# Introduction and Feature Detection

CS448V — Computational Video Manipulation

April 2019

---



**Raiders of the Lost Ark: The Adaptation** [Zala 82-89]
Shot-for-shot remake by three 12-year olds (took 7 years)

**How can we let people easily make such video?**

# YouTube

Search | Browse | Movies | Upload | Create Account | Sign In

**Join the largest worldwide video-sharing community!**

**Create Account ›** Already have an account? **Sign In**

## Music

DRAKE HEADLINES 3:57
Drake - Headlines (Take Care 201... by WeAreYoungMoneyEntTM

5:46
Kids React To Rebecca Black - My... by TheFineBros

1:52
2NE1 - Hate You : ComeBack Stage by CrazyCarrot370

UKF BASS CULTURE 10:00
UKF Bass Culture (Dubstep Megamix) by UKFDubstep

## Entertainment

2:42
Wingman by kevjumba

YOGSCAST TEACHES ATHENE MINECRAFT 13:58
Minecraft - Yogscast Teaches Ath... by BlueXephos

0:59
X Games 17: Moto X Enduro Women'... by XGames

MINECRAFT Mountain of Kikatchu #1 17:00
Minecraft - Mountain of Kikatchu... by BlueXephos

## Sports

0:29

0:14

2:35

0:31

Show Ad

Spotlight

**Music Tuesday: My Morning Jacket**

It's day two of Lollapalooza Week here on YouTube. My Morning Jacket will be performing live in Chicago and on YouTube this Saturday. To whet your appetite enjoy their new video, plus a playlist of their favorites.
Presented by: lollapalooza

4:26
My Morning Jacket "Holdin On To Black Metal" by MMJofficial 3,557 views

0:42
My Morning Jacket: Exclusive Video Playlist by MMJofficial 8,867 views

1:37
"Time Piece" by HensonCompany 40,655 views

9:05
Marvin Gaye "What's Going On / What's Happening... by brainchild9 8,612,838 views

Trends

2:20
President Obama's Message on the Debt Agreement by BarackObamadotcom 303,836 views

0:14
Kobe Bryant scores own goal at Mia Hamm charity... by jmoynihanpatch 188,695 views

---

# Snap

Tap to take a Snap, then send it to a friend!

# Memories

A personal collection of your favorite Snaps and Stories

Memories
ALL SNAPS STORIES CAMERA ROLL MY EYES O

# Challenge

**People want to create and share stories**

26% of all Internet users post original videos [Pew 13]

3,500,000,000 snaps/day uploaded to Snap [The Verge 17]

300 hours video/minute uploaded to YouTube [Youtube FAQ 18]

**But raw video rarely tells a compelling story**

Content not well thought out

Poor composition, lighting, etc.

Often too long

**Best stories are planned, edited and produced**

Current tools force users to work with low-level controls

**Need higher-level tools for manipulating video**

# Course Goals

1. *Gain overview of* algorithmic techniques used to manipulate video
2. *Present research paper* and *lead discussion* on a research paper
3. *Capture and edit* video manually and using algorithmic techniques
4. *Develop* substantial video manipulation project

# Instructor: Maneesh Agrawala



**Visual Rhythm and Beat.** Abe Davis and Maneesh Agrawala, SIGGRAPH 2018.

# Instructor: Ohad Fried



Text-based Editing of Talking-head Video

**Text-Based Editing of Talking Head Video.** Ohad Fried, Ayush Tewari, Michael Zollhoefer, Adam Finkelstein, Eli Shectman, Dan B Goldman, Kyle Genova, Zeyu Jin, Christian Theobolt and Maneesh Agrawala, SIGGRAPH 2019.

# Instructor: Michael Zollhöfer



**Deep Video Portraits**  H. Kim, P. Garrido, A. Tewari, W. Xu, J. Thies, M. Nießner, P. Perez, C. Richardt, M. Zollhöfer, C. Theobalt SIGGRAPH 2018

# Course Mechanics

# Readings, Discussions, Presentations

**Required to read about one paper per class**

We will provide prompts to guide reading

You are responsible for written response to prompt

*Due on paper at beginning of class, 2 free passes for the quarter*

**Required to present a paper and lead discussion once in the quarter**

Usually Mon will be student presentations

You will meet with us (instructors) in week before presentation to go over 1st draft

---

# Website

https://magrawala.github.io/cs448v-sp19/

# Requirements

**Participation (15%)**

Attendance with prompt response is mandatory (but 2 free passes)
Also must engage in discussion in class

**Presentation (15%)**

Deeply engage with at least one paper and help others understand it

**Assignments (20%)**

Will help you learn about manual editing and the programmatic toolkits (e.g. OpenCV)
available to implement algorithms

**Final Project (50%)**

Implement a research project on video manipulation

# A1: Manual Manipulation

**Interview a classmate and capture on video for at least 15 minutes**

Plan the interview questions ahead of time
Capture on video (at least 15 minutes) – Do **not** hold camera, use a stand

**Edit raw footage into a short video (< 2min) you would be proud to share**

Use any video editing software you wish (e.g. Premiere, FinalCut Pro, iMovie)

**Write down your reflections (half page PDF)**

What was difficult in capturing and especially editing?
List all the pain points.
Describe ways video editing could be improved

**Due Wed Apr 10 at 1:30pm**

# Feature Detection

# Image Matching



by Diva Sian



by scgbt

Slide credit: Seitz

# Local Measures of Distinctiveness

**Suppose we only consider a small window of pixels**
What defines whether a feature is a good or bad candidate?

# Feature Detection

**Local measure of feature uniqueness**
- How does the window change when you shift it?
- Shifting the window in *any direction* causes a *big change*

"flat" region:
no change in all
directions

"edge":
no change along the
edge direction

"corner":
significant change in
all directions

# Feature Detection: Math

**Consider shifting the window W by (u,v)**

- How do the pixels in **W** change?
- Compare each pixel before and after by summing up the squared differences (SSD)
- This defines an SSD "error" of *E(u,v)*:

**W**

$$E(u, v) = \sum_{(x,y) \in W} [I(x + u, y + v) - I(x, y)]^2$$

# Small Motion Assumption

Taylor Series expansion of I:

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

If the motion (u,v) is small, then first order approx is good

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

$$\approx I(x, y) + [I_x \ I_y] \begin{bmatrix} u \\ v \end{bmatrix}$$

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

Plugging this into the formula on the previous slide…

# Feature Detection: Math

**Consider shifting the window W by (u,v)**

- How do the pixels in **W** change?
- Compare each pixel before and after by summing up the squared differences (SSD)
- This defines an SSD "error" of *E(u,v)*:

$$E(u,v) = \sum_{(x,y)\in W} [I(x+u, y+v) - I(x,y)]^2$$

$$\approx \sum_{(x,y)\in W} \left[ I(x,y) + [I_x \ I_y] \begin{bmatrix} u \\ v \end{bmatrix} - I(x,y) \right]^2$$

$$\approx \sum_{(x,y)\in W} \left[ [I_x \ I_y] \begin{bmatrix} u \\ v \end{bmatrix} \right]^2$$

# Feature Detection:  Math

**This can be rewritten:**

$$E(u,v) = \sum_{(x,y)\in W} [u \ v] \underbrace{\begin{bmatrix} I_x^2 & I_x I_y \\ I_y I_x & I_y^2 \end{bmatrix}}_{H} \begin{bmatrix} u \\ v \end{bmatrix}$$

**Suppose you can move the center of the blue window in any direction**

- Which directions will result in the largest and smallest E values?
- We can find these directions by looking at the eigenvectors of **H**

# Eigenvalues & Eigenvectors

The **eigenvectors** of a matrix **A** are the vectors **x** that satisfy:

$$Ax = \lambda x$$

The scalar $\lambda$ is the **eigenvalue** corresponding to **x**
   The eigenvalues are found by solving:

$$det(A - \lambda I) = 0$$

- In our case, **A** = **H** is a 2x2 matrix, so we have

$$det \begin{bmatrix} h_{11} - \lambda & h_{12} \\ h_{21} & h_{22} - \lambda \end{bmatrix} = 0$$

- The solution:

$$\lambda_{\pm} = \tfrac{1}{2}\left[ (h_{11} + h_{22}) \pm \sqrt{4h_{12}h_{21} + (h_{11} - h_{22})^2} \right]$$

Once you know $\lambda$, you find **x** by solving

$$\begin{bmatrix} h_{11} - \lambda & h_{12} \\ h_{21} & h_{22} - \lambda \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = 0$$

# Feature Detection:  Math

**This can be rewritten:**

$$E(u,v) = \sum_{(x,y) \in W} [u\ v] \underbrace{\begin{bmatrix} I_x^2 & I_x I_y \\ I_y I_x & I_y^2 \end{bmatrix}}_{H} \begin{bmatrix} u \\ v \end{bmatrix}$$



**Eigenvalues and eigenvectors of H**

- Define shifts with the smallest and largest change (E value)
- $x_+$ = direction of **largest** increase in E.
- $\lambda_+$ = amount of increase in direction $x_+$
- $x_-$ = direction of **smallest** increase in E.
- $\lambda$- = amount of increase in direction $x_+$

$$Hx_+ = \lambda_+ x_+$$
$$Hx_- = \lambda_- x_-$$

# Feature Detection:  Math

How are $\lambda_+$, $\mathbf{x}_+$, $\lambda_-$, and $\mathbf{x}_+$ relevant for feature detection?

• What's our feature scoring function?

# Feature detection:  the math

How are $\lambda_+$, $\mathbf{x}_+$, $\lambda_-$, and $\mathbf{x}_+$ relevant for feature detection?

• What's our feature scoring function?

Want *E(u,v)* to be ***large*** for small shifts in ***all*** directions

• the *minimum* of *E(u,v)* should be large, over all unit vectors [u v]
• this minimum is given by the smaller eigenvalue ($\lambda_-$) of ***H***



$$I \qquad\qquad \lambda_+ \qquad\qquad \lambda_-$$

# Feature Detection Summary

**Here's what you do**

- Compute the gradient at each point in the image
- Create the **H** matrix from the entries in the gradient
- Compute the eigenvalues.
- Find points with large response ($\lambda_- >$ threshold)
- Choose those points where $\lambda_-$ is a local maximum as features



$$I \qquad \lambda_+ \qquad \lambda_-$$

Slide credit: Seitz, Frovola, Simakov

# Feature Detection Summary

**Here's what you do**

- Compute the gradient at each point in the image
- Create the **H** matrix from the entries in the gradient
- Compute the eigenvalues.
- Find points with large response ($\lambda_- >$ threshold)
- Choose those points where $\lambda_-$ is a local maximum as features



$$\lambda_-$$

Slide credit: Seitz, Frovola, Simakov

# The Harris Operator

$\lambda_-$ is a variant of the "Harris operator" for feature detection

$$f = \frac{\lambda_1 \lambda_2}{\lambda_1 + \lambda_2}$$

$$= \frac{determinant(H)}{trace(H)}$$

- The *trace* is the sum of the diagonals, i.e., *trace(H) = h₁₁ + h₂₂*
- Very similar to $\lambda_-$ but less expensive (no square root)
- Called the "Harris Corner Detector" or "Harris Operator"
- Lots of other detectors, this is one of the most popular

Slide credit: Seitz, Frovola, Simakov

# The Harris Operator



Harris operator

$\lambda_-$

Slide credit: Seitz, Frovola, Simakov

# Harris Operator Example

# f value (red high, blue low)

# Threshold (f > value)



Slide credit: Seitz, Frovola, Simakov

# Find Local Maxima of f



Slide credit: Seitz, Frovola, Simakov

# Harris Features (in red)



Slide credit: Seitz, Frovola, Simakov

# Invariance with Harris Corners

- Translation invariance
- Rotation invariance
- Scale invariance?



Corner

All points will be classified as **edges**!

**Not** invariant to image scale!

Slide credit: Kristen Grauman

# Scale Invariant Detection

Consider regions (e.g. circles) of different sizes around a point

Find regions of corresponding sizes that will look the same in both images?



# Scale Invariant Detection

The problem: how do we choose corresponding circles *independently* in each image?

# Difference of Gaussians



$$G(k^2\sigma)*I$$
$$G(k\sigma)*I$$
$$G(\sigma)*I$$

$$D(\sigma)\equiv(G(k\sigma)-G(\sigma))*I$$

Gaussian

Difference of Gaussian (DOG)

Scale (first octave)

# Scale-Space Extrema

**Choose all extrema within 3x3x3 neighborhood**



$$D(k^2\sigma)$$
$$D(k\sigma)$$
$$D(\sigma)$$

Scale

X is selected if it is larger or smaller than all 26 neighbors

# Invariant Local Features

**Find features that are invariant to transformations**
- geometric invariance: translation, rotation, scale
- photometric invariance: brightness, exposure, …



**Feature Descriptors**

# Becoming Rotation Invariant

- We are given a keypoint and its scale from DoG

- We select a characteristic orientation for the keypoint (based on the most prominent gradient in local region)

- We describe all features **relative** to this orientation

- Causes features to be rotation invariant!

  - If the keypoint appears rotated in another image, the features will be the same, because they're **relative** to the characteristic orientation

# SIFT Descriptor Formation



**Use the blurred image associated with the keypoint's scale**

**Take image gradients over the keypoint 16x16 neighborhood (put in 36 bin histogram)**
- Treat max bin as keypoint orientation θ

**To become rotation invariant, rotate the gradient directions AND locations by (- θ)**
- Now we've cancelled out rotation and have gradients expressed at locations **relative** to keypoint orientation θ
- We could also have just rotated the whole image by -θ, but that would be slower

Slide credit: Niebles and Krishna

# SIFT Descriptor Formation



**Using precise gradients & locations is fragile**

**For robustness create array of orientation histograms**

**Put the rotated gradients into their local orientation histograms**
- A gradient's contribution is divided among the nearby histograms based on distance. If it's halfway between two histogram locations, it gives a half contribution to both.
- Also, scale down gradient contributions for gradients far from the center

**The SIFT authors found that best results were with 8 orientation bins per histogram**

Slide credit: Niebles and Krishna

# SIFT Descriptor Formation



Image gradients

Keypoint descriptor

**Using precise gradients & locations is fragile**

**For robustness create array of orientation histograms**

**Put the rotated gradients into their local orientation histograms**

- A gradient's contribution is divided among the nearby histograms based on distance. If it's halfway between two histogram locations, it gives a half contribution to both.
- Also, scale down gradient contributions for gradients far from the center

**The SIFT authors found that best results were with 8 orientation bins per histogram and 4x4 histogram array**

# SIFT Descriptor Formation



Image gradients

Keypoint descriptor

**8 orientation bins per histogram, and  4x4 histogram array: 8 x 4x4 = 128 numbers**

**So a SIFT descriptor is a length 128 vector, which is invariant to rotation (because we rotated the descriptor) and scale (because we worked with the scaled image from DoG)**

**We can compare each vector from image A to each vector from image B to find matching keypoints!**

Euclidean "distance" between descriptor vectors gives a good measure of keypoint similarity

# SIFT Descriptor Formation



Image gradients → Keypoint descriptor

**Adding robustness to illumination changes:**

- Descriptor is made of gradients (differences between pixels), so it's already invariant to changes in brightness (e.g. adding 10 to all image pixels yields the exact same descriptor)

- A higher-contrast photo will increase the magnitude of gradients linearly. So, to correct for contrast changes, normalize the vector (scale to length 1.0)

- Very large image gradients are usually from unreliable 3D illumination effects (glare, etc). So, to reduce their effect, clamp all values in the vector to be ≤ 0.2 (an experimentally tuned value). Then normalize the vector again.

Slide credit: Niebles and Krishna

# SIFT Keypoints Detection

**Threshold on value at DOG peak and on ratio of principle curvatures**



**(a)** 233x189 image
**(b)** 832 DOG extrema
**(c)** 729 left after peak value threshold
**(d)** 536 left after testing ratio of principle curvatures

**Vectors indicate scale, orientation and location**

Slide credit: Niebles and Krishna

# Properties of SIFT

**Extraordinarily robust matching technique**
- Can handle changes in viewpoint  (Up to about 60 degree out of plane rotation)
- Can handle significant changes in illumination (Sometimes even day vs. night (see below))
- Fast and efficient—can run in real time
- Lots of code available
    - http://people.csail.mit.edu/albert/ladypack/wiki/index.php/Known_implementations_of_SIFT



Slide credit: Seitz

# Feature Matching

Given a feature in I$_1$, how to find the best match in I$_2$?

1. Define distance function that compares two descriptors
    (e.g. Euclidean distance between SIFT descriptors)

2. Test all the features in I$_2$, find the one with min distance

Slide credit: Seitz

# Feature Distance

How to define the difference between two features $f_1$, $f_2$?

- Simple approach is SSD($f_1$, $f_2$)
  - sum of square differences between entries of the two descriptors
  - can give good scores to very ambiguous (bad) matches



$I_1$    $I_2$

Slide credit: Seitz

# Feature Distance

How to define the difference between two features $f_1$, $f_2$?

- Better approach:  ratio distance = SSD($f_1$, $f_2$) / SSD($f_1$, $f_2$')
  - $f_2$ is best SSD match to $f_1$ in $I_2$
  - $f_2$' is  2nd best SSD match to $f_1$ in $I_2$
  - gives small values for ambiguous matches



$I_1$    $I_2$

Slide credit: Seitz

# Evaluating the Results

**How can we measure the performance of a feature matcher?**



50
75
200

feature distance

# True/False Positives



50
75
**True match**
200
**False match**

feature distance

**The distance threshold affects performance**

- True positives = # of detected matches that are correct
  - Suppose we want to maximize these—how to choose threshold?
- False positives = # of detected matches that are incorrect
  - Suppose we want to minimize these—how to choose threshold?

# Evaluating the Results



$$\frac{\text{\# true positives}}{\text{\# matching features (positives)}}$$ *true positive rate*

$$\frac{\text{\# false positives}}{\text{\# unmatched features (negatives)}}$$ *false positive rate*

---

# Evaluating the Results

**ROC curve** ("Receiver Operator Characteristic")



$$\frac{\text{\# true positives}}{\text{\# matching features (positives)}}$$ *true positive rate*

$$\frac{\text{\# false positives}}{\text{\# unmatched features (negatives)}}$$ *false positive rate*

**ROC Curves**

• Generated by counting # current/incorrect matches, for different thresholds
• Want to maximize area under the curve (AUC)
• Useful for comparing different feature matching methods
• For more info: http://en.wikipedia.org/wiki/Receiver_operating_characteristic

# Application: Mosaicing



http://www.cs.ubc.ca/~mbrown/autostitch/autostitch.html

# Application: Wide Baseline Stereo



[Image from T. Tuytelaars ECCV 2006 tutorial]

# Application: Object/Scene Recognition



Schmid and Mohr 1997



Sivic and Zisserman, 2003



Rothganger et al. 2003



Lowe 2002