# Setting up Filtering

D.A. Forsyth, UIUC

# Knowing motion is a cue

- Localization:
  - Match observation to map - this tells us where we are
  - Now imagine we move; there are two cues to location:
    - what the map says;
    - how we moved (from, eg, shaft encoders, etc.)
- Tracking:
  - Estimate state of object in one frame
  - Now we receive a second frame; there are two cues to object state:
    - estimate in new frame;
    - use a motion model (eg moving under gravity)
- Key technical question:
  - Fuse these two cues into a single estimate recursively

# Tracking: Why?

- Motion capture
  - build models of moving people from video
- Recognition from motion
  - eg cyclists move differently than runners
- Surveillance
  - who is doing what?
    - for security (eg keep people out of sensitive areas in airports)
    - for HCI (eg kinect, eyetoy, etc.)

# Tracking

- Establish state of object using time sequence
  - state could be:
    - position; position+velocity; position+velocity+acceleration
    - or more complex, eg all joint angles for a person
  - Biggest problem -- Data Association
    - which image pixels are informative, which are not?
- Key ideas
  - Tracking by detection
    - if we know what an object looks like, that selects the pixels to use
  - Tracking through flow
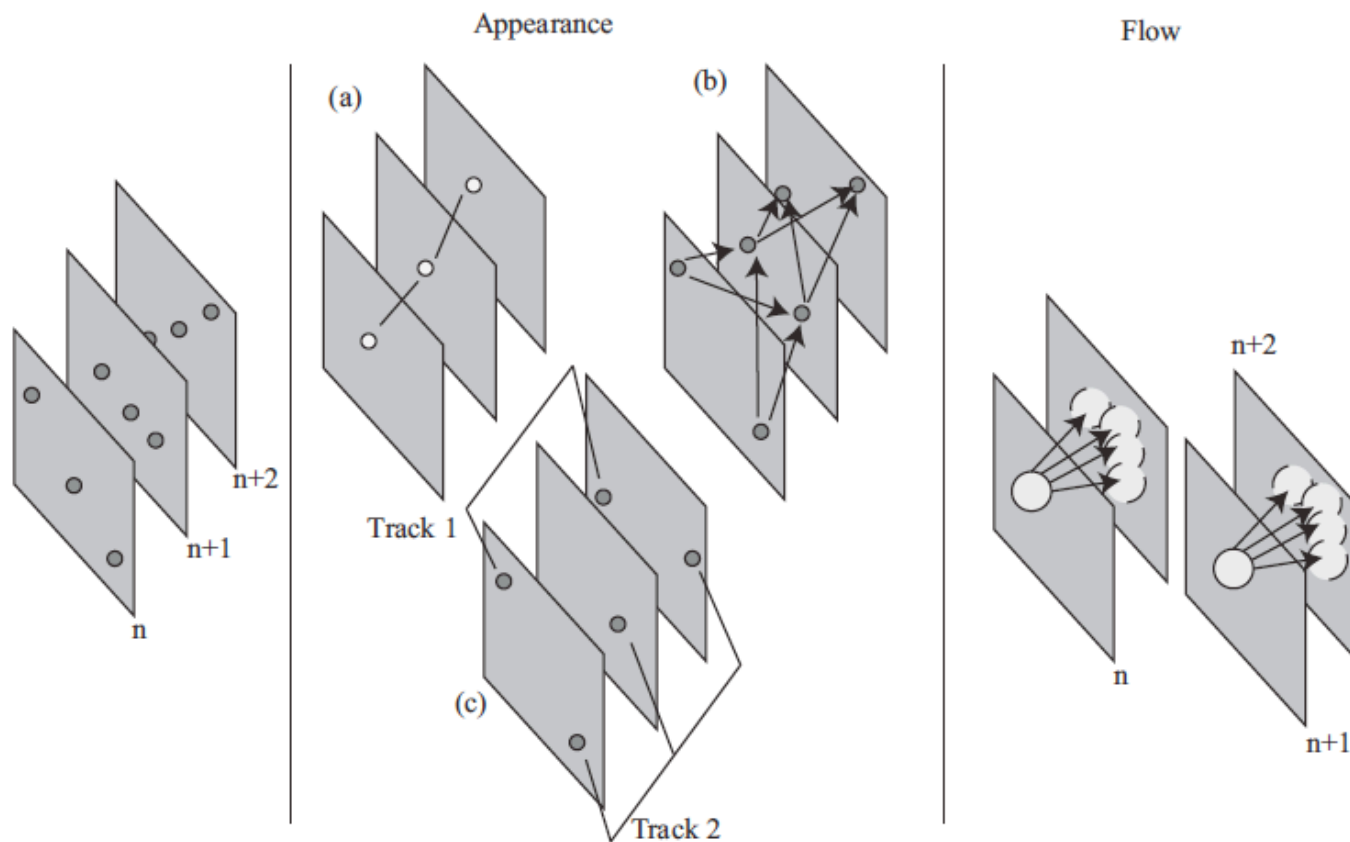    - if we know how an object moves, that selects the pixels to use

Appearance      Flow

(a) (b) (c)

Track 1

Track 2

n, n+1, n+2

**FIGURE 11.1:** In tracking problems, we want to build space time paths followed by tokens—which might be objects, or regions, or interest points, or image windows—in an image sequence (**left**). There are two important sources of information; carefully used, they can resolve many tracking problems without further complexity. One is the appearance of the token being tracked. If there is only one token in each frame with a distinctive appearance, then we could detect it in each frame, then link the detector responses (**a**). Alternatively, if there is more than one instance per frame, a cost function together with weighted bipartite matching could be enough to build the track (**b**). If some instances drop out, we will need to link detector responses to abstract tracks (**c**); in the figure, track 1 has measurements for frames $n$ and $n + 2$, but does not have a measurement for frame $n + 1$. Another important source of information is the motion of the token; if we have a manageable model of the flow, we could search for the flow that generates the best match in the next frame. We choose that match as the next location of the token, then iterate this procedure (**right**).

# Track by detection (simple form)

- Assume
  - a very reliable detector (e.g. faces; back of heads)
  - detections that are well spaced in images (or have distinctive properties)
    - e.g. news anchors; heads in public
- Link detects across time
  - only one - easy
  - multiple - weighted bipartite matching
  - but what if one is missing?
- Better: create abstract tracks
  - link detects to track
  - create tracks, reap tracks as required
  - clean up spacetime paths

**Notation:**
Write $\mathbf{x}_k(i)$ for the $k$'th response of the detector in the $i$th frame
Write $t(k, i)$ for the $k$'th track in the $i$th frame
Write $*t(k, i)$ for the detector response attached to the $k$'th track in the $i$th frame
(Think C pointer notation)

**Assumptions:** We have a detector which is reasonably reliable.
We know some distance $d$ such that $d(*t(k, i-1), *t(k, i))$ is always small.

**First frame:** Create a track for each detector response.

**N'th frame:**
**Link** tracks and detector responses by solving a bipartite matching problem.
**Spawn** a new track for each detector response not allocated to a track.
**Reap** any track that has not received a detector response for some number of frames.

**Cleanup:** We now have trajectories in space time. Link anywhere this is justified (perhaps by a more sophisticated dynamical or appearance model, derived from the candidates for linking).

**Algorithm 11.1:** Tracking by Detection.

# Tracking by known appearance

- Even if we don't have a detector
- Know rectangle in image (n)
- Want to find corresponding rectangle in (n+1)
- Search over nearby rectangles
  - to find one that minimizes SSD error (Sum of Squared Differences)

$$\sum_{i,j}(\mathcal{R}_{ij}^{(n)} - \mathcal{R}_{ij}^{(n+1)})^2.$$

  - where sum is over pixels in rectangle

- Application
  - stabilize players in TV sport

FIGURE 11.2: A useful application of tracking is to stabilize an image box around a more interesting structure, in this case a football player in a television-resolution video. A frame from the video is shown on the **left**. Inset is a box around a player, zoomed to a higher resolution. Notice that the limbs of the player span a few pixels, are blurry, and are hard to resolve. A natural feature for inferring what the player is doing can be obtained by stabilizing the box around the player, then measuring the motion of the limbs with respect to the box. Players move relatively short distances between frames, and their body configuration changes a relatively small amount. This means the new box can be found by searching all nearby boxes of the same size to get the box whose pixels best match those of the original. On the **right**, a set of stabilized boxes; the strategy is enough to center the player in a box. *This figure was originally published as Figure 7 of "Recognizing Action at a Distance," A. Efros, A.C. Berg, G. Mori, and J. Malik, Proc. IEEE ICCV, 2003, © IEEE, 2003.*

# But what if the patch deforms?

- Eg a football player's jersey
    - Colors are "similar" but SSD won't work
- Idea:  patch histogram is stable

- To track:
    - repeat
        - predict location of new patch
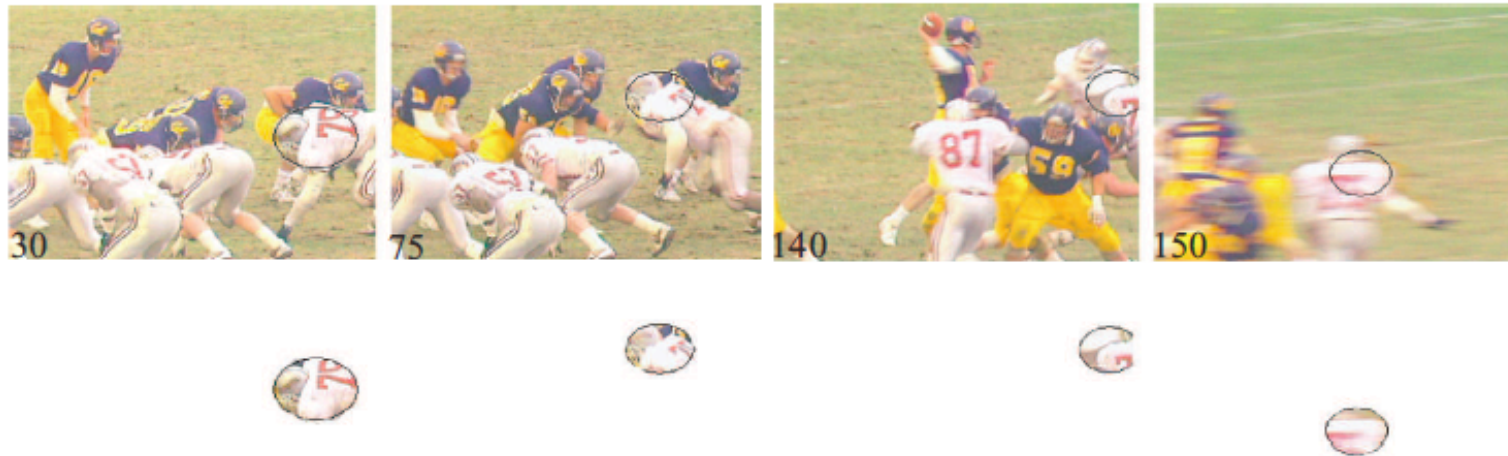        - search nearby for patch whose histogram matches original the best

**FIGURE 11.5**: Four frames from a sequence depicting football players, with superimposed domains. The object to be tracked is the blob on top of player 78 (at the center right in frame 30). We have masked off these blobs (**below**) to emphasize just how strongly the pixels move around in the domain. Notice the motion blur in the final frame. These blobs can be matched to one another, and this is done by comparing histograms (in this case, color histograms), which are less affected by deformation than individual pixel values. *This figure was originally published as Figure 1 of "Kernel-Based Object Tracking" by D. Comaniciu, V. Ramesh, and P. Meer, IEEE Transactions on Pattern Analysis and Machine Intelligence, 2003, © IEEE 2003.*

# Q: Why no deep network methods?

- Silly A: I wrote these slides in 2011

- Better A: They've changed some feature mechanics
  - but nothing major
  - I'll work through some examples later

# When are motions "easy"?

- Current procedure
  - predict state
  - obtaining measurement from prediction by search
  - correct state
- Easy
  - When the object is close to where you expect it to be
    - eg  Object guaranteed to move a little
- Large motions can be easy
  - When they're "predictable"
    - e.g. ballistic motion
    - e.g. constant velocity
- Need a theory to fuse this procedure with motion model

# Q: why care about dynamics?

- Surely deep network methods make this go away?

- A: No they didn't
  - Some problems are "easy" with dynamical models, hard without
  - Example:
    - many similar, fast-moving cars