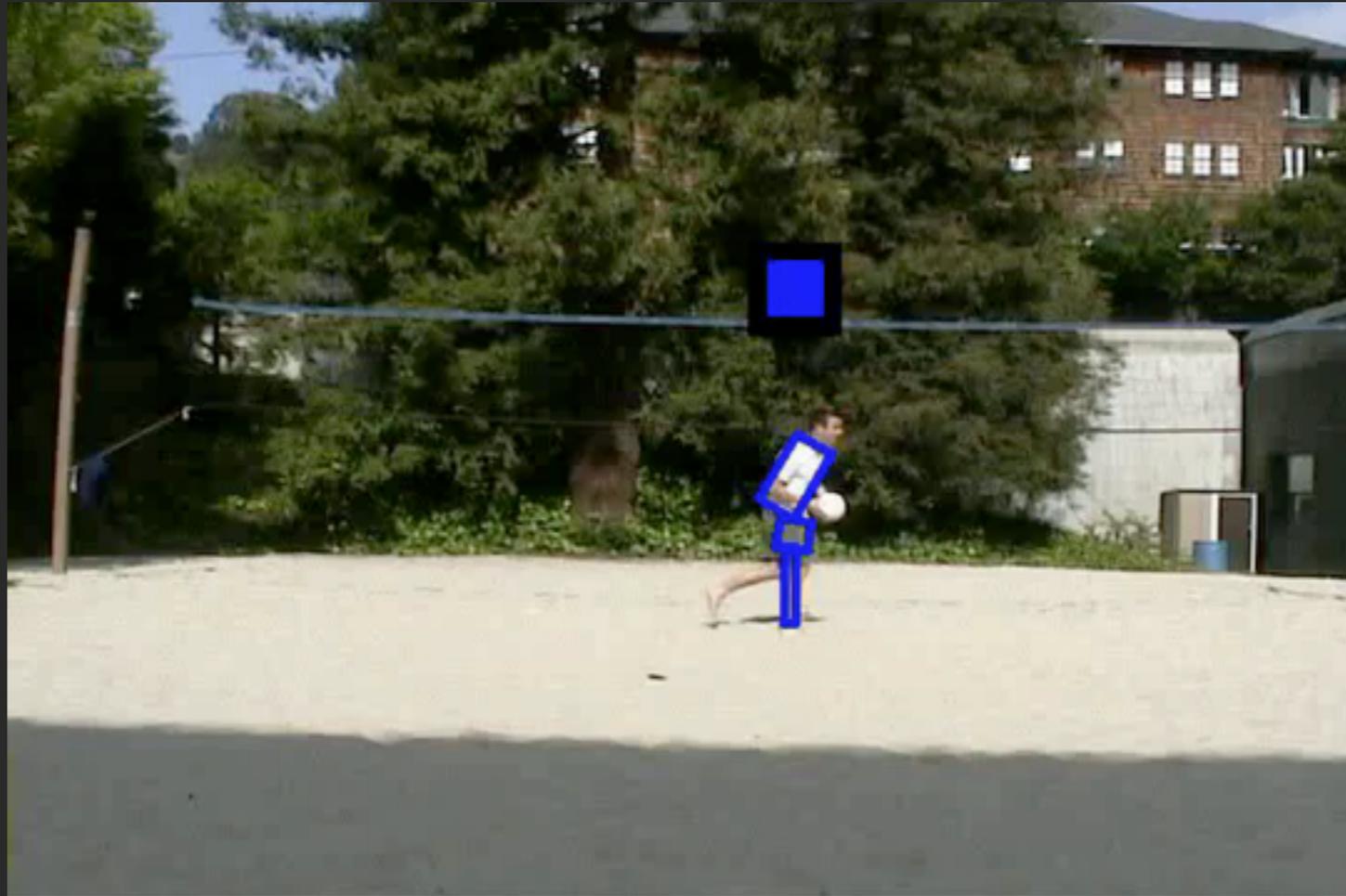


Unexpected/unlikely motions often **very** important
The motion prior $P(z_{t+1}|z_t)$ may smooth out such subtleties

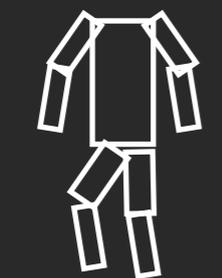


Tracking multiple people

Independently track each figure



Clothing appearance is no longer a nuisance



person
detector



Deva
detector



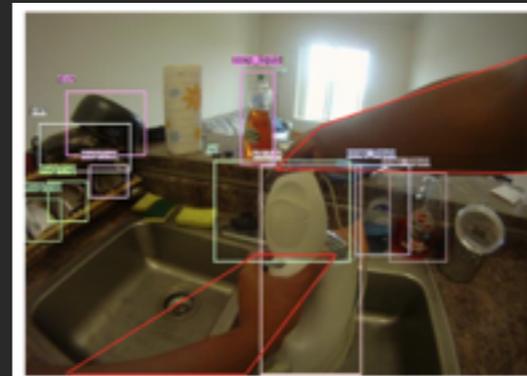
Bryan
detector



John
detector

Roadmap

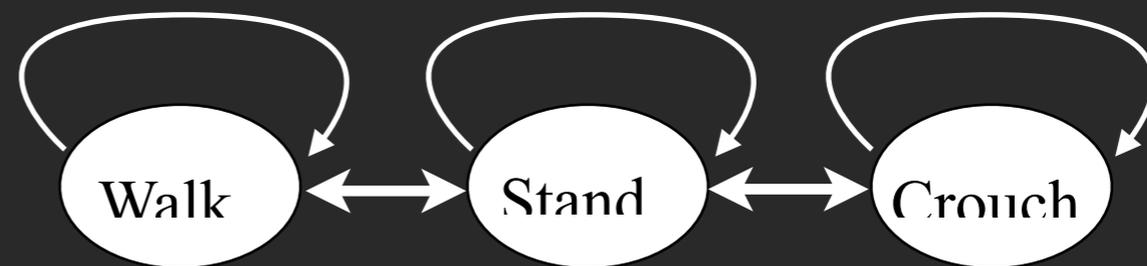
Data/benchmark analysis



Spatiotemporal features



Spatiotemporal models

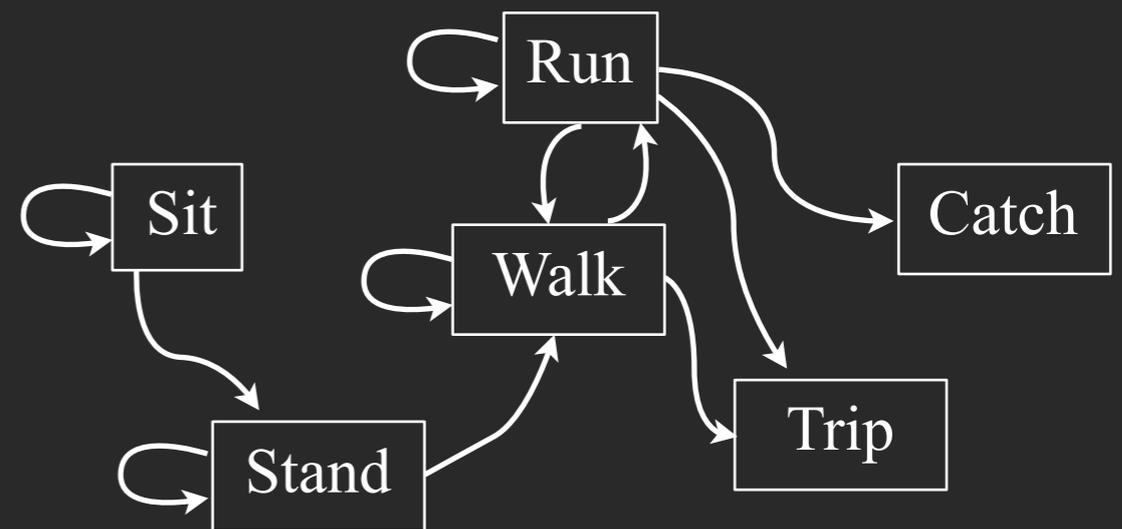


Spatiotemporal models

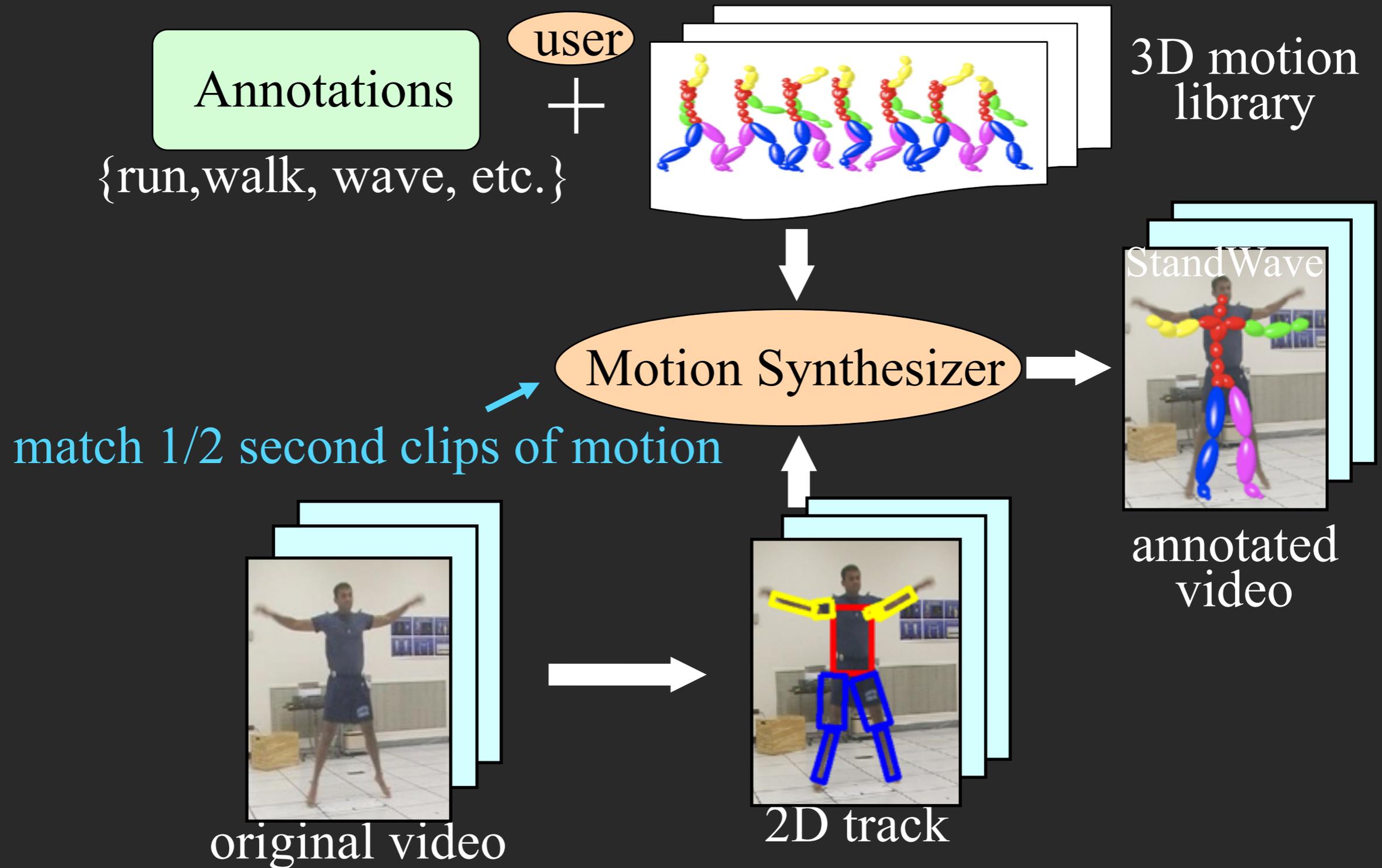
Data-driven



Model-driven

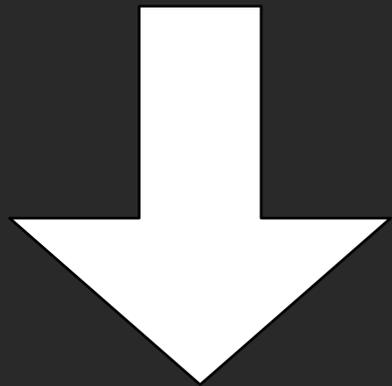


Data-driven action recognition

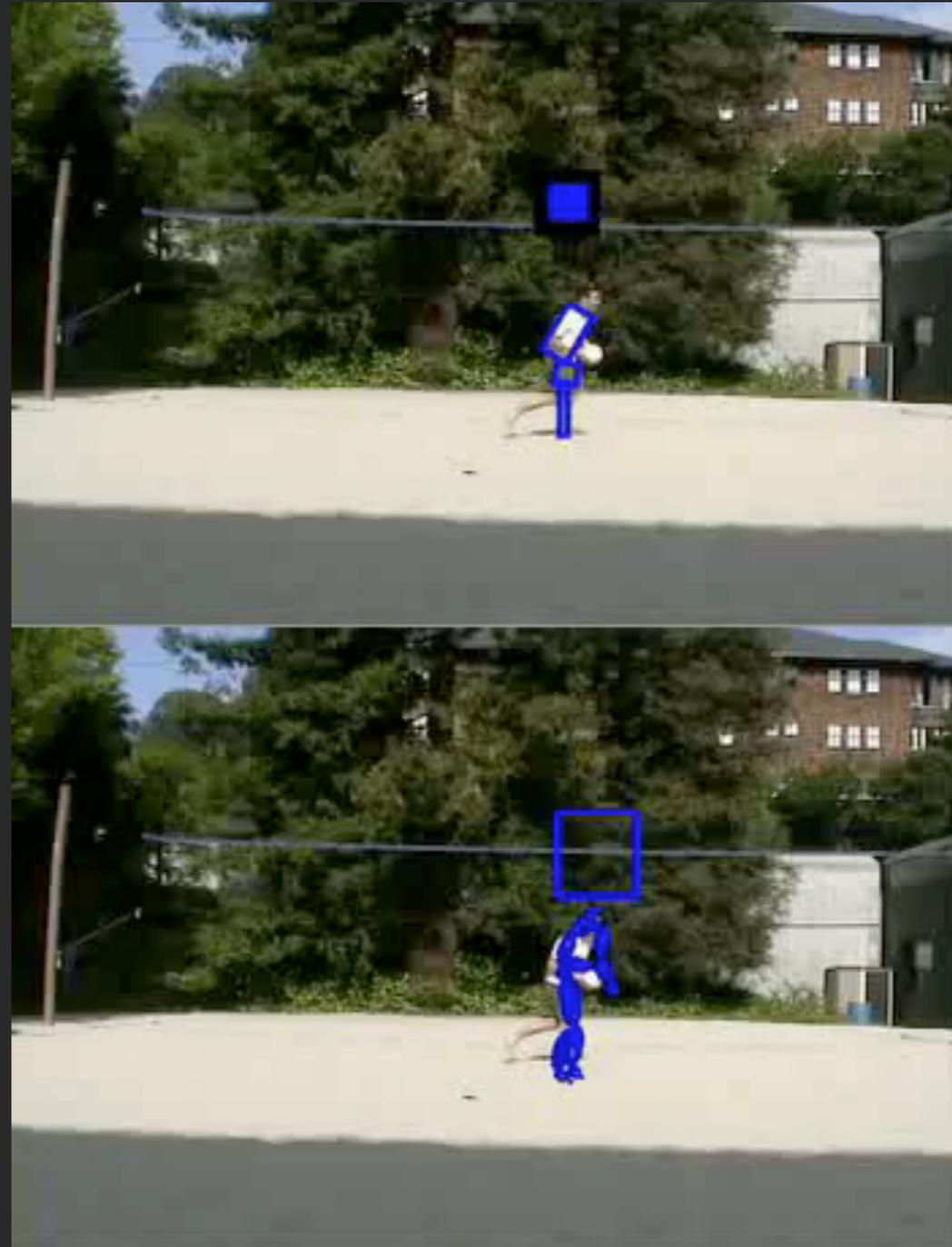


Results

(Low-level)
tracking



(High-level)
spatiotemporal
models



Pipeline surprisingly rare (e.g., doesn't work on TrecVid)

Recognizing structured actions

Making tea from a wearable camera



Start boiling
water

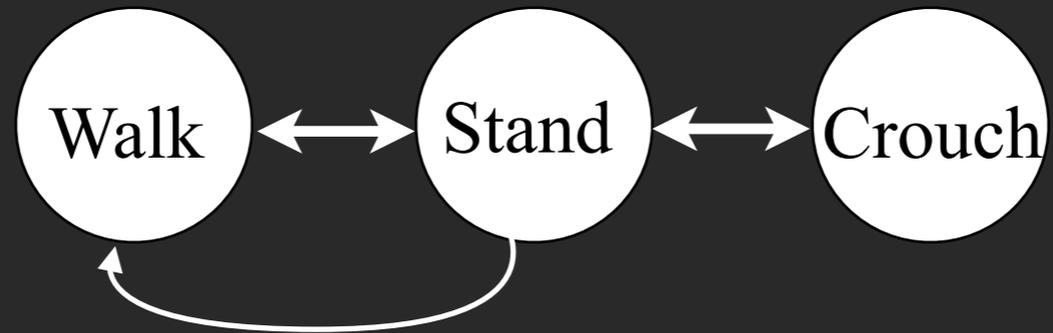
Do other things
(while waiting)

Pour in cup

Drink tea

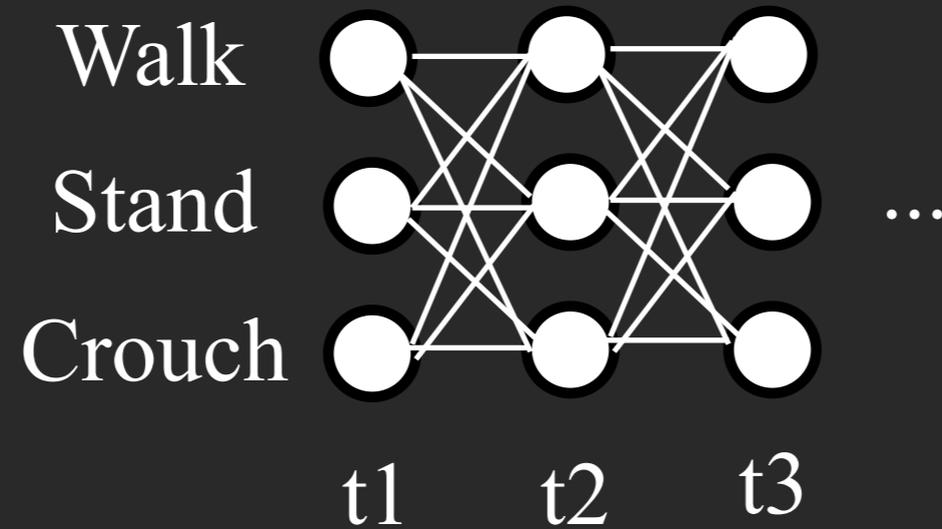
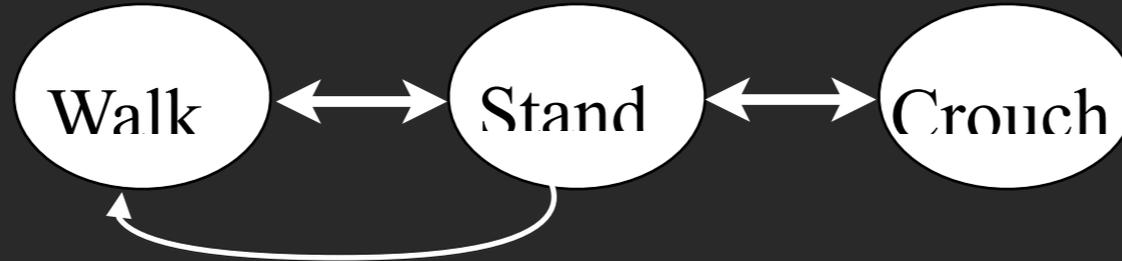


How do we capture long-term structure?

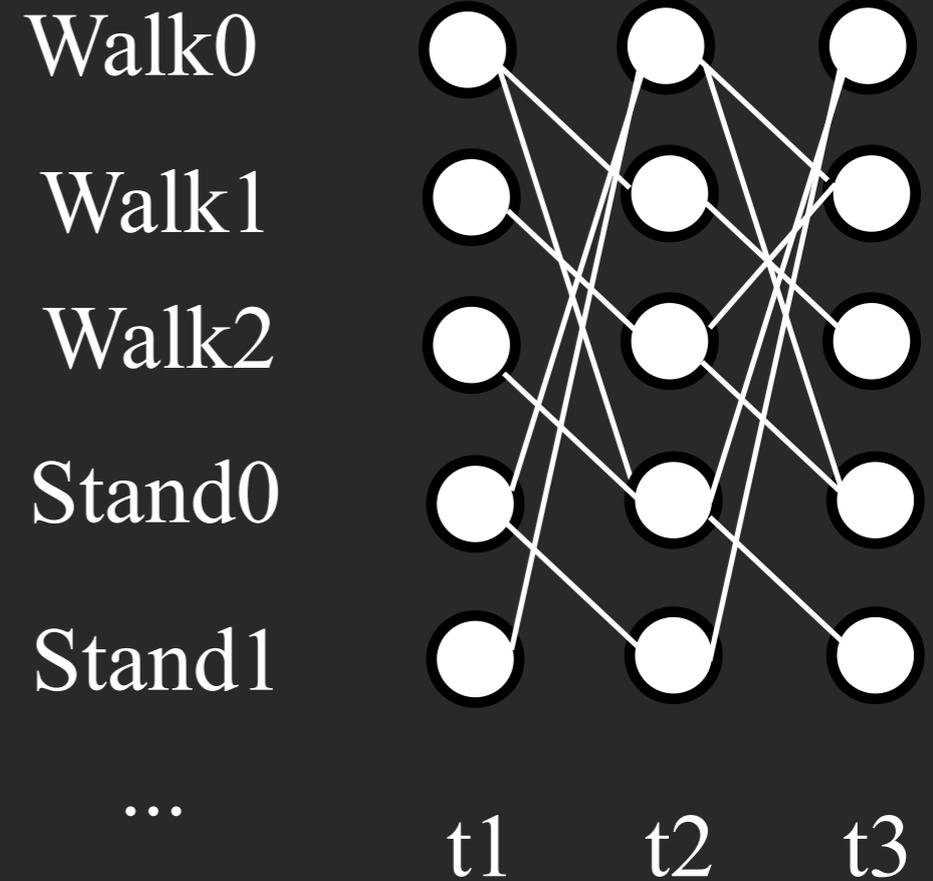
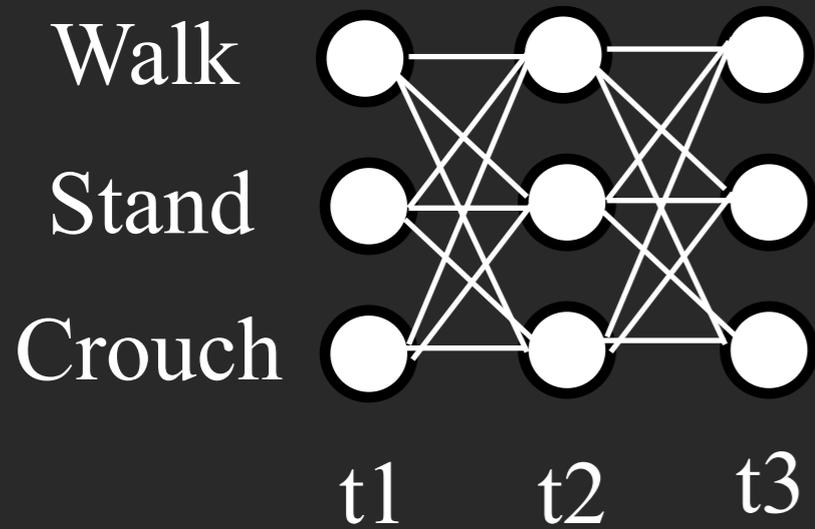


Markov models

What's magic behind semi-markov models?

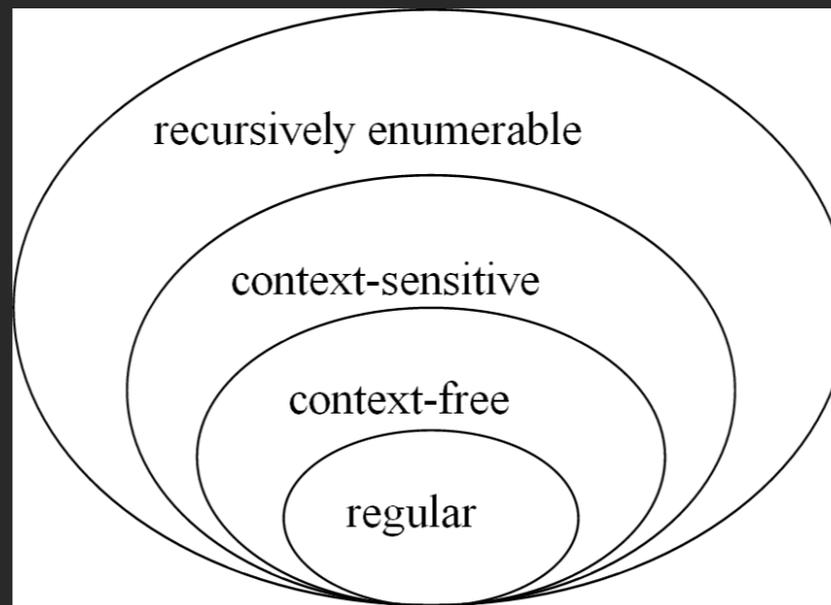


Semi-markov models

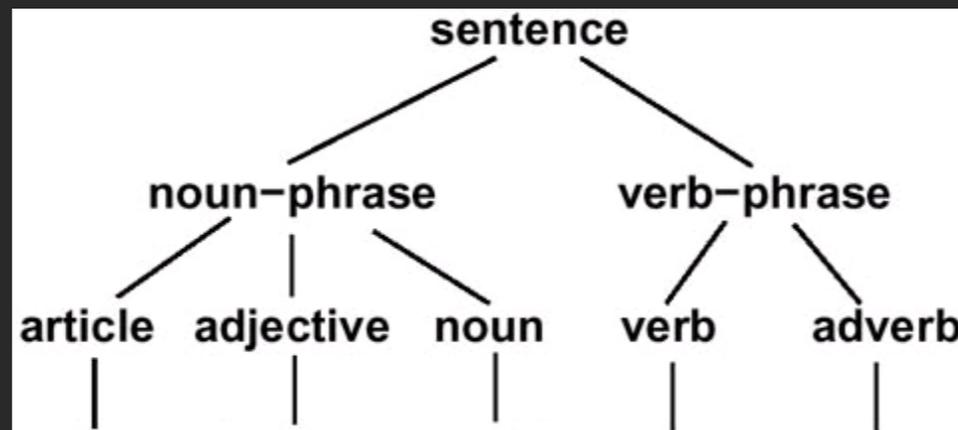


Add counting states and force sparse transitions (Walk0 to Walk1)
Counting state costs can model arbitrary priors over segment lengths

How do we capture long-term structure?



Exploit models for language



“The hungry rabbit eats quickly”

Context-free grammar

Example grammar



“yank”



“pause”



“press”



“background”

Clean&Jerk action = 

Snatch action = 

$S \rightarrow \square$

$S \rightarrow S \text{    $

$S \rightarrow S \text{   $

Example parse



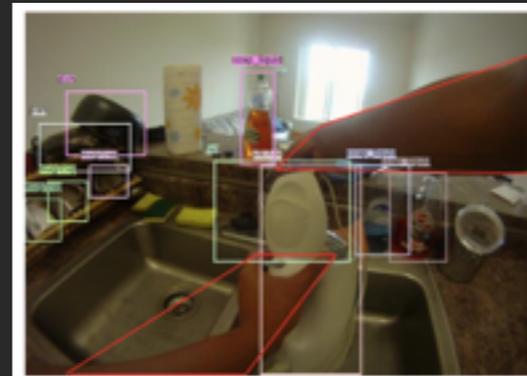
Zhu et al
Bobick et al

time

Grammars

A look back

Data/benchmark analysis



Spatiotemporal features



Spatiotemporal models

