

Motion and optical flow

Tues Feb 5, 2018
Kristen Grauman
UT Austin

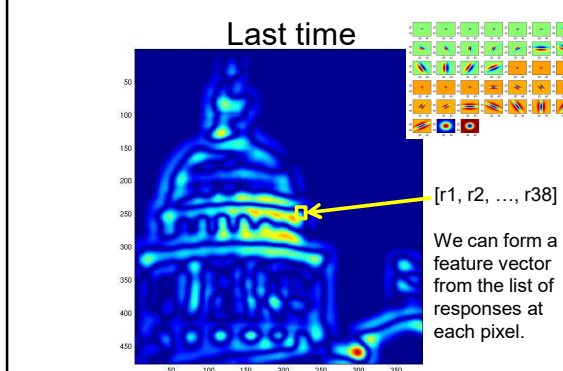
Announcements

- A1 due this Friday

Last time

- Texture is a useful property that is often indicative of materials, appearance cues
- **Texture representations** attempt to summarize repeating patterns of local structure
- **Filter banks** useful to measure redundant variety of structures in local neighborhood
 - Feature spaces can be multi-dimensional
- Neighborhood statistics can be exploited to “sample” or **synthesize** new texture regions
 - Example-based technique

Last time




We can form a feature vector from the list of responses at each pixel.

$$D(a,b) = \sqrt{\sum_{i=1}^d (a_i - b_i)^2} \text{ Euclidean distance (L}_2\text{)}$$

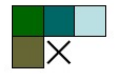
Texture synthesis: intuition

Before, we inserted the next word based on existing nearby words...

Now we want to insert **pixel intensities** based on existing nearby pixel values.



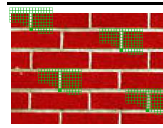
Sample of the texture
("corpus")



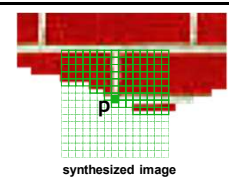
Place we want to insert next

Distribution of a value of a pixel is conditioned on its neighbors alone.

Synthesizing One Pixel



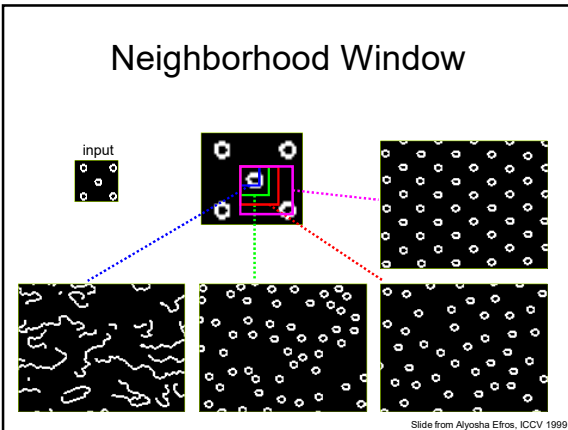
input image



synthesized image

- What is $P(x|\text{neighborhood of pixels around } x)$?
- Find all the windows in the image that match the neighborhood
- To synthesize x
 - pick one matching window at random
 - assign x to be the center pixel of that window
- An **exact** neighbourhood match might not be present, so find the **best** matches using **SSD error** and randomly choose between them, preferring better matches with higher probability

Slide from Ajaysha Efros, ICCV 1999



How to evaluate texture generation?

Quantitative Evaluation of Near Regular Texture Synthesis Algorithms

Wen-Chieh Lin^{*} James Hays^{**} Chenyu Wu^{**} Vivek Kwatra^{***} Yanxi Liu^{**}

^{*} College of Computer Science, National Chiao-Tung University
^{**} School of Computer Science, Carnegie Mellon University
^{***} Department of Computer Science, University of North Carolina at Chapel Hill

Figure 2. A near-regular texture overlaid with its lattice (left) and its geometrically regular counterpart (right), where L and L_r are the underlying lattices, and \vec{r}_1 and \vec{r}_2 are the generating vectors of the regular lattice L_r .

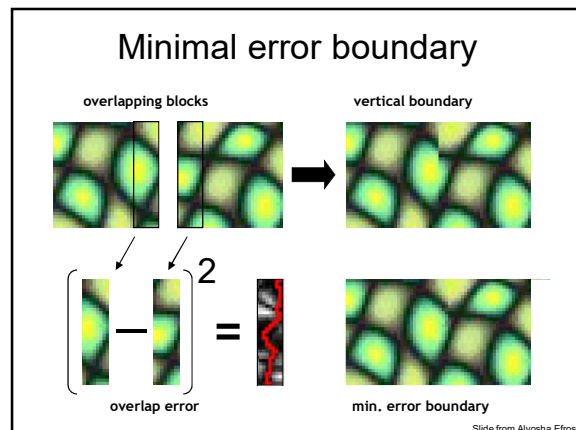
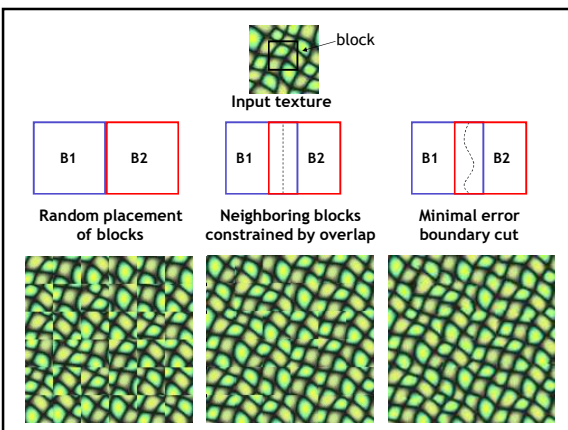
Figure 3. A sample set of near-regular textures used in this study

- The Efros & Leung algorithm
 - Simple
 - Surprisingly good results
 - Synthesis is easier than analysis!
 - ...but very slow

Image Quilting [Efros & Freeman 2001]


- **Observation:** neighbor pixels are highly correlated
- Idea:** unit of synthesis = block
 - Exactly the same but now we want $P(B|N(B))$
 - Much faster: synthesize all pixels in a block at once

Slide from Alyosha Efros, ICCV 1999









Texture Transfer

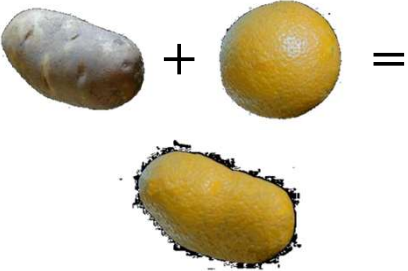
- Take the texture from one object and "paint" it onto another object
 - This requires separating texture and shape
 - That's HARD, but we can cheat
 - Assume we can capture shape by boundary and rough shading
- Then, just add another constraint when sampling: similarity to underlying image at that spot




Slide credit: Freeman & Efros


+

=



+

=



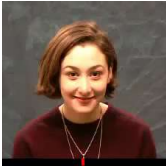


(Manual) texture synthesis in the media





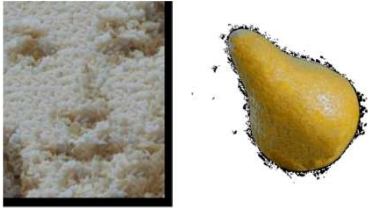
Video textures

Arno Schodl, Richard Szeliski, David H. Salesin, and Irfan Essa. [Video textures](#). *Proceedings of SIGGRAPH 2000*, pages 489-498, July 2000.

Style transfer

<https://www.youtube.com/watch?v=Khuj4ASldmU>

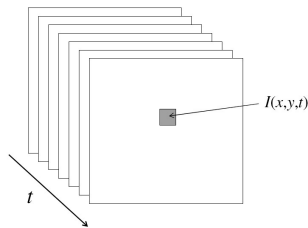


Today

- Optical flow: estimating motion in video
- Background subtraction

Video

- A video is a sequence of frames captured over time
- Now our image data is a function of space (x, y) and time (t)

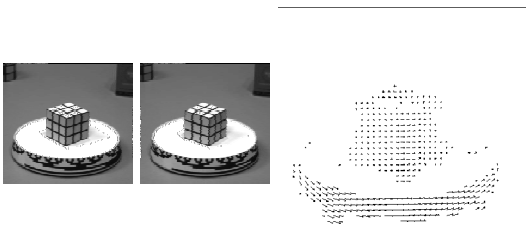


Uses of motion

- Estimating 3D structure
- Segmenting objects based on motion cues
- Learning dynamical models
- Recognizing events and activities
- Improving video quality (motion stabilization)

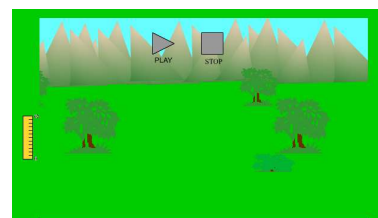
Motion field

- The motion field is the projection of the 3D scene motion into the image



Motion parallax

<http://psych.hanover.edu/KRANTZ/MotionParallax/MotionParallax.html>



Motion field + camera motion

Length of flow vectors inversely proportional to depth Z of 3d point

points closer to the camera move more quickly across the image plane

Figure 1.2. Two images taken from a helicopter flying through a canyon and the computed optical flow field.

Figure from Michael Black, Ph.D. Thesis

Motion field + camera motion

Zoom out Zoom in Pan right to left

Motion estimation techniques

- **Direct methods**
 - Directly recover image motion at each pixel from spatio-temporal image brightness variations
 - Dense motion fields, but sensitive to appearance variations
 - Suitable for video and when image motion is small
- **Feature-based methods**
 - Extract visual features (corners, textured areas) and track them over multiple frames
 - Sparse motion fields, but more robust tracking
 - Suitable when image motion is large (10s of pixels)

Optical flow

- Definition: optical flow is the *apparent* motion of brightness patterns in the image
- Ideally, optical flow would be the same as the motion field
- Have to be careful: apparent motion can be caused by lighting changes without any actual motion

Apparent motion != motion field

Figure 12.2. The optical flow is not always equal to the motion field. In (a) a smooth sphere is rotating under constant illumination—the image does not change, yet the motion field is nonzero. In (b) a fixed sphere is illuminated by a moving source—the shading in the image changes, yet the motion field is zero.

Figure from Horn book

Problem definition: optical flow

$H(x, y)$ $I(x, y)$

How to estimate pixel motion from image H to image I ?

- Solve pixel correspondence problem
 - given a pixel in H , look for nearby pixels of the same color in I

Key assumptions

- **color constancy**: a point in H looks the same in I
 - For grayscale images, this is **brightness constancy**
- **small motion**: points do not move very far

This is called the **optical flow problem**

Color/brightness constancy

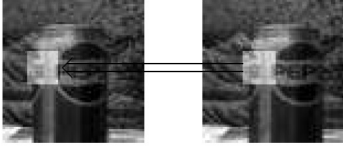
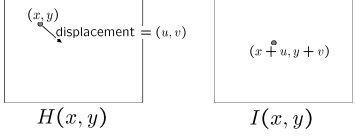


Figure 1.5: Data conservation assumption. The highlighted region in the right image looks roughly the same as the region in the left image, despite the fact that it has moved.

Figure by Michael Black

Optical flow constraints



Let's look at these constraints more closely

- brightness constancy: Q: what's the equation?

$$H(x, y) = I(x + u, y + v)$$
- small motion:

$$I(x + u, y + v) = I(x, y) + \frac{\partial I}{\partial x}u + \frac{\partial I}{\partial y}v + \text{higher order terms}$$

$$\approx I(x, y) + \frac{\partial I}{\partial x}u + \frac{\partial I}{\partial y}v$$

Optical flow equation

shorthand: $I_x = \frac{\partial I}{\partial x}$

Combining these two equations

$$0 = I(x + u, y + v) - H(x, y)$$

$$\approx I(x, y) + I_x u + I_y v - H(x, y)$$

$$\approx (I(x, y) - H(x, y)) + I_x u + I_y v$$

$$\approx I_t + I_x u + I_y v$$

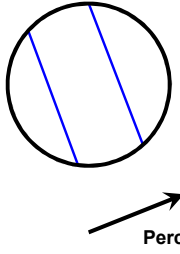
$$\approx I_t + \nabla I \cdot [u \ v]$$

Optical flow equation

$$0 = I_t + \nabla I \cdot [u \ v]$$

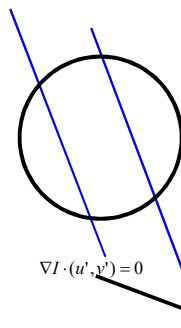
Q: how many unknowns and equations per pixel?

The aperture problem



Perceived motion

The aperture problem



Actual motion

$$\nabla I \cdot (u', y') = 0$$

The barber pole illusion

Barber's pole Motion field Optical flow

http://en.wikipedia.org/wiki/Barberpole_illusion

Solving the aperture problem

- How to get more equations for a pixel?
- Spatial coherence constraint:** pretend the pixel's neighbors have the same (u,v)

Figure 1.7: Spatial coherence assumption. Neighboring points in the image are assumed to belong to the same surface in the scene.

Figure by Michael Black

Solving the aperture problem

- How to get more equations for a pixel?
- Spatial coherence constraint:** pretend the pixel's neighbors have the same (u,v)
 - If we use a 5x5 window, that gives us 25 equations per pixel

$$0 = I_t(p_i) + \nabla I(p_i) \cdot [u \ v]$$

$$\begin{bmatrix} I_x(p_1) & I_y(p_1) \\ I_x(p_2) & I_y(p_2) \\ \vdots & \vdots \\ I_x(p_{25}) & I_y(p_{25}) \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix} = - \begin{bmatrix} I_t(p_1) \\ I_t(p_2) \\ \vdots \\ I_t(p_{25}) \end{bmatrix}$$

$$A \ d = b$$

25x2 2x1 25x1

Slide credit: Steve Seitz

Solving the aperture problem

Prob: we have more equations than unknowns

$$A \ d = b \quad \xrightarrow{\text{minimize } \|Ad - b\|^2}$$

25x2 2x1 25x1

Solution: solve least squares problem

- minimum least squares solution given by solution (in d) of:

$$(A^T A) \ d = A^T b$$

2x2 2x1 2x1

$$\begin{bmatrix} \sum I_x I_x & \sum I_x I_y \\ \sum I_x I_y & \sum I_y I_y \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix} = - \begin{bmatrix} \sum I_x I_t \\ \sum I_y I_t \end{bmatrix}$$

$A^T A$ $A^T b$

- The summations are over all pixels in the K x K window
- This technique was first proposed by Lucas & Kanade (1981)

Slide credit: Steve Seitz

Conditions for solvability

$$\begin{bmatrix} \sum I_x I_x & \sum I_x I_y \\ \sum I_x I_y & \sum I_y I_y \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix} = - \begin{bmatrix} \sum I_x I_t \\ \sum I_y I_t \end{bmatrix}$$

$A^T A$ $A^T b$

When is this solvable?


- $A^T A$ should be invertible
- $A^T A$ should not be very small
 - eigenvalues λ_1 and λ_2 of $A^T A$ should not be very small
- $A^T A$ should be well-conditioned
 - λ_1 / λ_2 should not be too large (λ_1 = larger eigenvalue)

Slide by Steve Seitz, UW

Edge


- gradients very large or very small
- large λ_1 , small λ_2

Low-texture region




- gradients have small magnitude
- small λ_1 , small λ_2

High-texture region



- gradients are different, large magnitudes
- large λ_1 , large λ_2

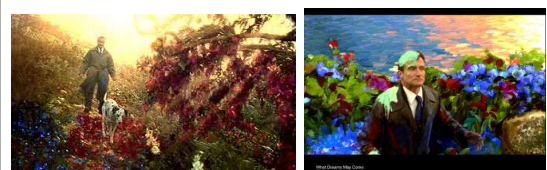
Example use of optical flow: facial animation



<http://www.fxguide.com/article333.html>

Example use of optical flow: Motion Paint

Use optical flow to track brush strokes, in order to animate them to follow underlying scene motion.

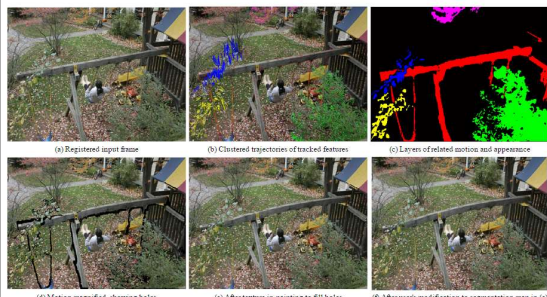


<http://www.fxguide.com/article333.html>

Motion estimation techniques

- Direct methods
 - Directly recover image motion at each pixel from spatio-temporal image brightness variations
 - Dense motion fields, but sensitive to appearance variations
 - Suitable for video and when image motion is small
- Feature-based methods
 - Extract visual features (corners, textured areas) and track them over multiple frames
 - Sparse motion fields, but more robust tracking
 - Suitable when image motion is large (10s of pixels)

Motion magnification



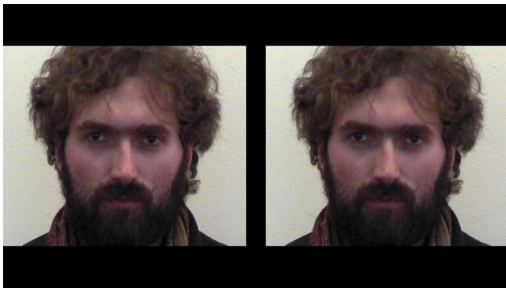
(a) Motion magnified, showing holes (b) After texture re-painting to fill holes (c) After user's modifications to segmentation map see (c)

Liu et al. SIGGRAPH 2005

original



magnified



Motion magnification

- <http://people.csail.mit.edu/mrub/vidmag/>
- Hao-Yu Wu, Michael Rubinstein, Eugene Shih, John Guttag, Frédo Durand, William T. Freeman
Eulerian Video Magnification for Revealing Subtle Changes in the World ACM Transactions on Graphics, Volume 31, Number 4 (Proc. SIGGRAPH), 2012
- Ce Liu, Antonio Torralba, William T. Freeman, Frédo Durand, Edward H. Adelson. **Motion Magnification** ACM Transactions on Graphics, Volume 24, Number 3 (Proc. SIGGRAPH), 2005

Today

- Optical flow: estimating motion in video
- Background subtraction



Video as an "Image Stack"

Can look at video data as a spatio-temporal volume

- If camera is stationary, each line through time corresponds to a single ray in space

Alvisha Efros, CMU

Input Video

Alvisha Efros, CMU

Average Image

Alvisha Efros, CMU

Background Subtraction

▶ Given an image (mostly likely to be a video frame), we want to identify the **foreground objects** in that image!

Motivation

- ▶ In most cases, objects are of interest, not the scene.
- ▶ Makes our life easier: less processing costs, and less room for error.

Slide credit: Birgi Tamersov

Background subtraction

- Simple techniques can do ok with static camera
- ...But hard to do perfectly
- Widely used:
 - Traffic monitoring (counting vehicles, detecting & tracking vehicles, pedestrians),
 - Human action recognition (run, walk, jump, squat),
 - Human-computer interaction
 - Object tracking

Simple Approach

Image at time t : $I(x, y, t)$ Background at time t : $B(x, y, t)$

1. Estimate the background for time t .
2. Subtract the estimated background from the input frame.
3. Apply a threshold, Th , to the absolute difference to get the **foreground mask**.

Slide credit: Birgi Tamersov


Frame Differencing

- Background is estimated to be the previous frame. Background subtraction equation then becomes:

$$B(x, y, t) = I(x, y, t - 1)$$

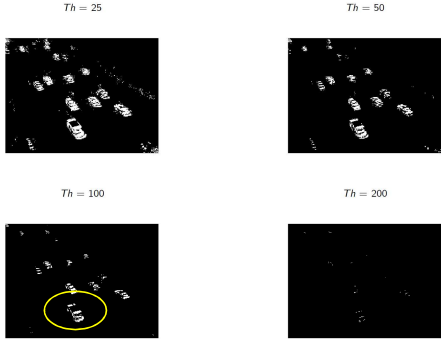
$$\downarrow$$

$$|I(x, y, t) - I(x, y, t - 1)| > Th$$
- Depending on the object structure, speed, frame rate and global threshold, this approach may or may not be useful (usually **not**).



Slide credit: Birgi Tamersov

Frame Differencing



Slide credit: Birgi Tamersov

Mean Filter


- In this case the background is the mean of the previous n frames:

$$B(x, y, t) = \frac{1}{n} \sum_{i=0}^{n-1} I(x, y, t - i)$$


$$\downarrow$$

$$|I(x, y, t) - \frac{1}{n} \sum_{i=0}^{n-1} I(x, y, t - i)| > Th$$
- For $n = 10$:

Estimated Background



Foreground Mask



Slide credit: Birgi Tamersov

Median Filter


- Assuming that the background is more likely to appear in a scene, we can use the median of the previous n frames as the background model:

$$B(x, y, t) = \text{median}\{I(x, y, t - i)\}$$


$$\downarrow$$

$$|I(x, y, t) - \text{median}\{I(x, y, t - i)\}| > Th \text{ where } i \in \{0, \dots, n - 1\}.$$
- For $n = 10$:

Estimated Background





Foreground Mask





Slide credit: Birgi Tamersov

Average/Median Image





Alyosha Efros, CMU

Background Subtraction

=



Alyosha Efros, CMU

Pros and cons

Advantages:

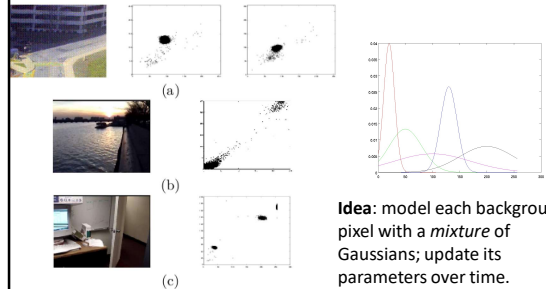
- Extremely easy to implement and use!
- All pretty fast.
- Corresponding background models need not be constant, they change over time.

Disadvantages:

- Accuracy of frame differencing depends on object speed and frame rate
- Median background model: relatively high memory requirements.
- Setting global threshold Th...

When will this basic approach fail?

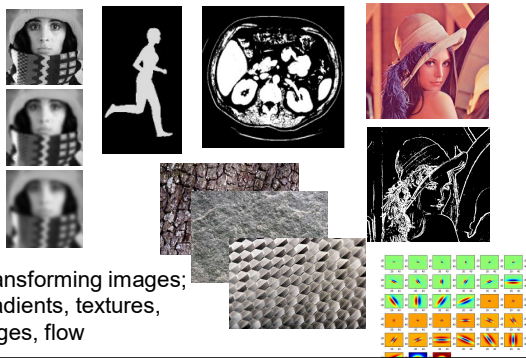
Background mixture models



Idea: model each background pixel with a *mixture* of Gaussians; update its parameters over time.

Adaptive Background Mixture Models for Real-Time Tracking, Chris Stauffer & W.E.L. Grimson

So far: features and filters



Transforming images; gradients, textures, edges, flow