



University of Cyprus – MSc Artificial Intelligence

# MAI644 – COMPUTER VISION

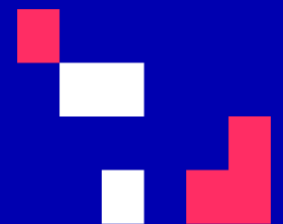
## Lecture 1: Introduction to Computer Vision

**Melinos Averkiou**

CYENS Centre of Excellence

University of Cyprus - Department of Computer Science

[m.averkiou@cyens.org.cy](mailto:m.averkiou@cyens.org.cy)





# Today's Agenda

- Who we are
- Introduction to Computer Vision
  - What is Computer Vision
  - How hard is Computer Vision
  - Why is Computer Vision so hard
  - How to organize Computer Vision
  - Why study Computer Vision
  - Applications
- What we do



# Who we are

## Visual Computing Group at CYENS Centre of Excellence



**Melinos Averkiou**  
MRG Leader



**Yiangos Georgiou**  
Research Associate  
(DTP)



**Marios Loizou**  
Research Associate



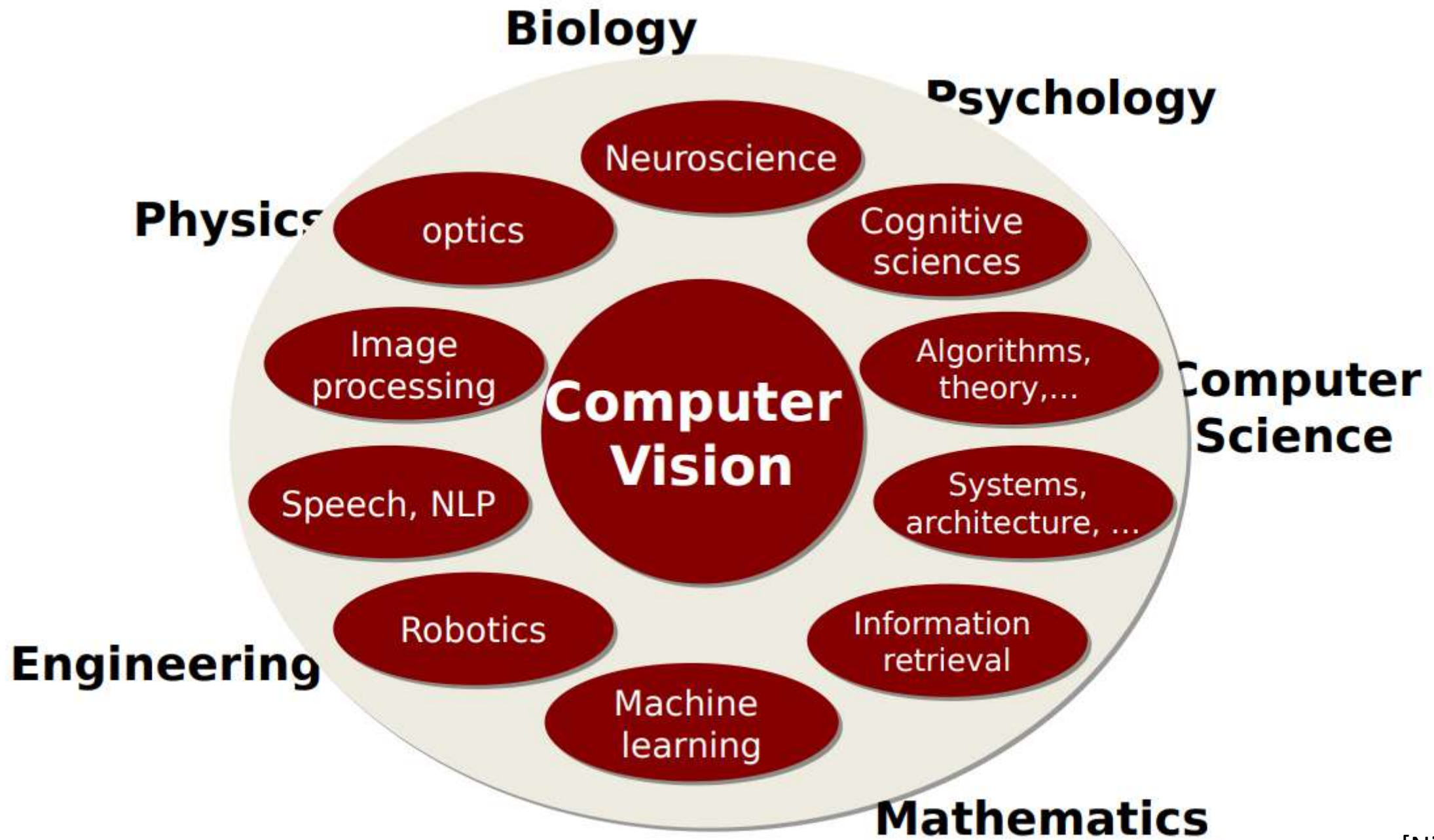
**Yeshwanth Kumar  
Adimoolam**  
Research Associate  
(DTP)

# Today's Agenda

- Course Overview
- Introduction to Computer Vision
  - What is Computer Vision
  - How hard is Computer Vision
  - Why is Computer Vision so hard
  - How to organize Computer Vision
  - Why study Computer Vision
  - Applications
- What we do



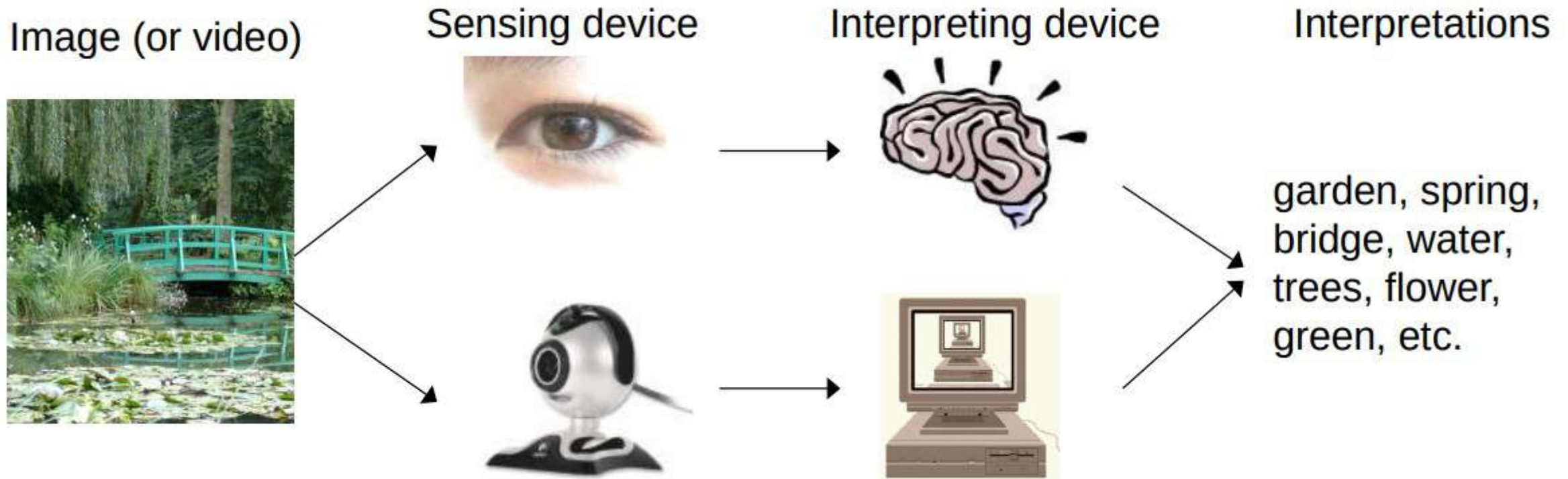




[Niebles]



# What is (Computer) Vision ?



[Niebles, Fergus]

# What is (Computer) Vision?

- *Vision* is about discovering from images
  - **what** is present in the scene, and
  - **where** it is
- In *Computer Vision* a **camera** (or several **cameras**) is linked to a computer
- The computer interprets **images** of a scene to obtain information
- Useful for tasks such as navigation, manipulation and recognition



# What is Computer Vision?

Computer Vision's **goal** is to obtain a high-level **understanding** of the world using images as input



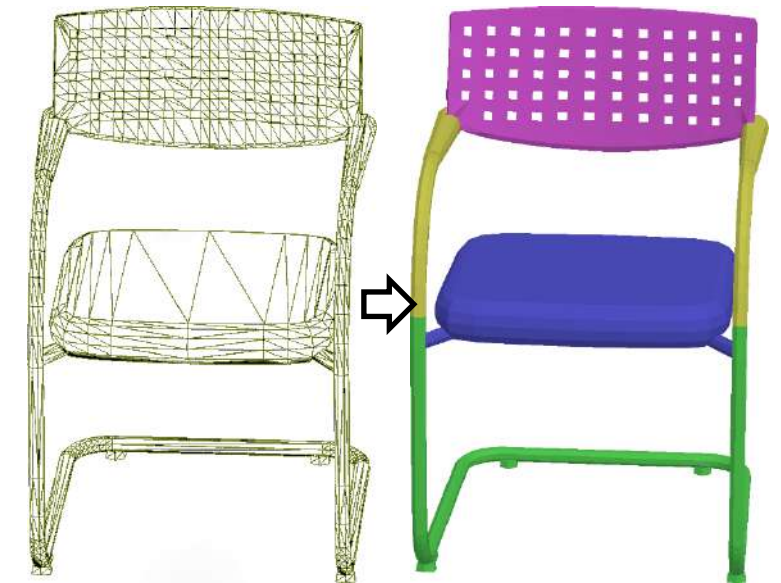
# Understand = Obtain Semantics & Geometry



3D Object layout  
Input: RGBD Image



Building facade segmentation  
Input: RGB Image



Object semantic segmentation  
Input: 3D Mesh





# Today's Agenda

- Course Overview
- Introduction to Computer Vision
  - What is Computer Vision
  - How hard is Computer Vision
  - Why is Computer Vision so hard
  - How to organize Computer Vision
  - Why study Computer Vision
  - Applications
- What we do



# How hard is Computer Vision ?

- The Summer Vision Project – MIT AI Memo 100, 1966
  - ‘Solve vision in a summer project’ - almost an [urban legend](#)
    - Basic foreground/background segmentation,
    - Analyse scenes with simple non-overlapping objects,
    - Extend the system to more complex objects.

MASSACHUSETTS INSTITUTE OF TECHNOLOGY  
PROJECT MACArtificial Intelligence Group  
Vision Memo. No. 100.

July 7, 1966

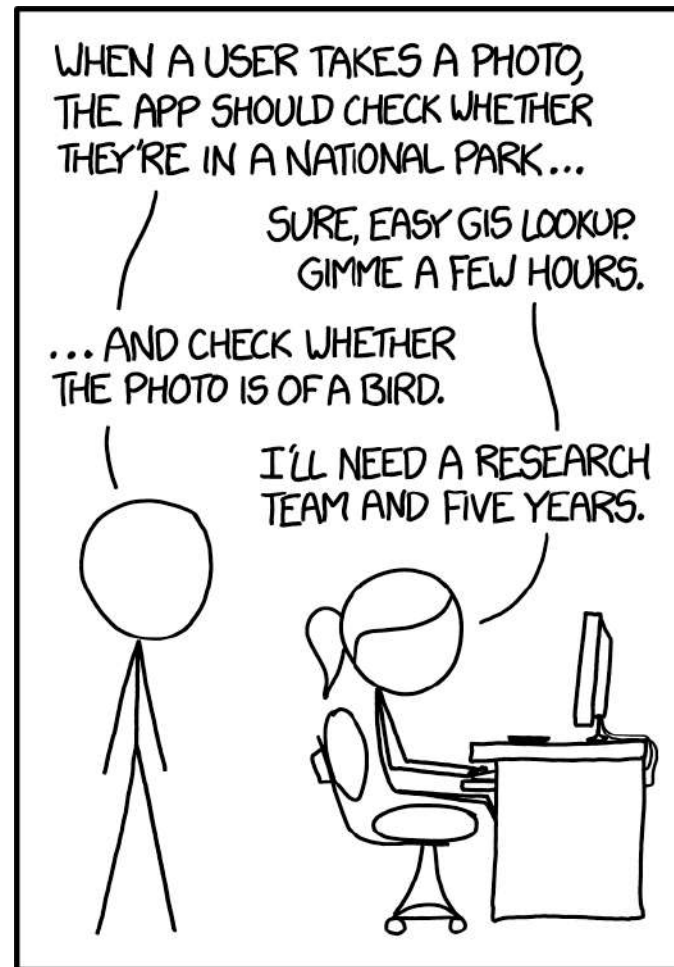
THE SUMMER VISION PROJECT

Seymour Papert

The summer vision project is an attempt to use our summer workers effectively in the construction of a significant part of a visual system. The particular task was chosen partly because it can be segmented into sub-problems which will allow individuals to work independently and yet participate in the construction of a system complex enough to be a real landmark in the development of "pattern recognition".



# How hard is Computer Vision ?



[XKCD](#)

IN CS, IT CAN BE HARD TO EXPLAIN THE DIFFERENCE BETWEEN THE EASY AND THE VIRTUALLY IMPOSSIBLE.

# How hard is Computer Vision ?

**PARK or BIRD**

Want to know if your photo is from a U.S. national park? Want to know if it contains a bird? Just drag it into the box to the left, and we'll tell you. We'll use the GPS embedded in your photo (if it's there) to see whether it's from a park, and we'll use our super-cool computer vision skills to try to see whether it's a bird (which is a hard problem, but we do a pretty good job at it).

To try it out, just drag any photo from your desktop into the upload box, or try dragging any of our example images. We'll give you your answers below!

Want to know more about PARK or BIRD, including why the heck we did this? Just [click here for more info](#) → ⓘ

**PARK? YES**  
Hey, yeah! I went to Bryce Canyon once!

**BIRD? NO**  
Beautiful clouds, but I don't see any birds flying up there.

EXAMPLE PHOTOS

Photo credits

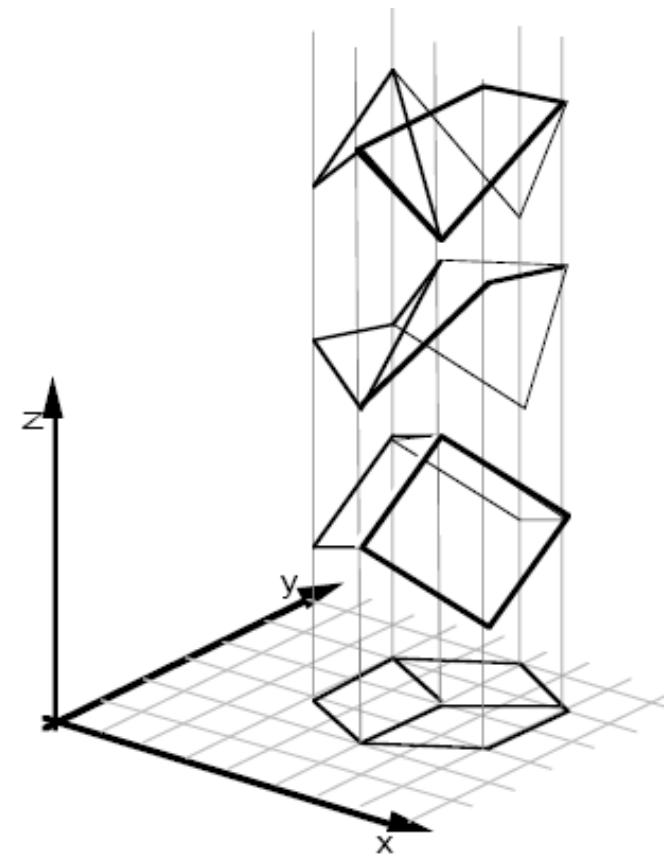
[Flickr 'solved' it](#)

# Today's Agenda

- Course Overview
- Introduction to Computer Vision
  - What is Computer Vision
  - How hard is Computer Vision
  - Why is Computer Vision so hard
  - How to organize Computer Vision
  - Why study Computer Vision
  - Applications
- What we do

# Why is Computer Vision so hard?

Because it is an ill-posed problem



[Sinha and Adelson 1993]



# Challenge 1: viewpoint variation



Madonna della Pietà,  
Michelangelo Buonarroti, 1498-99



[Fei Fei, Fergus & Torralba]

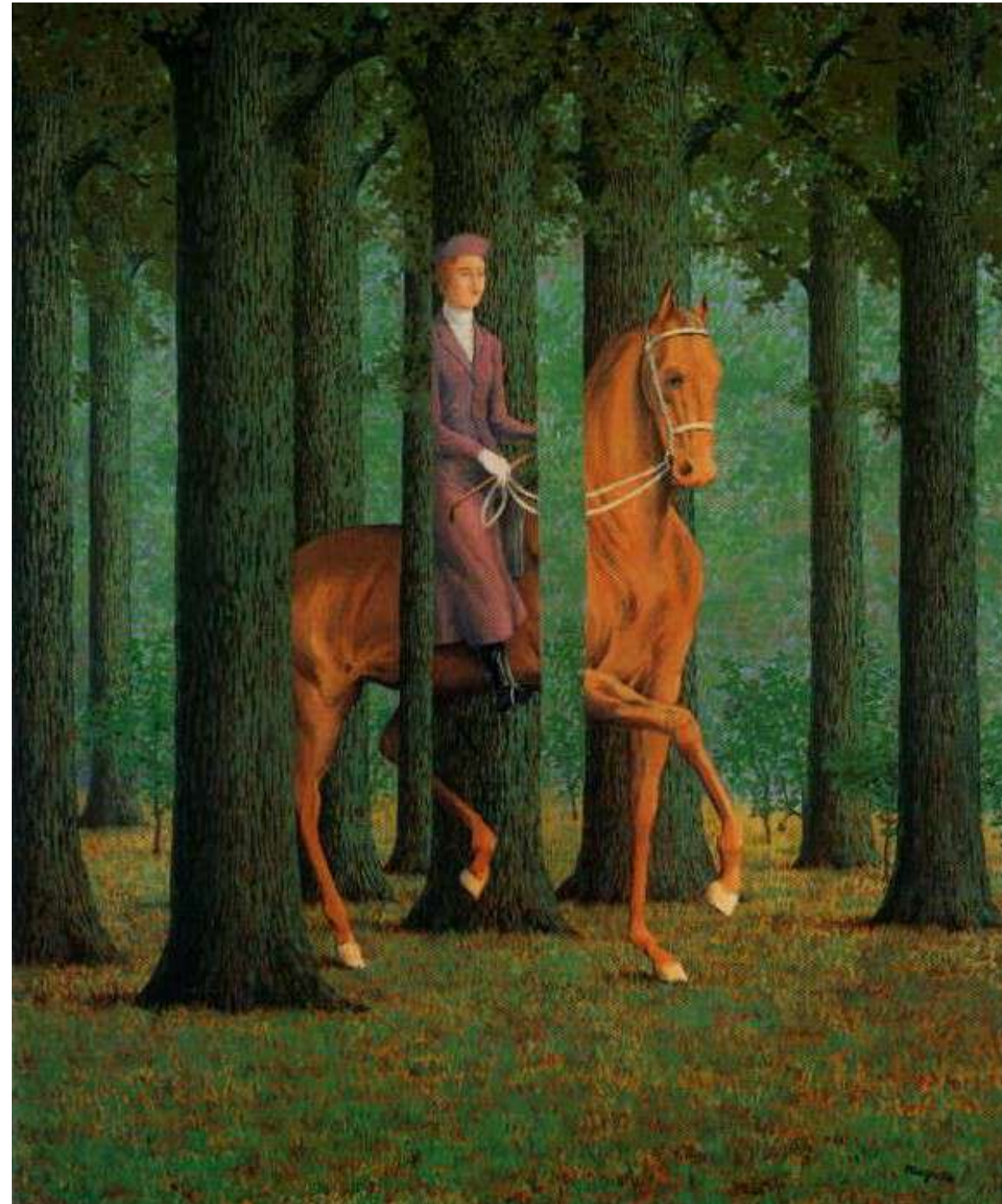
## Challenge 2: illumination



[S. Ullman]



# Challenge 3: occlusion



The Blank Signature,  
Rene Magritte, 1965

[Fei Fei, Fergus & Torralba]

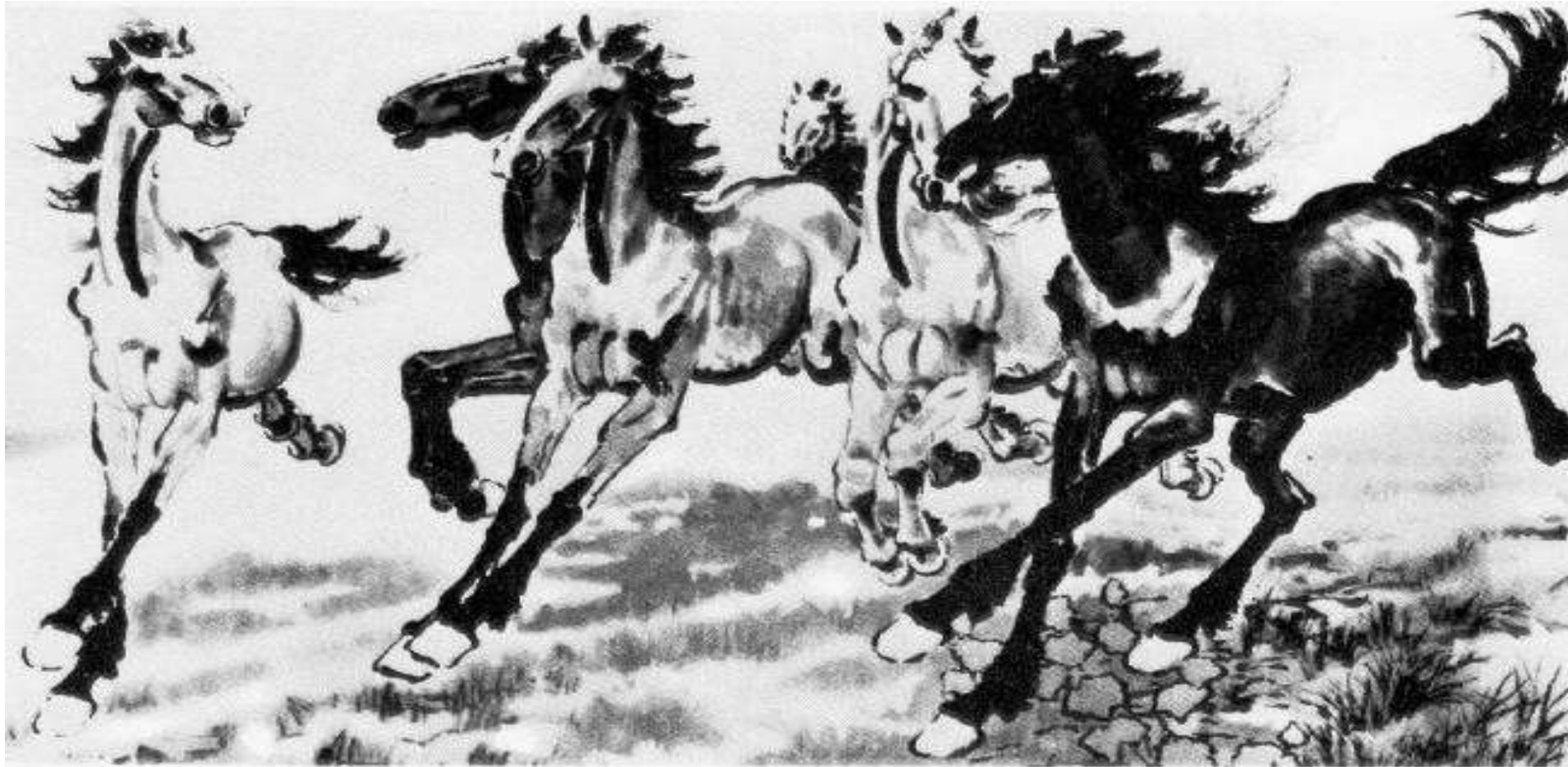
## Challenge 4: scale



[Fei Fei, Fergus & Torralba]



## Challenge 5: deformation



Six Galloping Horses,  
Xu Beihong, 1942

[Fei Fei, Fergus & Torralba]



## Challenge 6: background clutter



The Maiden,  
Gustav Klimt, 1913

[Fei Fei, Fergus &  
Torralba]

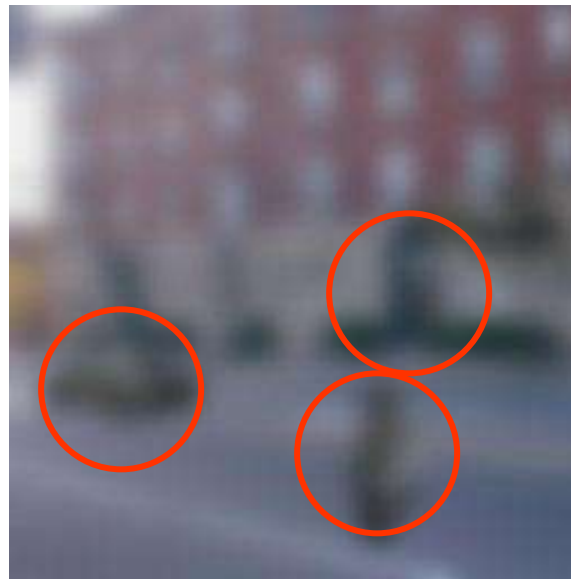
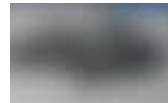
## Challenge 7: object intra-class variation



[Fei Fei, Fergus & Torralba]



## Challenge 8: local ambiguity



[Fei Fei, Fergus & Torralba]

# Today's Agenda

- Course Overview
- Introduction to Computer Vision
  - What is Computer Vision
  - How hard is Computer Vision
  - Why is Computer Vision so hard
  - How to organize Computer Vision
  - Why study Computer Vision
  - Applications
- What we do

# How to organize Computer Vision ?

0	3	2	5	4	7	6	9	8	0	3	2	5	4	7	6	9	8	0	3	2	5	4	7	6	9	8	0	3	2	5	4	7	6	9	8
3	0	1	2	3	4	5	6	7	3	0	1	2	3	4	5	6	7	3	0	1	2	3	4	5	6	7	3	0	1	2	3	4	5	6	7
2	1	0	3	2	5	4	7	6	2	1	0	3	2	5	4	7	6	2	1	0	3	2	5	4	7	6	2	1	0	3	2	5	4	7	6
5	2	3	0	1	2	3	4	5	5	2	3	0	1	2	3	4	5	5	2	3	0	1	2	3	4	5	5	2	3	0	1	2	3	4	5
4	3	2	1	0	3	2	5	4	4	3	2	1	0	3	2	5	4	4	3	2	1	0	3	2	5	4	4	3	2	1	0	3	2	5	4
7	4	5	2	3	0	1	2	3	7	4	5	2	3	0	1	2	3	7	4	5	2	3	0	1	2	3	7	4	5	2	3	0	1	2	3
6	5	4	3	2	1	0	3	2	6	5	4	3	2	1	0	3	2	6	5	4	3	2	1	0	3	2	6	5	4	3	2	1	0	3	2
9	6	7	4	5	2	3	0	1	9	6	7	4	5	2	3	0	1	9	6	7	4	5	2	3	0	1	9	6	7	4	5	2	3	0	1
8	7	6	5	4	3	2	1	0	8	7	6	5	4	3	2	1	0	8	7	6	5	4	3	2	1	0	8	7	6	5	4	3	2	1	0
0	3	2	5	4	7	6	9	8	0	3	2	5	4	7	6	9	8	0	3	2	5	4	7	6	9	8	0	3	2	5	4	7	6	9	8
3	0	1	2	3	4	5	6	7	3	0	1	2	3	4	5	6	7	3	0	1	2	3	4	5	6	7	3	0	1	2	3	4	5	6	7
2	1	0	3	2	5	4	7	6	2	1	0	3	2	5	4	7	6	2	1	0	3	2	5	4	7	6	2	1	0	3	2	5	4	7	6
5	2	3	0	1	2	3	4	5	5	2	3	0	1	2	3	4	5	5	2	3	0	1	2	3	4	5	5	2	3	0	1	2	3	4	5
4	3	2	1	0	3	2	5	4	4	3	2	1	0	3	2	5	4	4	3	2	1	0	3	2	5	4	4	3	2	1	0	3	2	5	4
7	4	5	2	3	0	1	2	3	7	4	5	2	3	0	1	2	3	7	4	5	2	3	0	1	2	3	7	4	5	2	3	0	1	2	3
6	5	4	3	2	1	0	3	2	6	5	4	3	2	1	0	3	2	6	5	4	3	2	1	0	3	2	6	5	4	3	2	1	0	3	2
9	6	7	4	5	2	3	0	1	9	6	7	4	5	2	3	0	1	9	6	7	4	5	2	3	0	1	9	6	7	4	5	2	3	0	1
8	7	6	5	4	3	2	1	0	8	7	6	5	4	3	2	1	0	8	7	6	5	4	3	2	1	0	8	7	6	5	4	3	2	1	0
0	3	2	5	4	7	6	9	8	0	3	2	5	4	7	6	9	8	0	3	2	5	4	7	6	9	8	0	3	2	5	4	7	6	9	8
3	0	1	2	3	4	5	6	7	3	0	1	2	3	4	5	6	7	3	0	1	2	3	4	5	6	7	3	0	1	2	3	4	5	6	7
2	1	0	3	2	5	4	7	6	2	1	0	3	2	5	4	7	6	2	1	0	3	2	5	4	7	6	2	1	0	3	2	5	4	7	6
5	2	3	0	1	2	3	4	5	5	2	3	0	1	2	3	4	5	5	2	3	0	1	2	3	4	5	5	2	3	0	1	2	3	4	5
4	3	2	1	0	3	2	5	4	4	3	2	1	0	3	2	5	4	4	3	2	1	0	3	2	5	4	4	3	2	1	0	3	2	5	4
7	4	5	2	3	0	1	2	3	7	4	5	2	3	0	1	2	3	7	4	5	2	3	0	1	2	3	7	4	5	2	3	0	1	2	3
6	5	4	3	2	1	0	3	2	6	5	4	3	2	1	0	3	2	6	5	4	3	2	1	0	3	2	6	5	4	3	2	1	0	3	2
9	6	7	4	5	2	3	0	1	9	6	7	4	5	2	3	0	1	9	6	7	4	5	2	3	0	1	9	6	7	4	5	2	3	0	1
8	7	6	5	4	3	2	1	0	8	7	6	5	4	3	2	1	0	8	7	6	5	4	3	2	1	0	8	7	6	5	4	3	2	1	0

[Shapiro]





[Shapiro]

# How to organize Computer Vision ?

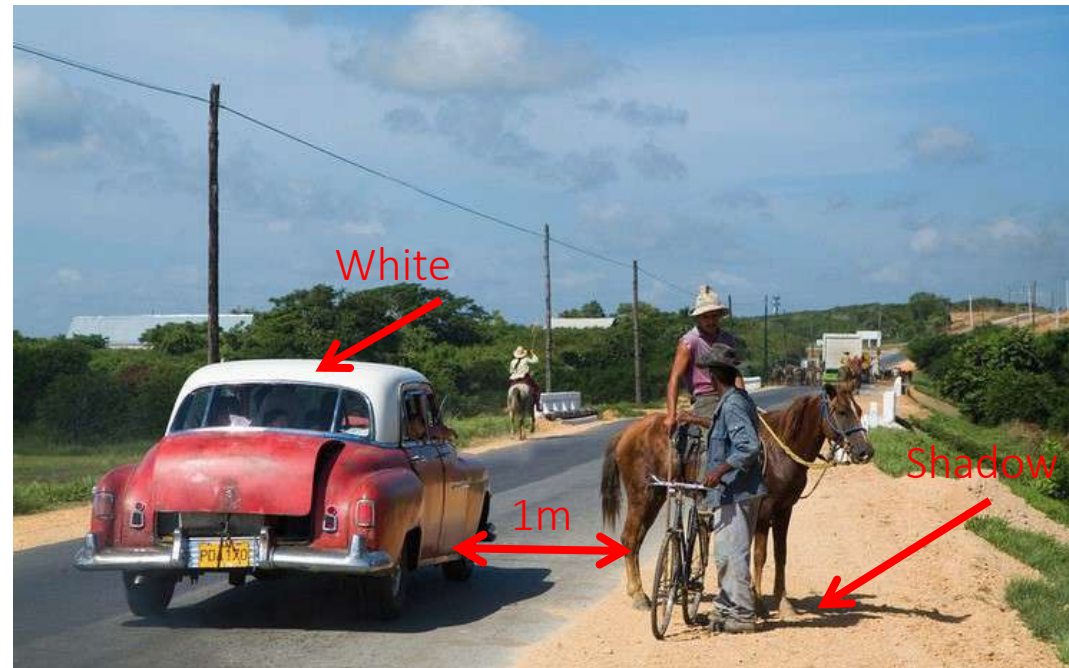
- Low Level Vision
  - Measurements
  - Enhancements
  - Region segmentation
  - Features
- Mid Level Vision
  - Reconstruction
  - Depth
  - Motion Estimation
- High Level Vision
  - Category detection
  - Activity recognition
  - Deep understanding



[Shapiro]

# How to organize Computer Vision ?

- Low Level Vision
  - Measurements
  - Enhancements
  - Region segmentation
  - Features
- Mid Level Vision
  - Reconstruction
  - Depth
  - Motion Estimation
- High Level Vision
  - Category detection
  - Activity recognition
  - Deep understanding



[Shapiro]



# Low-Level: Exposure

---



[Redmon]

# Low-Level: Edges

---



[Redmon]



# Low-Level: Segmentation (color)

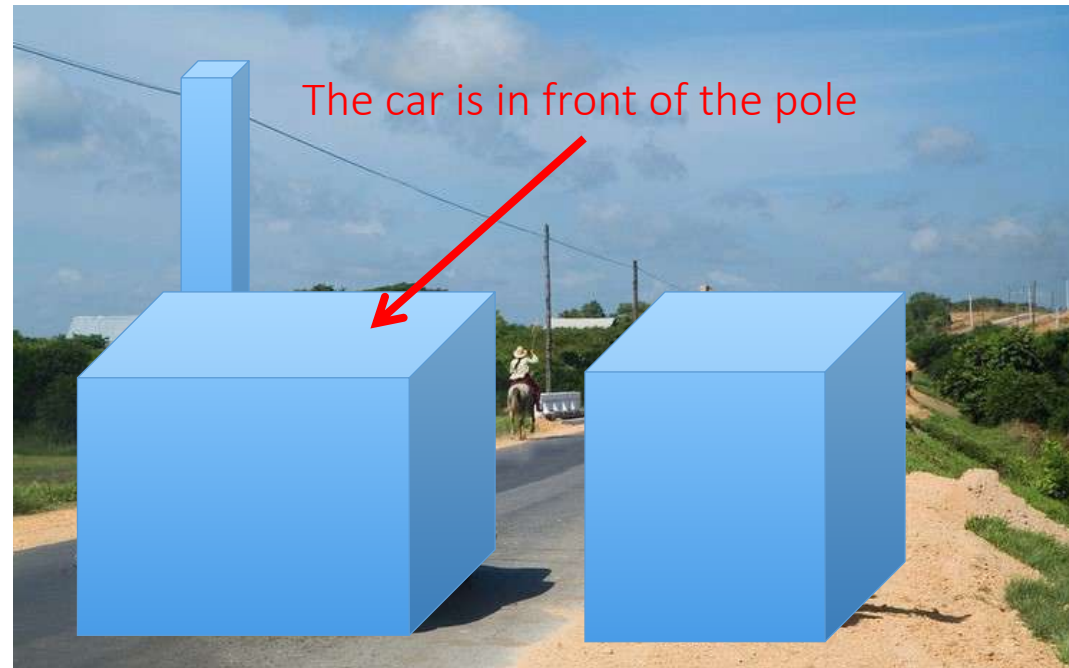
---



[Redmon]

# How to organize Computer Vision ?

- Low Level Vision
  - Measurements
  - Enhancements
  - Region segmentation
  - Features
- Mid Level Vision
  - Reconstruction
  - Depth
  - Motion Estimation
- High Level Vision
  - Category detection
  - Activity recognition
  - Deep understanding



[Shapiro]

# Mid-Level: Panorama Stitching

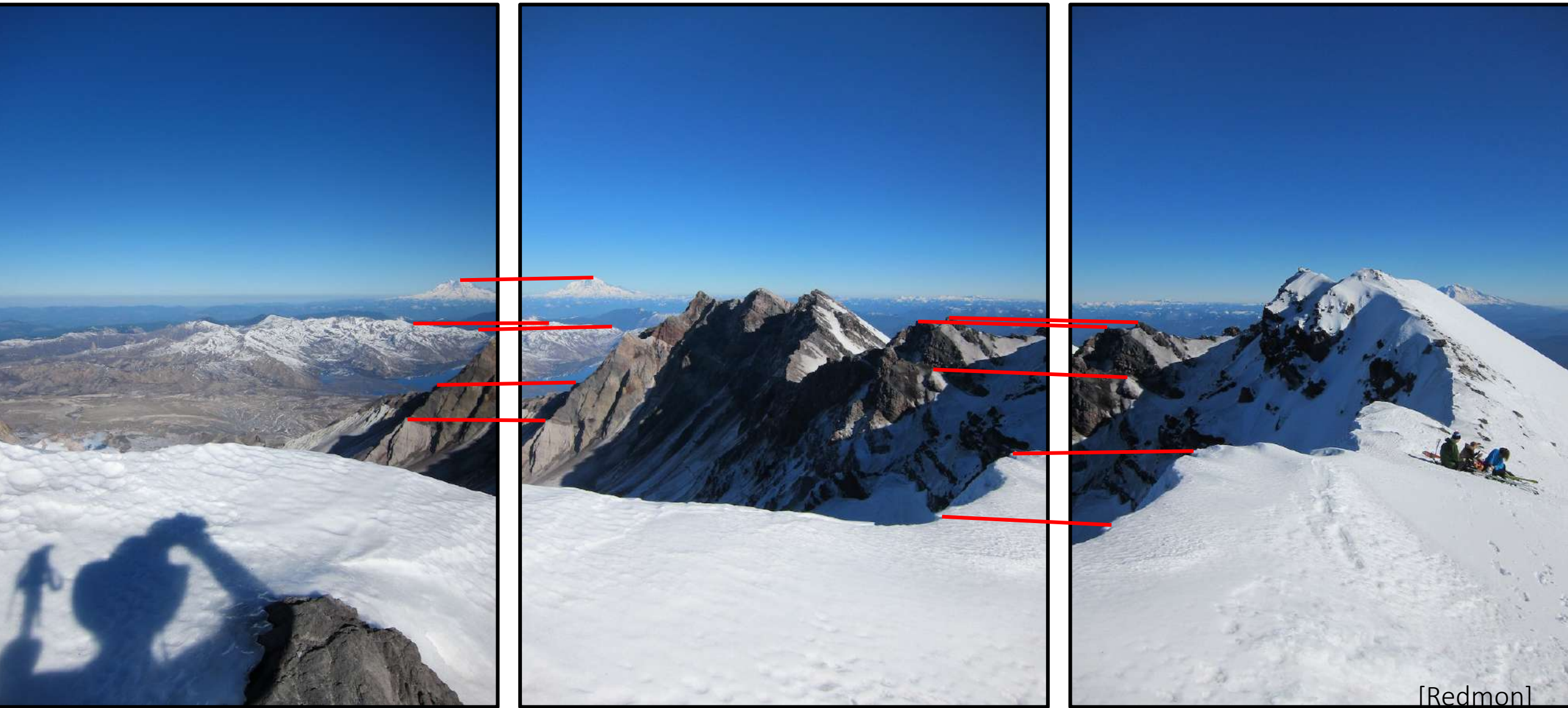
---



[Redmon]



## Mid-Level: Panorama Stitching



[Redmon]

# Mid-Level: Panorama Stitching

---



[Redmon]



# Mid-Level: Multi-View Stereo

---



[Redmon]

# Mid-Level: Multi-View Stereo



[Redmon]

# Mid-Level: Multi-View Stereo

---



[Building Rome in a Day,  
Agarwal et al., ICCV 2009]

The Colosseum, 2,106 images, 819,242 points

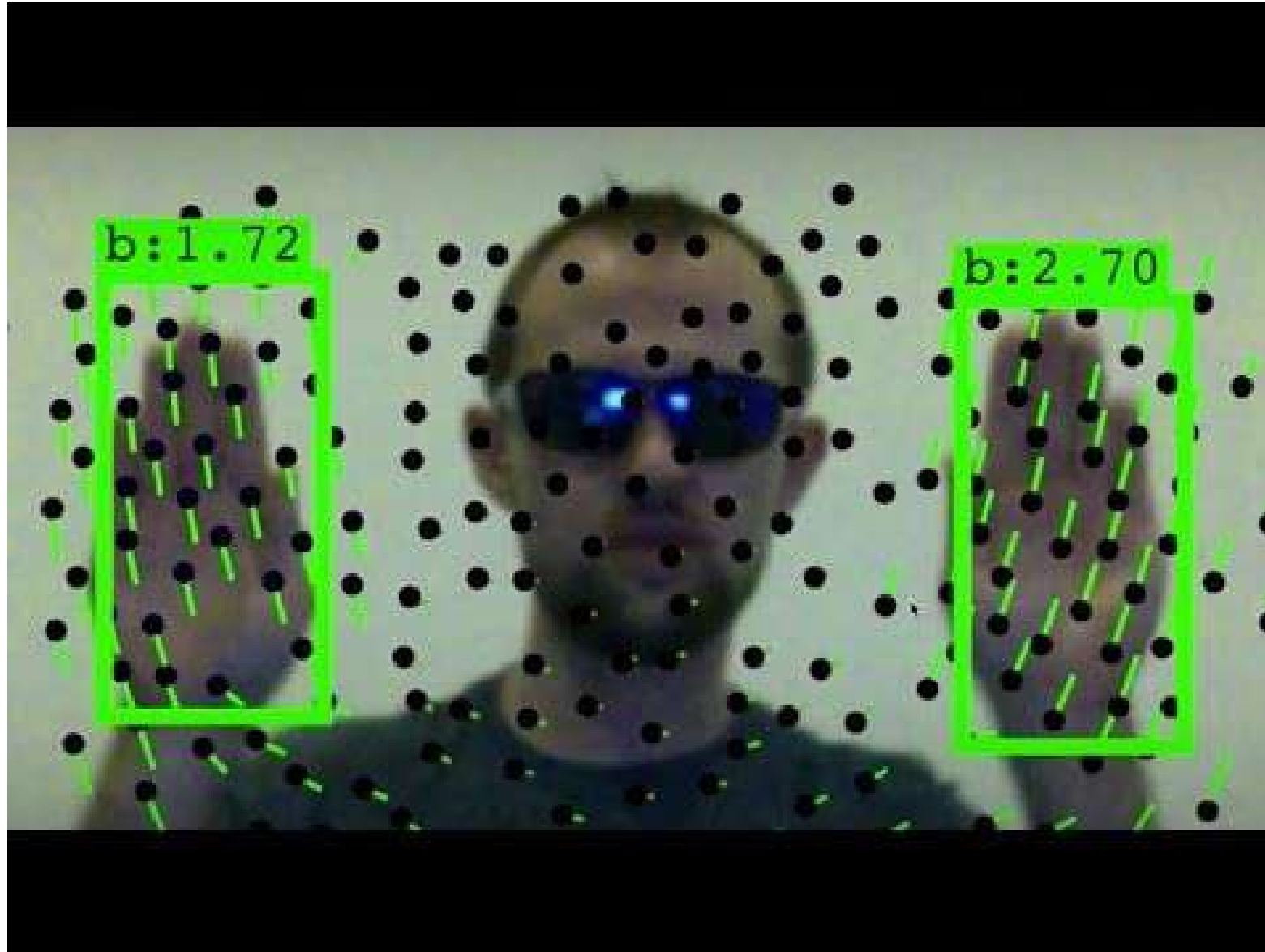
# Mid-Level: Optical Flow



[Redmon]



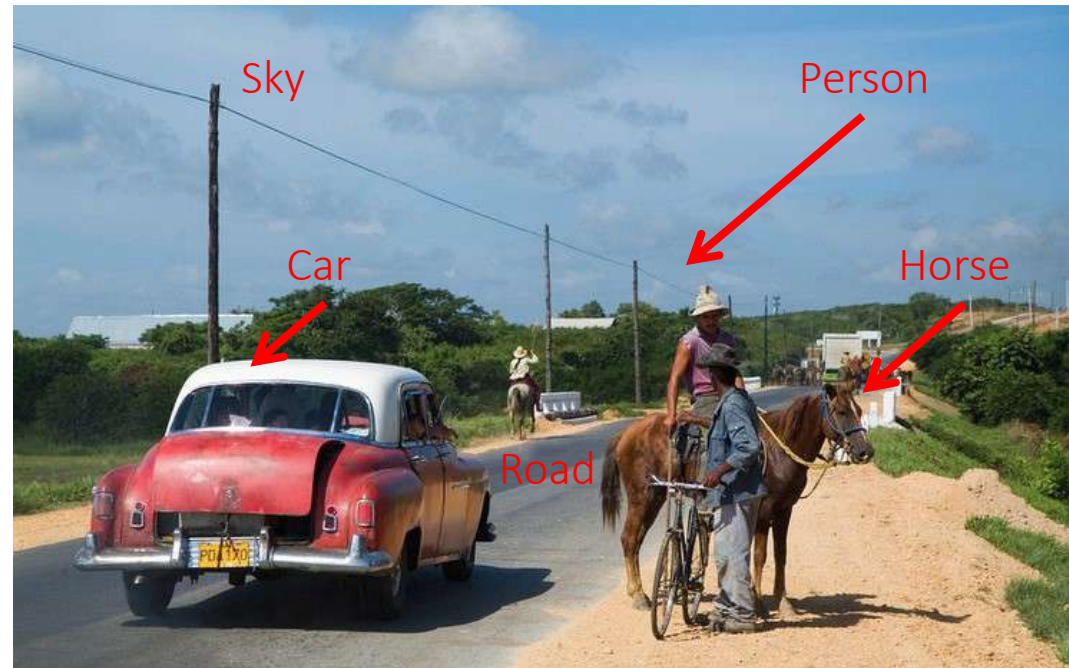
# Mid-Level: Optical Flow



[Redmon]

# How to organize Computer Vision ?

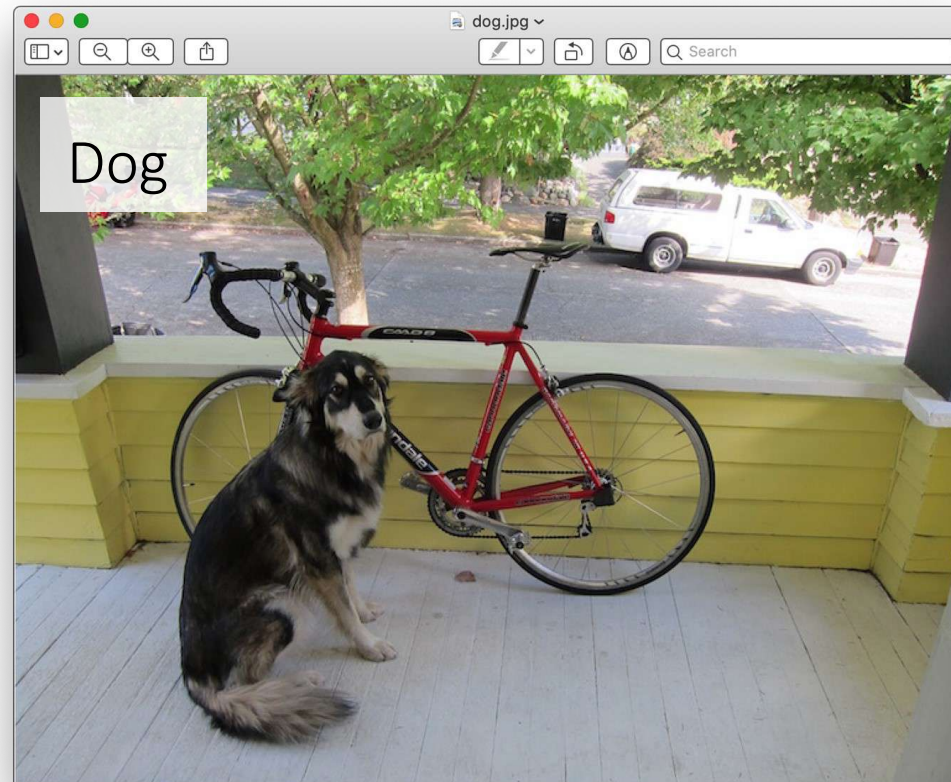
- Low Level Vision
  - Measurements
  - Enhancements
  - Region segmentation
  - Features
- Mid Level Vision
  - Reconstruction
  - Depth
  - Motion Estimation
- High Level Vision
  - Category detection
  - Activity recognition
  - Deep understanding



[Shapiro]

# High-Level: Classification

- What is in the image?

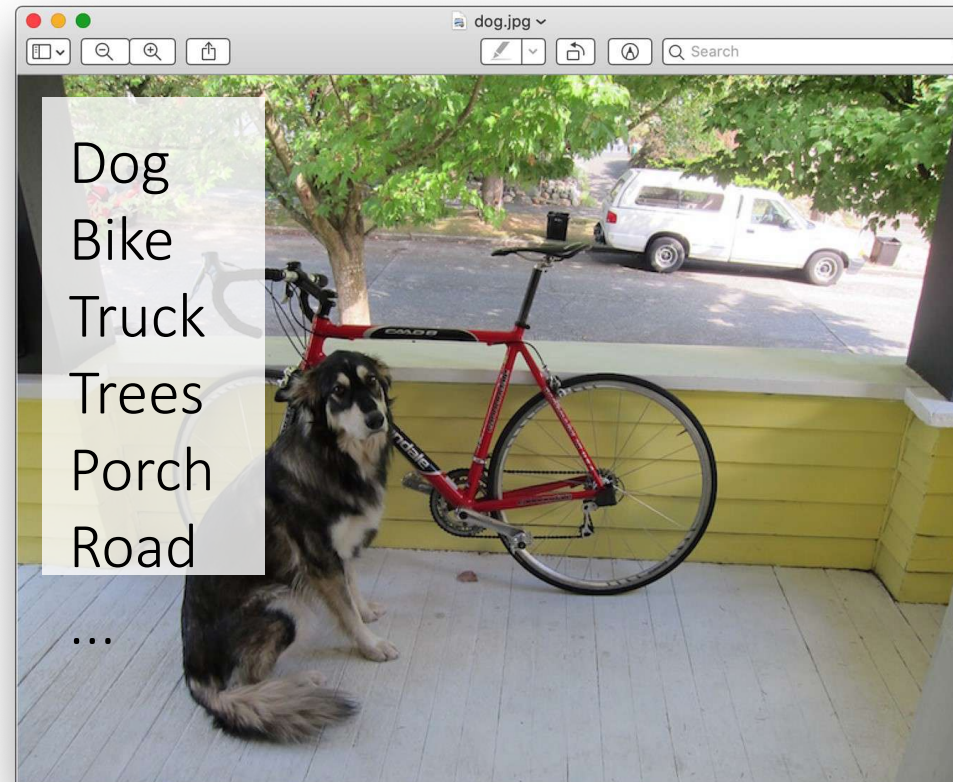


[Redmon]



# High-Level: Tagging

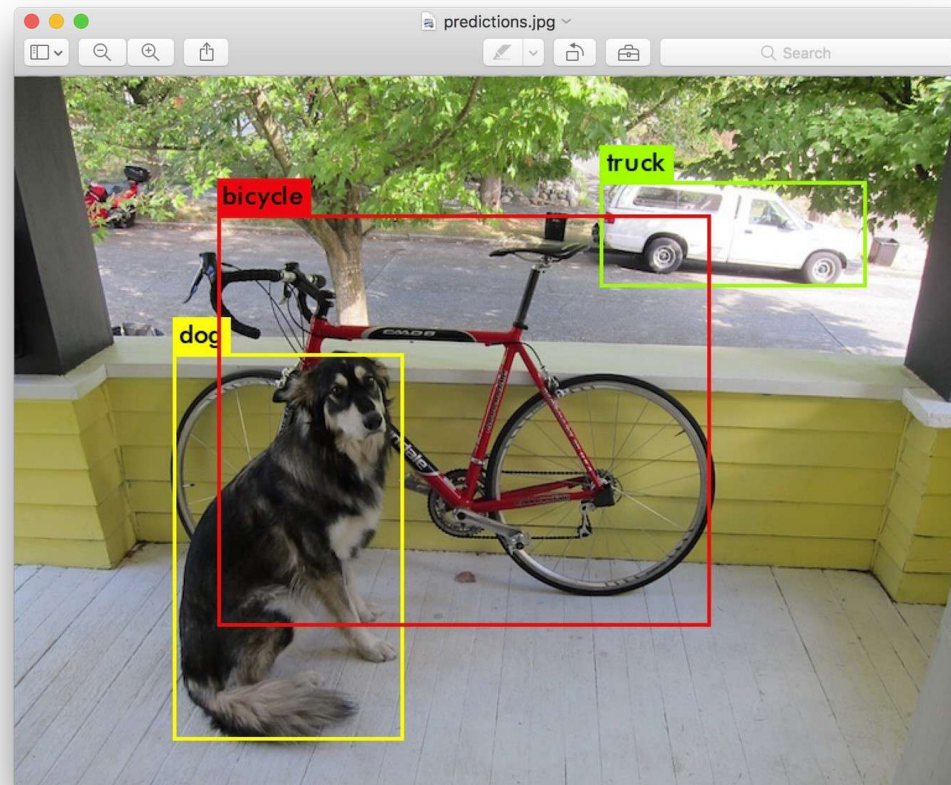
- What are ALL the things in the image?



[Redmon]

# High-Level: Detection

- What are ALL the things in the image?
- Where are they?



[Redmon]



# High-Level: Semantic Segmentation



[Redmon]



# High-Level: Instance Segmentation

---



<http://www.youtube.com/watch?v=OOT3UIXZztE>

[Redmon]

# Today's Agenda

- Course Overview
- Introduction to Computer Vision
  - What is Computer Vision
  - How hard is Computer Vision
  - Why is Computer Vision so hard
  - How to organize Computer Vision
  - Why study Computer Vision
  - Applications
- What we do

# Why study Computer Vision?

- *Match (or beat) human vision*
  - central to Artificial Intelligence, countless applications
- Understand human vision → neuroscience
- Do research with huge impact
- Get a job in the industry
- Timing is perfect: AI revolution – big data, faster hardware, deep learning



# Do research with huge impact

Google Scholar

Top publications

Categories ▾

Publication
1. Nature
2. The New England Journal of Medicine
3. Science
4. The Lancet
5. IEEE/CVF Conference on Computer Vision and Pattern Recognition
6. Advanced Materials
7. Nature Communications
8. Cell
9. Chemical Reviews
10. Chemical Society reviews
11. Journal of the American Chemical Society
12. Angewandte Chemie
13. Proceedings of the National Academy of Sciences
14. JAMA
15. Nucleic Acids Research

Google Scholar

Top publications

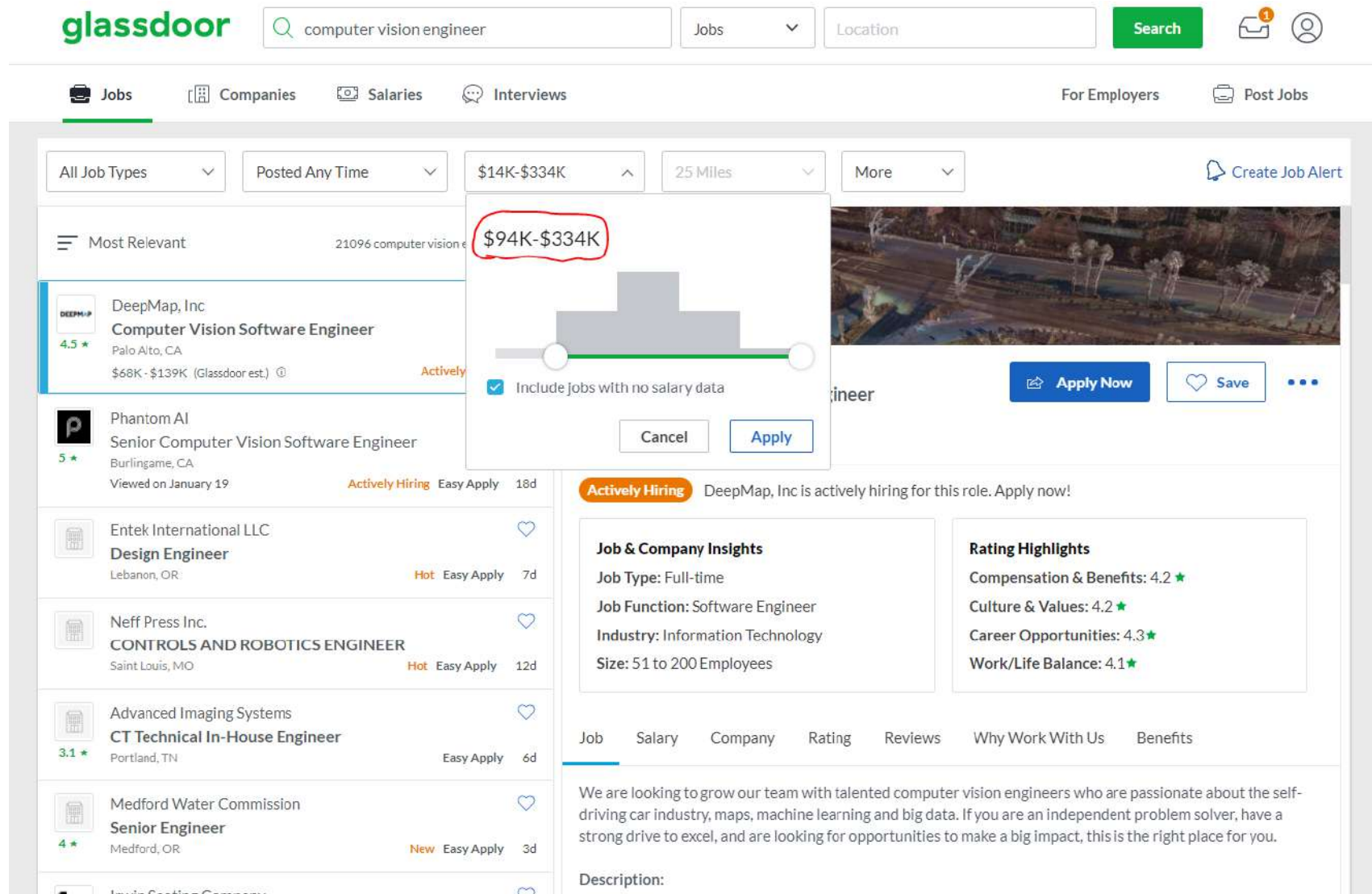
Categories ▾ English ▾

Publication	<u>h5-index</u>	<u>h5-median</u>
1. Nature	414	607
2. The New England Journal of Medicine	410	704
3. Science	391	564
4. IEEE/CVF Conference on Computer Vision and Pattern Recognition	356	583
5. The Lancet	345	600
6. Advanced Materials	294	406
7. Cell	288	459
8. Nature Communications	287	389
9. Chemical Reviews	270	434
10. International Conference on Learning Representations	253	470
11. JAMA	253	446
12. Neural Information Processing Systems	245	422
13. Proceedings of the National Academy of Sciences	245	337
14. Journal of the American Chemical Society	245	330
15. Angewandte Chemie	235	314

2021

2022

# Get a job in the industry



The screenshot shows a Glassdoor search for 'computer vision engineer' with filters for 'All Job Types', 'Posted Any Time', '\$14K-\$334K', and '25 Miles'. A salary filter overlay is active, showing a histogram and a selected range of '\$94K-\$334K'. The job list includes:

- DeepMap, Inc**: Computer Vision Software Engineer, Palo Alto, CA, \$68K - \$139K (Glassdoor est.), 4.5★, Actively Hiring.
- Phantom AI**: Senior Computer Vision Software Engineer, Burlingame, CA, 5★, Viewed on January 19, Actively Hiring, Easy Apply, 18d.
- Entek International LLC**: Design Engineer, Lebanon, OR, Hot, Easy Apply, 7d.
- Neff Press Inc.**: CONTROLS AND ROBOTICS ENGINEER, Saint Louis, MO, Hot, Easy Apply, 12d.
- Advanced Imaging Systems**: CT Technical In-House Engineer, Portland, TN, 3.1★, Easy Apply, 6d.
- Medford Water Commission**: Senior Engineer, Medford, OR, 4★, New, Easy Apply, 3d.

The job details for DeepMap, Inc are expanded, showing:

- Job & Company Insights**: Job Type: Full-time, Job Function: Software Engineer, Industry: Information Technology, Size: 51 to 200 Employees.
- Rating Highlights**: Compensation & Benefits: 4.2★, Culture & Values: 4.2★, Career Opportunities: 4.3★, Work/Life Balance: 4.1★.
- Description**: We are looking to grow our team with talented computer vision engineers who are passionate about the self-driving car industry, maps, machine learning and big data. If you are an independent problem solver, have a strong drive to excel, and are looking for opportunities to make a big impact, this is the right place for you.

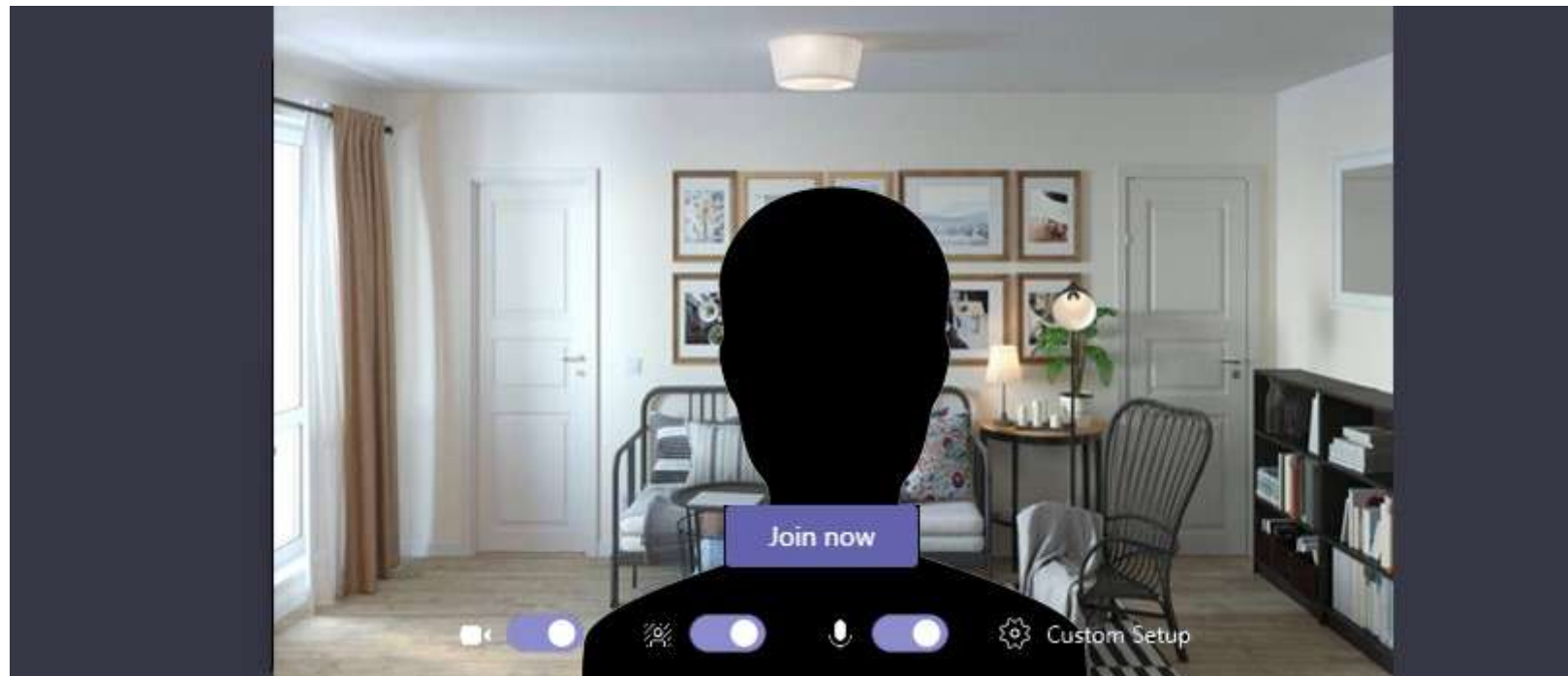


# Today's Agenda

- Course Overview
- Introduction to Computer Vision
  - What is Computer Vision
  - How hard is Computer Vision
  - Why is Computer Vision so hard
  - How to organize Computer Vision
  - Why study Computer Vision
  - Applications
- What we do



# Segmentation and Matting





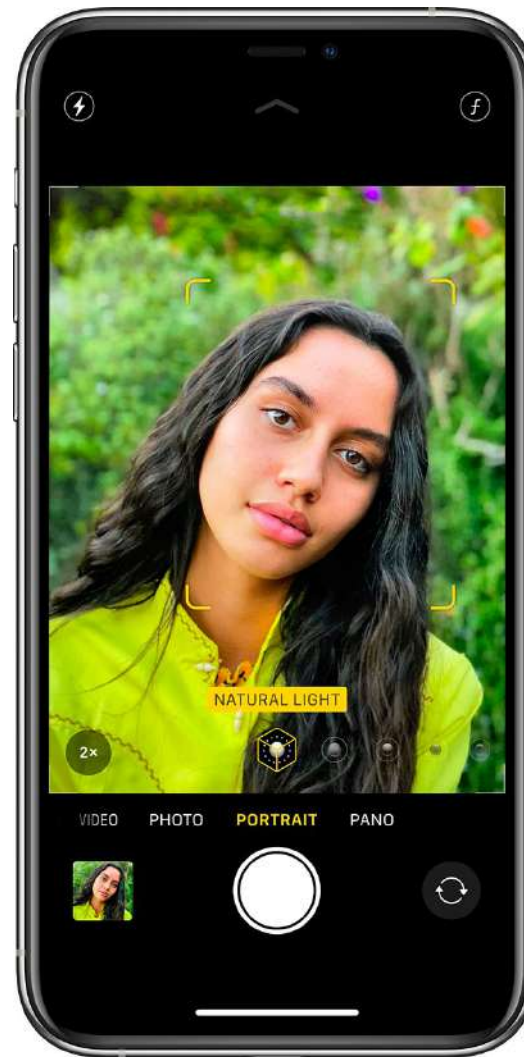
## 3D Maps



Apple Maps

[Seitz, Szeliski]

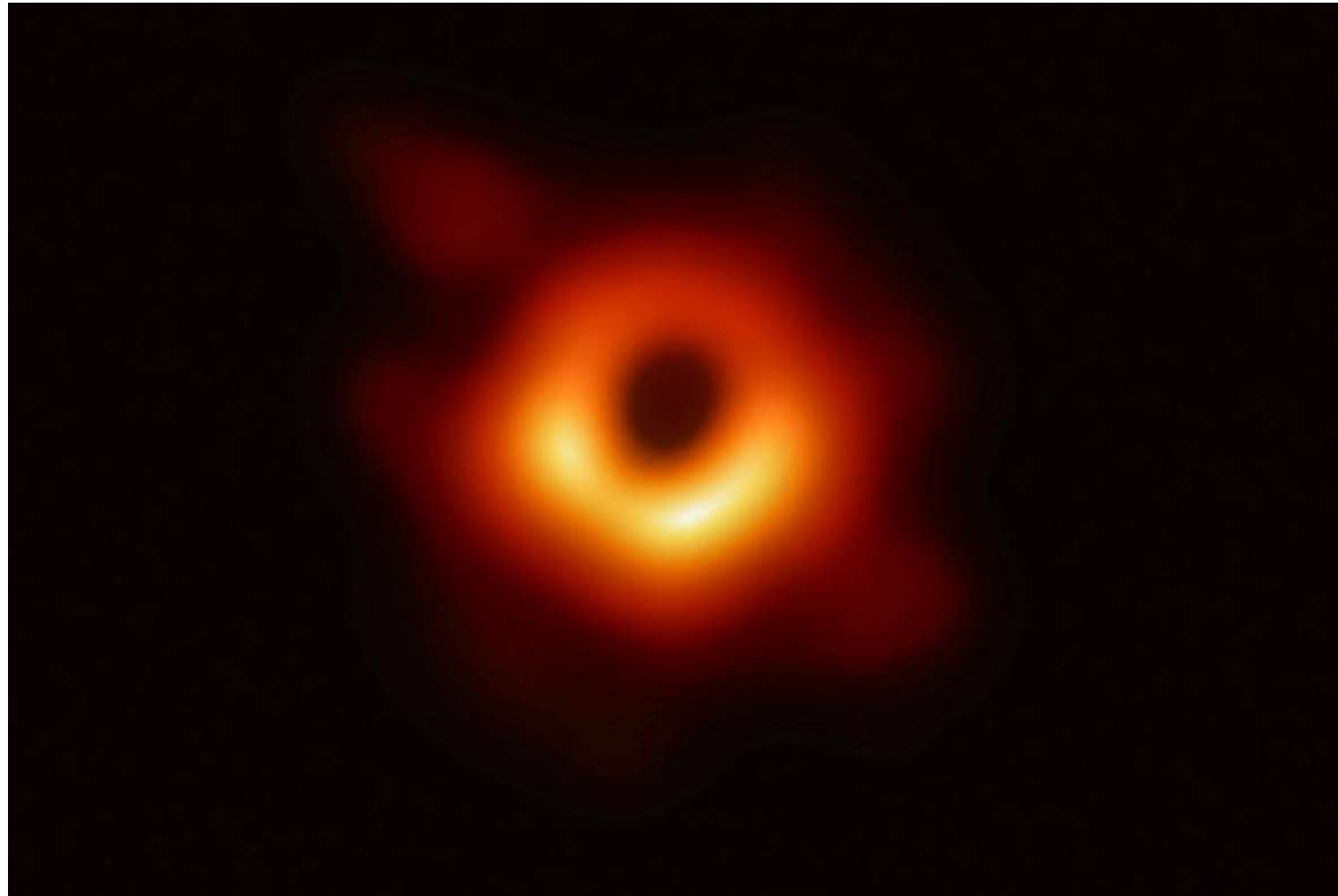
# Computational photography



Portrait mode  
simulating wider aperture

[Seitz, Szeliski]

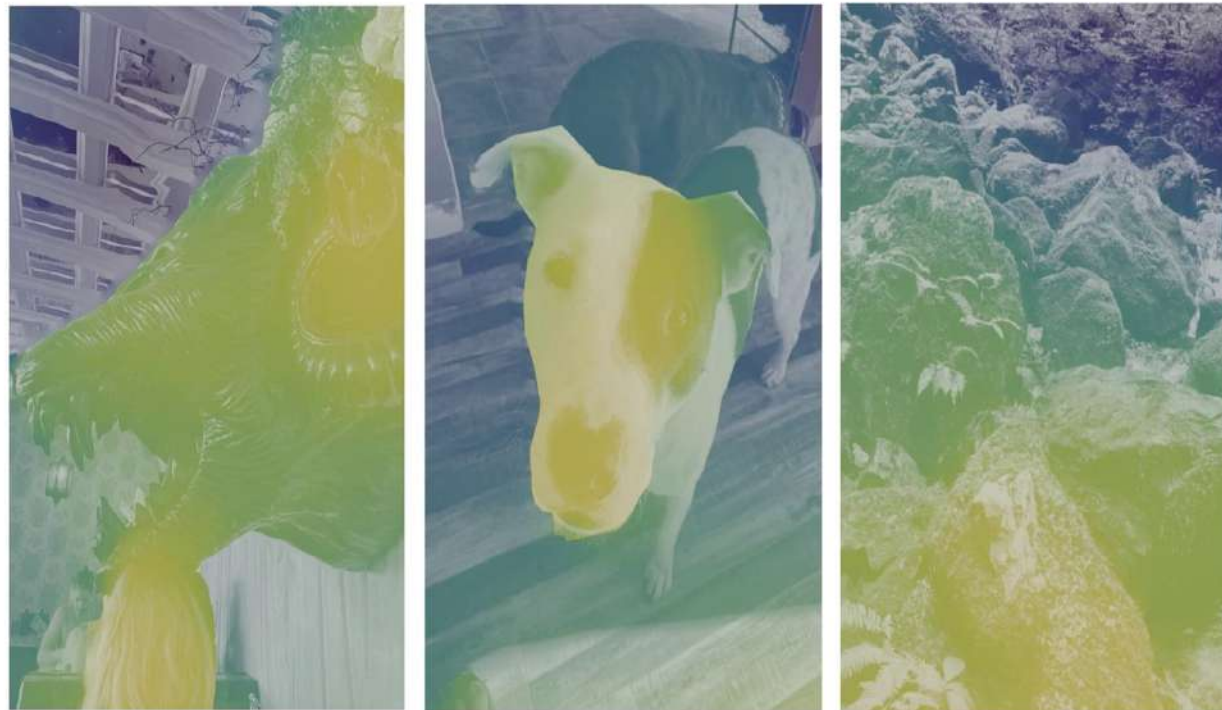
# Even wider aperture...



[How scientists captured the first image of a black hole, 2019](#)

[Seitz, Szeliski]

## 3D photos



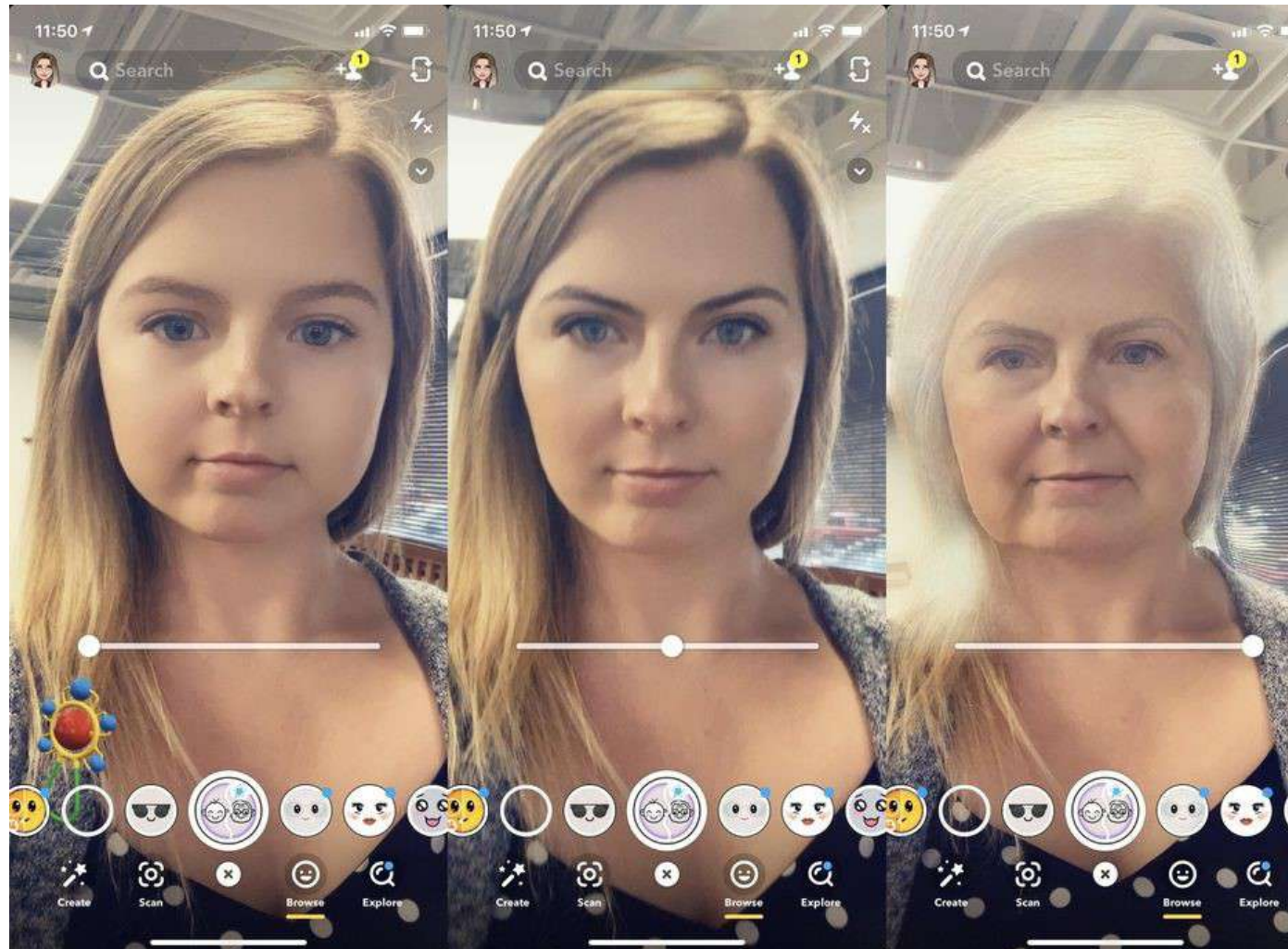
### [3D Photos on Facebook](#)

Estimate depth from photo to create animation

[Seitz, Szeliski]



# Age Simulation

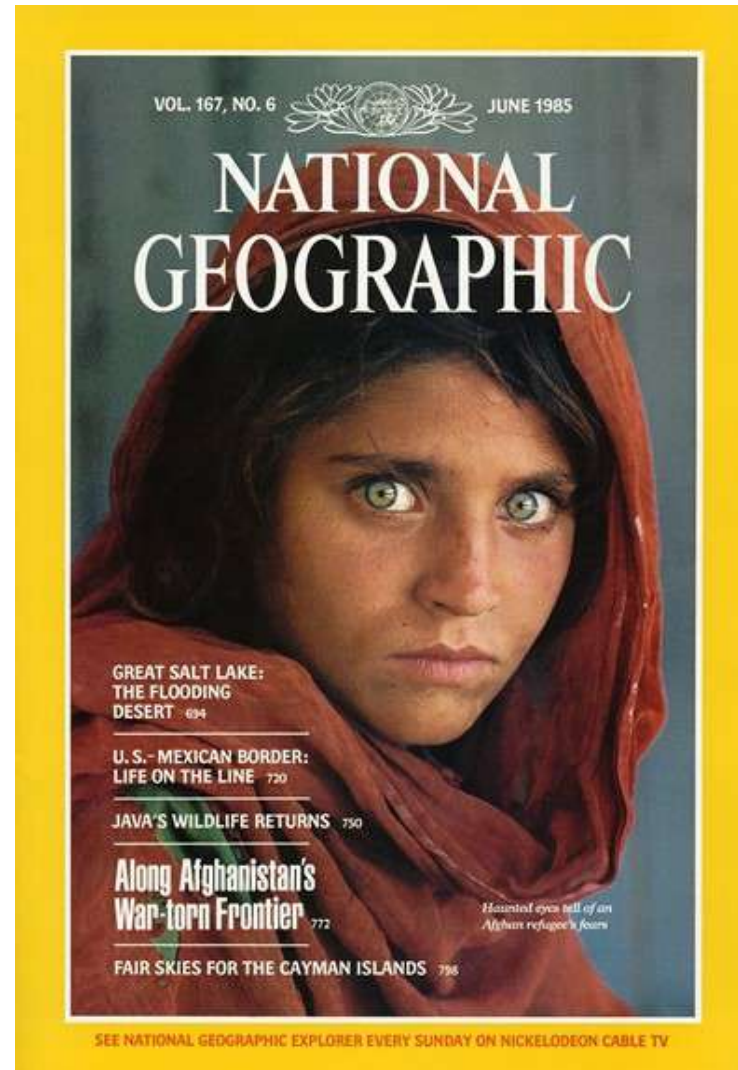


From [CNET](#)

Snapchat Time Machine

[Seitz, Szeliski]

# Face recognition

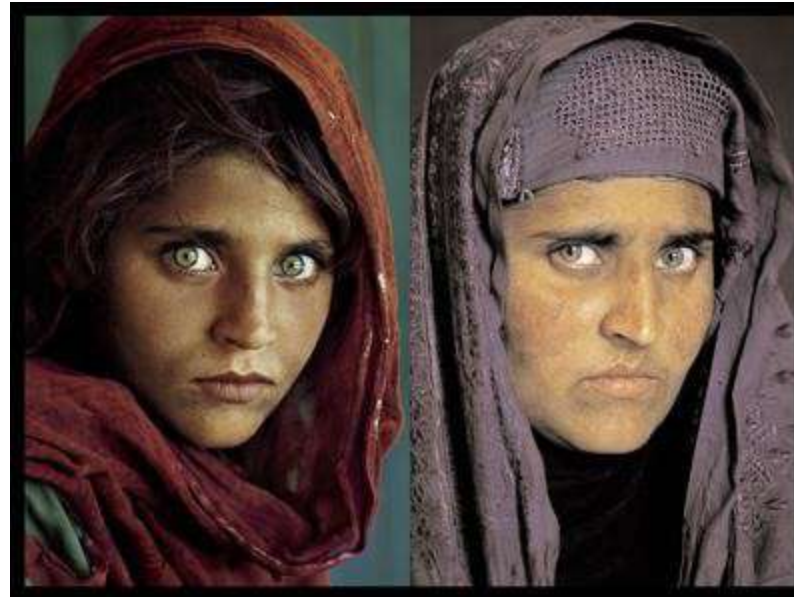


Who is she?

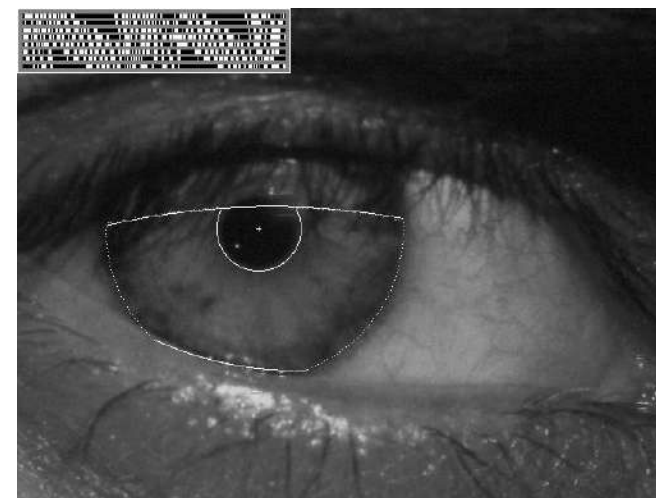
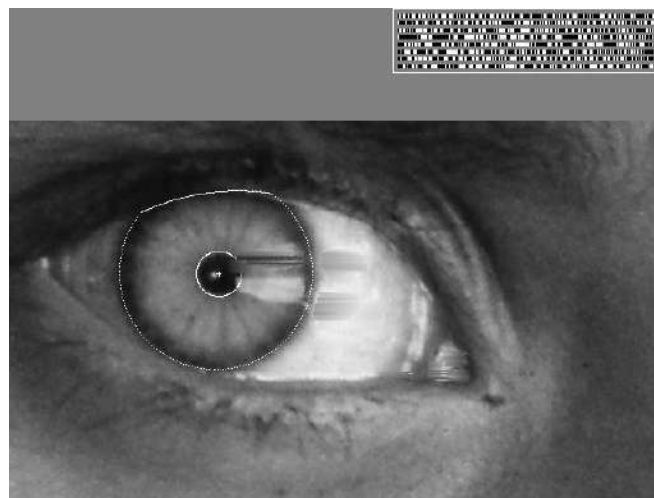
[Seitz, Szeliski]



# Vision-based biometrics

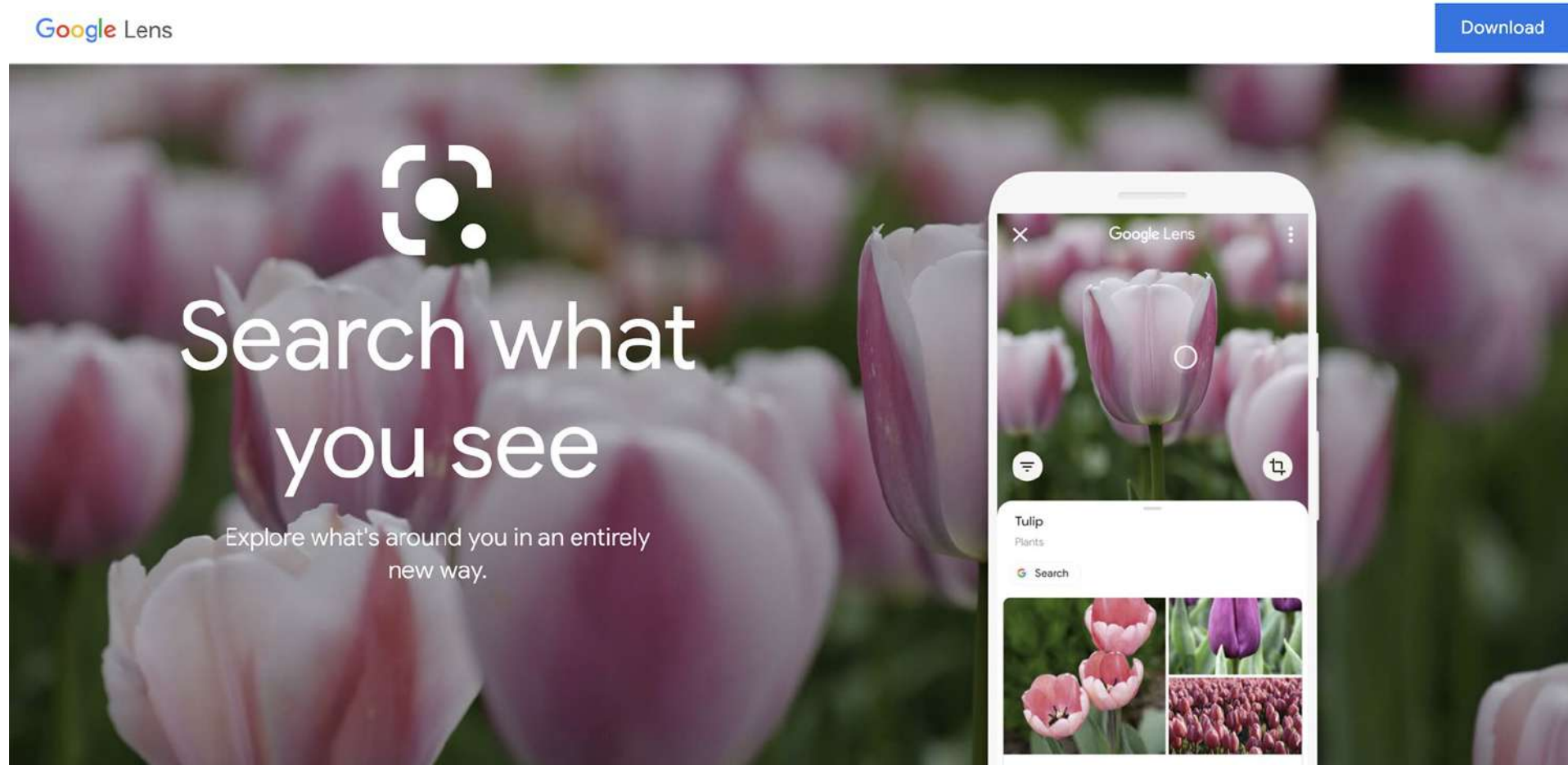


“How the Afghan Girl was Identified by Her Iris Patterns” Read the [story](#)



[Seitz, Szeliski]

# Object recognition



[Seitz, Szeliski]



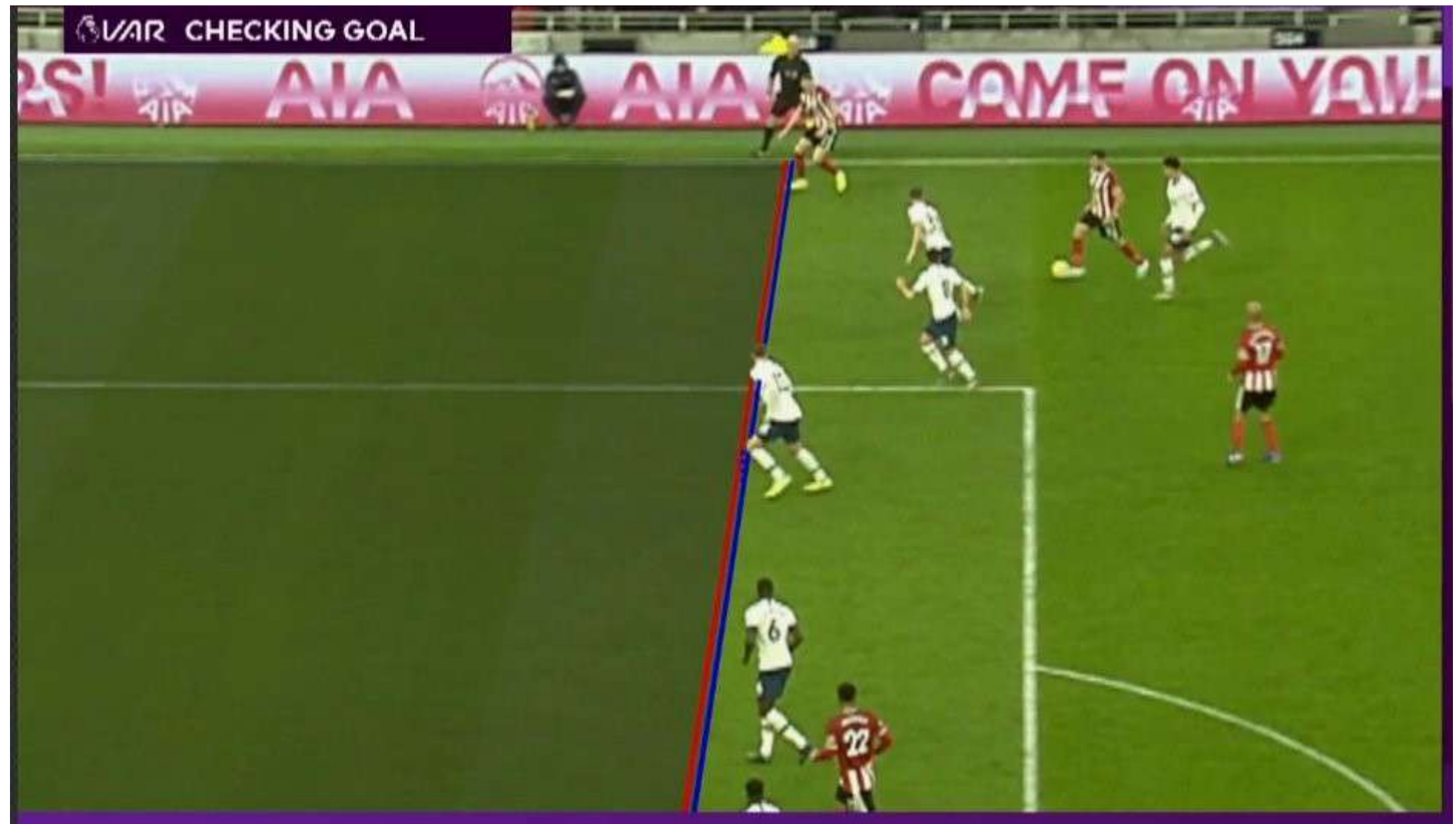
## Special effects: shape capture



*The Matrix* movies, ESC Entertainment, XYZRGB, NRC

[Seitz, Szeliski]

## Sports - VAR



## Games



Microsoft's XBox Kinect

[Seitz, Szeliski]



# Virtual Reality - Metaverse



Oculus Quest, Beat Saber

[Seitz, Szeliski]



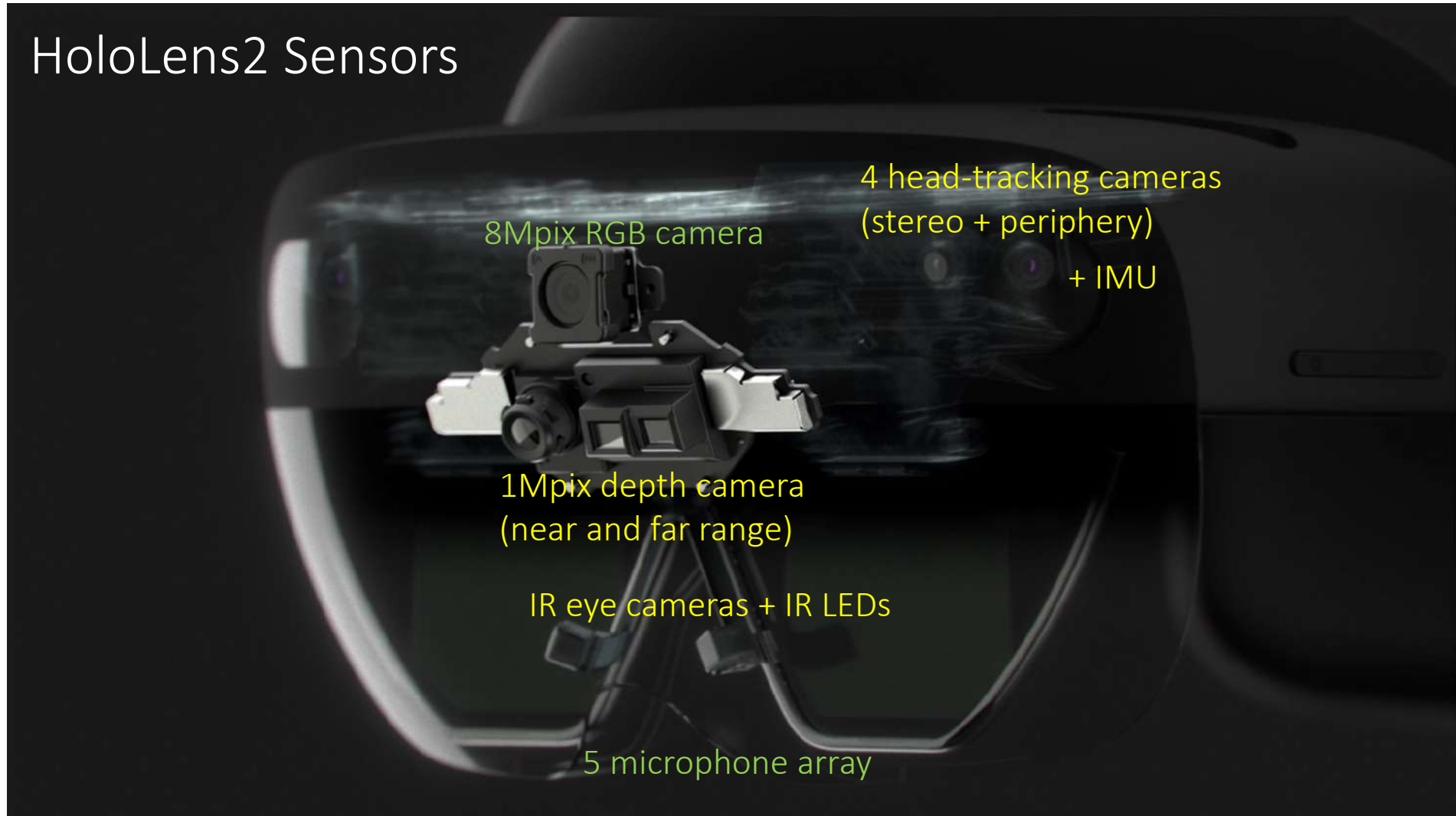
# Augmented Reality - Metaverse



Microsoft HoloLens 2

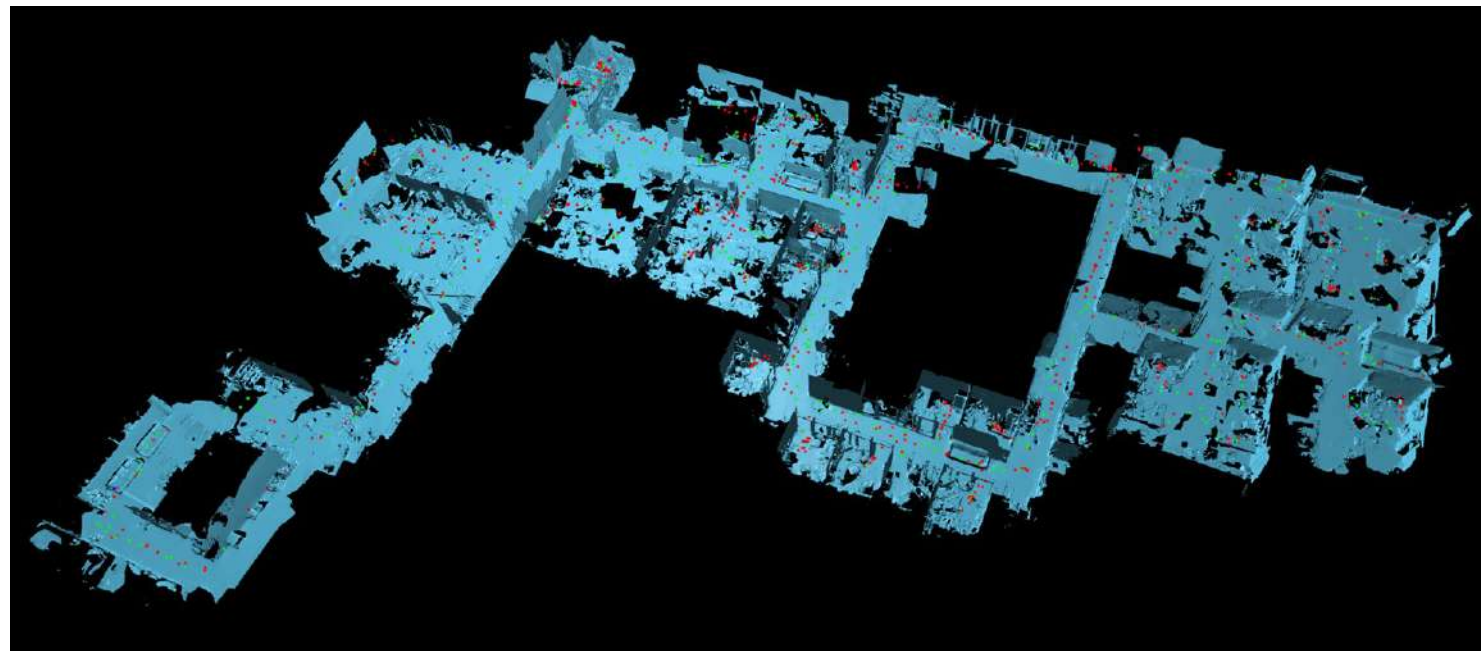
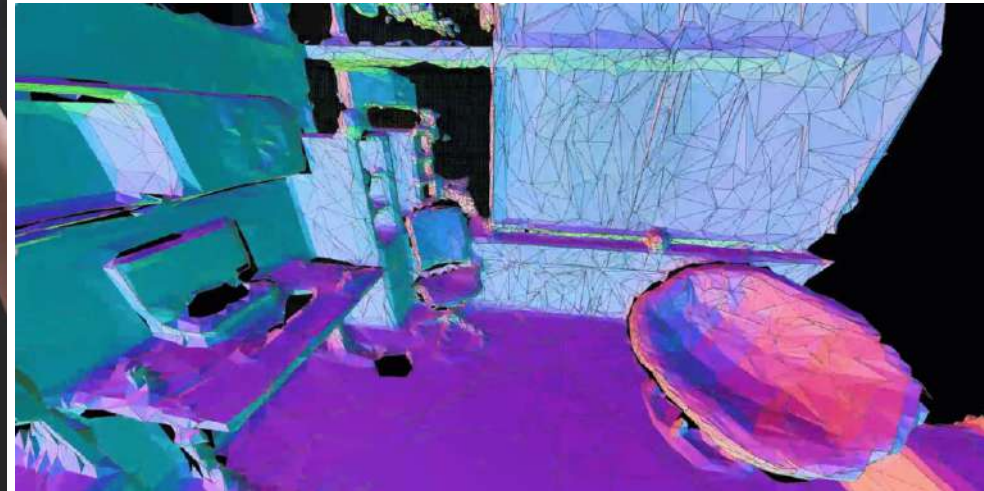
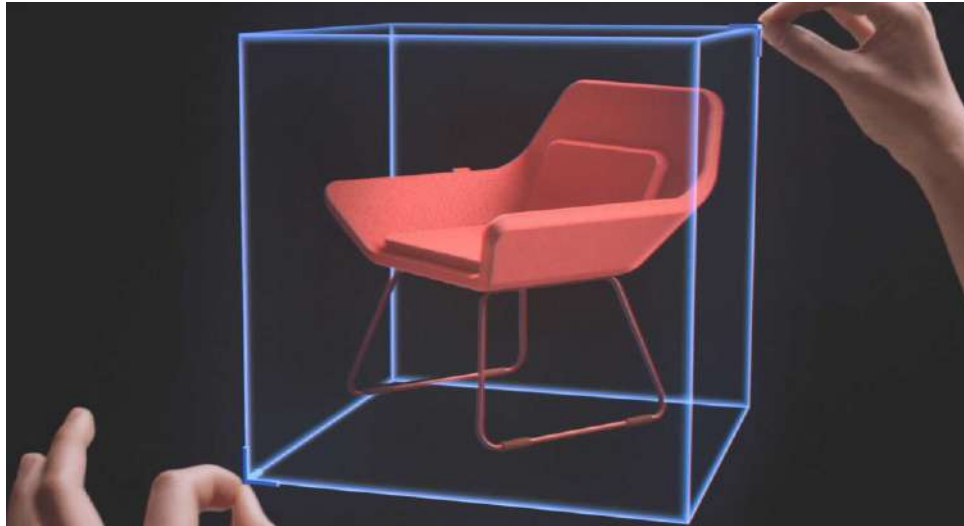
[Seitz, Szeliski]

# Augmented Reality - Metaverse



[Seitz, Szeliski]

# Augmented Reality



[Seitz, Szeliski]

## Phone-based AR

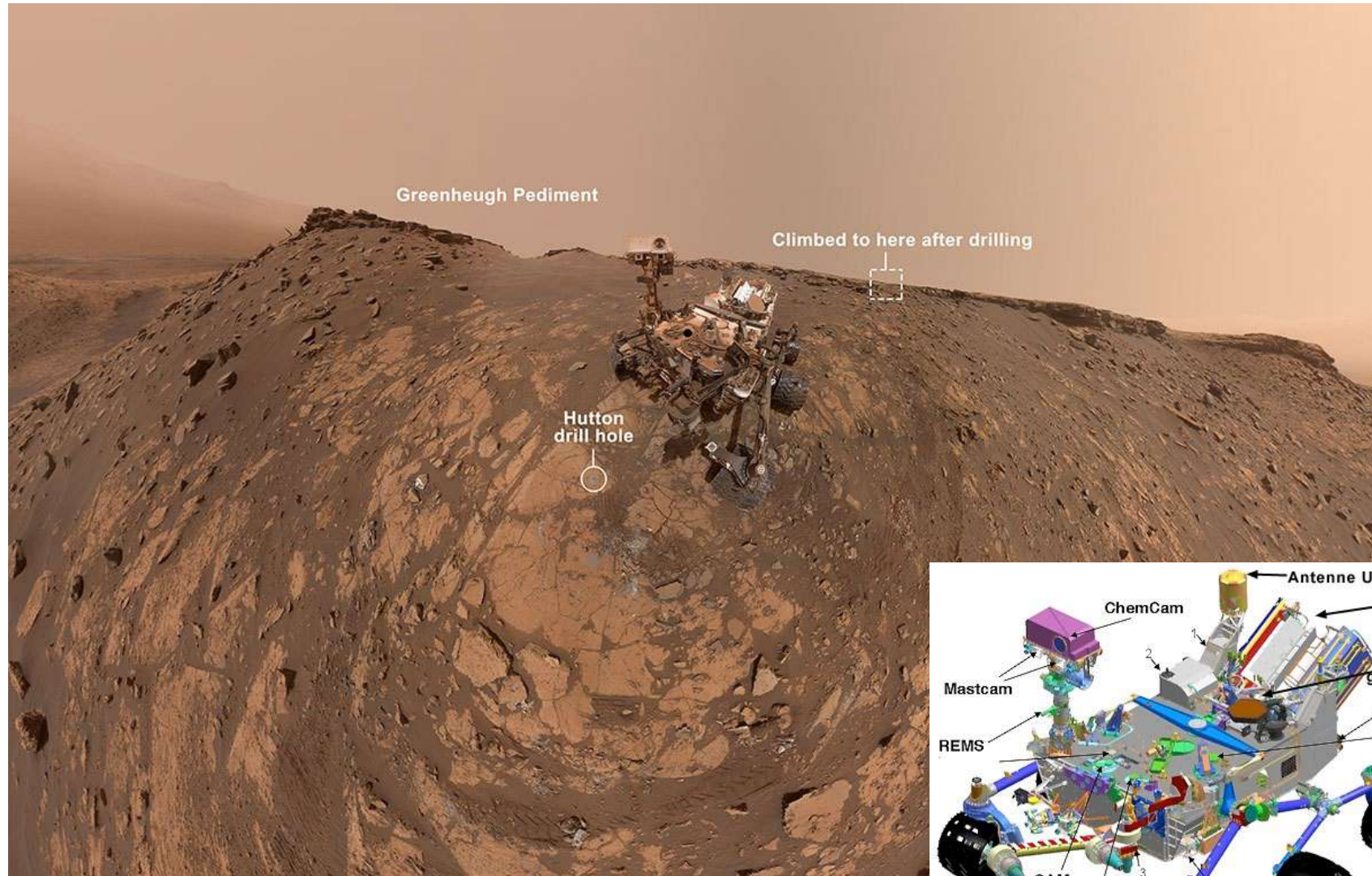


<http://www.youtube.com/watch?v=0Pj-jzy6ESE>

[Seitz, Szeliski]



# Robotics



NASA's Mars Curiosity Rover (self portrait)  
[https://en.wikipedia.org/wiki/Curiosity\\_\(rover\)](https://en.wikipedia.org/wiki/Curiosity_(rover))

[Seitz, Szeliski]

# Smart cars



manufacturer products | consumer products

## Our Vision. Your Safety.

rear looking camera | side looking camera | forward looking camera

**EyeQ** Vision on a Chip [read more](#)

**Vision Applications**  
Road, Vehicle, Pedestrian Protection and more [read more](#)

**AWS** Advance Warning System [read more](#)

**News**

- Mobileye Advanced Technologies Power Volvo Cars World First Collision Warning With Auto Brake System
- Volvo: New Collision Warning with Auto Brake Helps Prevent Rear-end

[all news](#)

**Events**

- Mobileye at Equip Auto, Paris, France
- Mobileye at SEMA, Las Vegas, NV

[read more](#)

## Mobileye

- Vision systems currently in high-end BMW, GM, Volvo models

[Seitz, Szeliski]

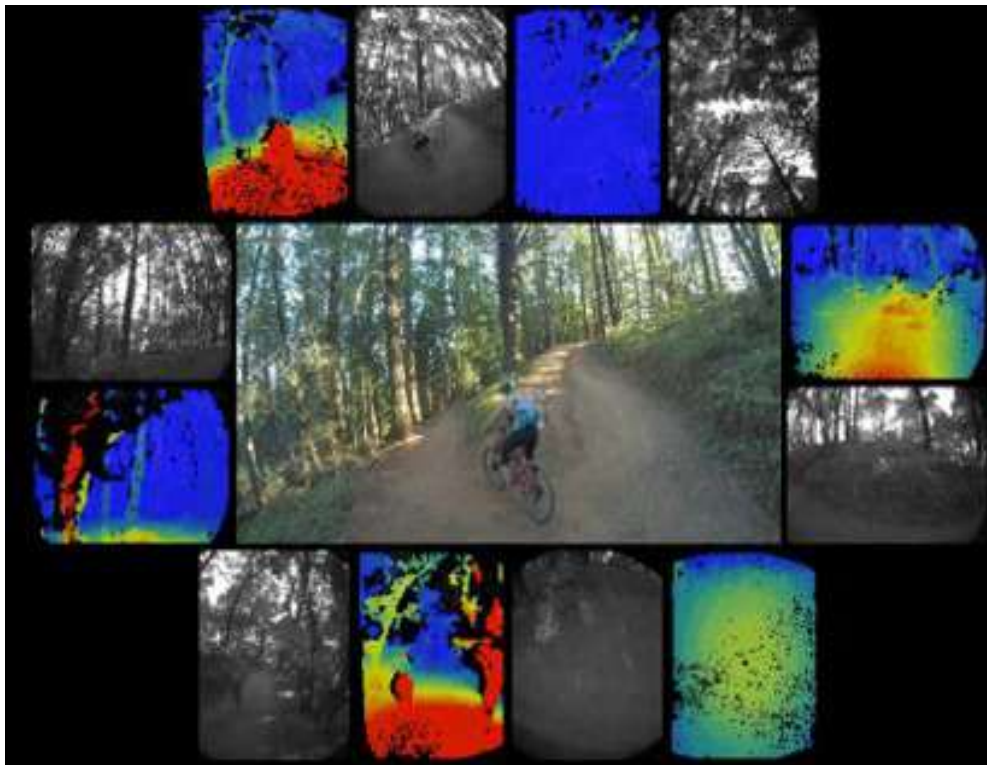


# Self-driving cars



<https://waymo.com/tech/>

## Drones



<https://www.skydio.com/>

[Seitz, Szeliski]



# Research: Yolo



<http://www.youtube.com/watch?v=MPU2HistivI>

[Seitz, Szeliski]

# Research: StyleGan



[http://www.youtube.com/watch?v=BIZg\\_PPuj\\_0](http://www.youtube.com/watch?v=BIZg_PPuj_0)

[Seitz, Szeliski]

# Today's Agenda

- Who are we?
- Introduction to Computer Vision
  - What is Computer Vision
  - How hard is Computer Vision
  - Why is Computer Vision so hard
  - How to organize Computer Vision
  - Why study Computer Vision
  - Applications
- What we do

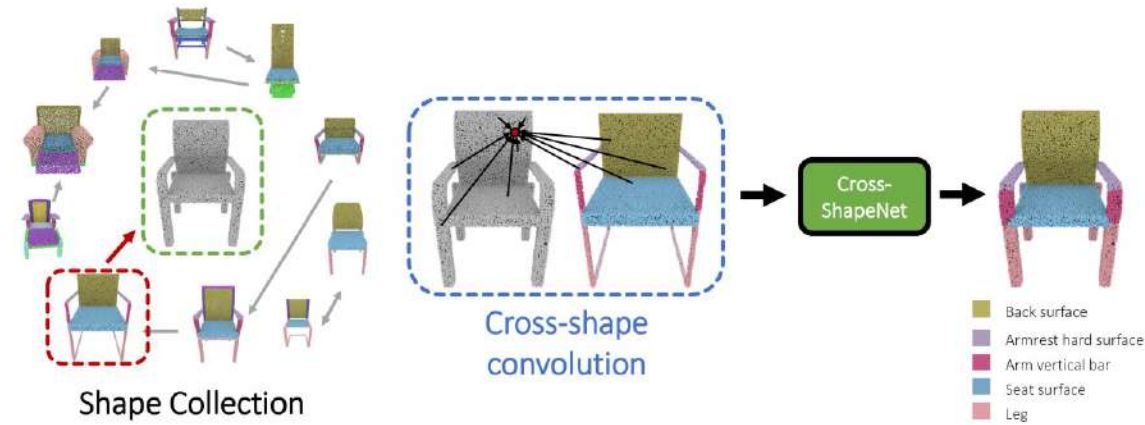


# 3D Shape Understanding

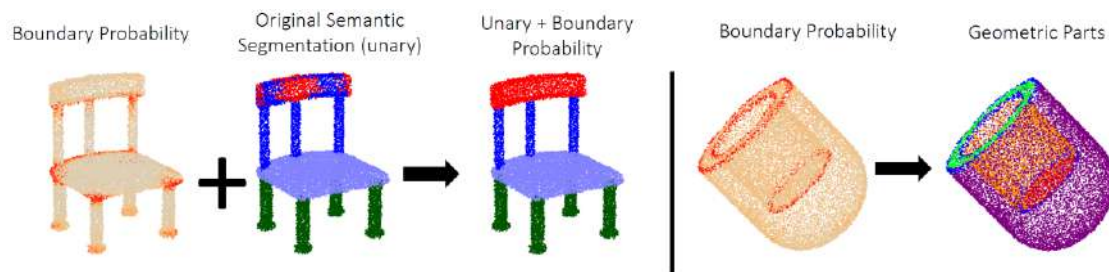
## 3D Building Semantic Understanding



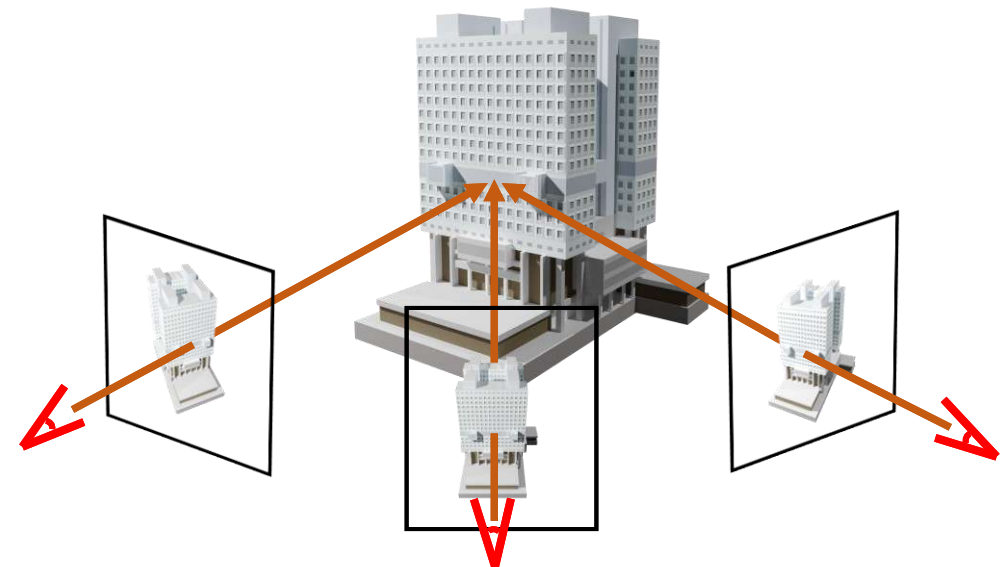
## Cross-shape semantic segmentation



## Geometric/Semantic Decomposition



## Neural 3D Reconstruction



# Texture Generation for 3D Data

## Single-View Guided Façade Synthesis



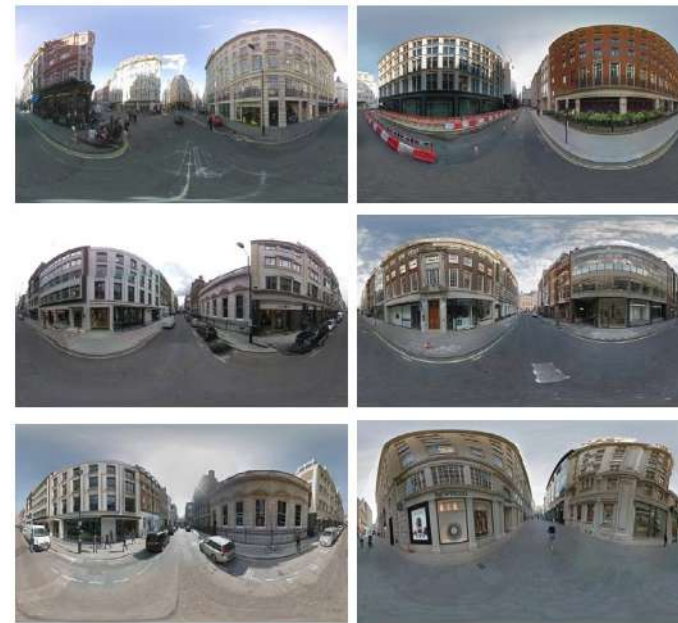
Reference Façade Images

3D Scene Renderings

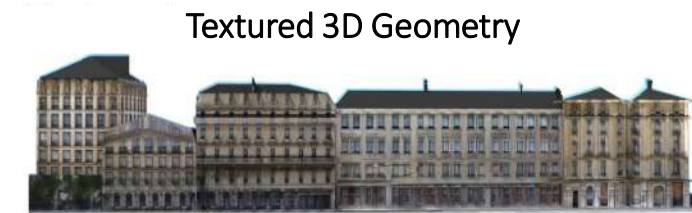


Interactive Texturing

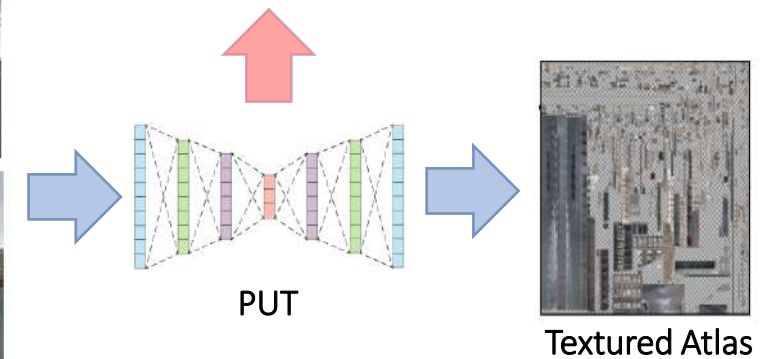
## Projective Urban Texturing



Street Level Panoramic Images



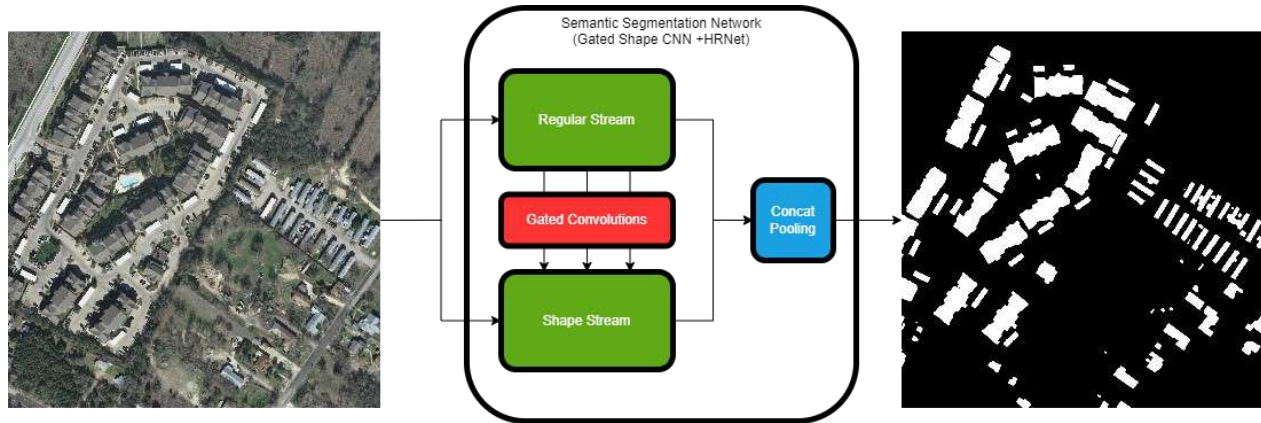
Textured 3D Geometry





# Urban Semantic Understanding from Remote Sensing Data Sources

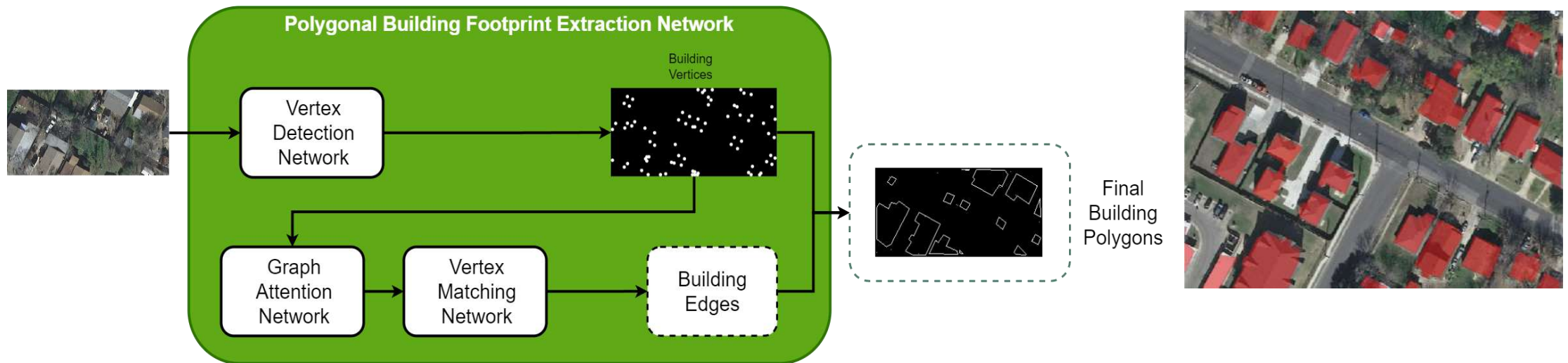
## Semantic Segmentation of Buildings



## Urban 3D Reconstruction



## Building Footprint Extraction



**MAI4CAREU**

Master programmes in Artificial  
Intelligence 4 Careers in Europe



# Thank you.







University of Cyprus – MSc Artificial Intelligence

# MAI644 – COMPUTER VISION

## Lecture 2: Fundamentals - Color

**Melinos Averkiou**

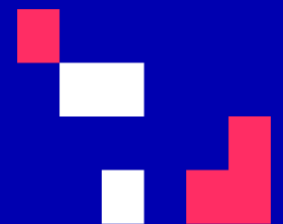
CYENS Centre of Excellence

University of Cyprus - Department of Computer Science

[m.averkiou@cyens.org.cy](mailto:m.averkiou@cyens.org.cy)



**CYENS**  
CENTRE OF EXCELLENCE



## Last time

- Course Overview
- Introduction to Computer Vision
  - What is Computer Vision
  - How hard is Computer Vision
  - Why is Computer Vision so hard
  - How to organize Computer Vision
  - Why study Computer Vision
  - Applications

# Today's Agenda - Overview of Color

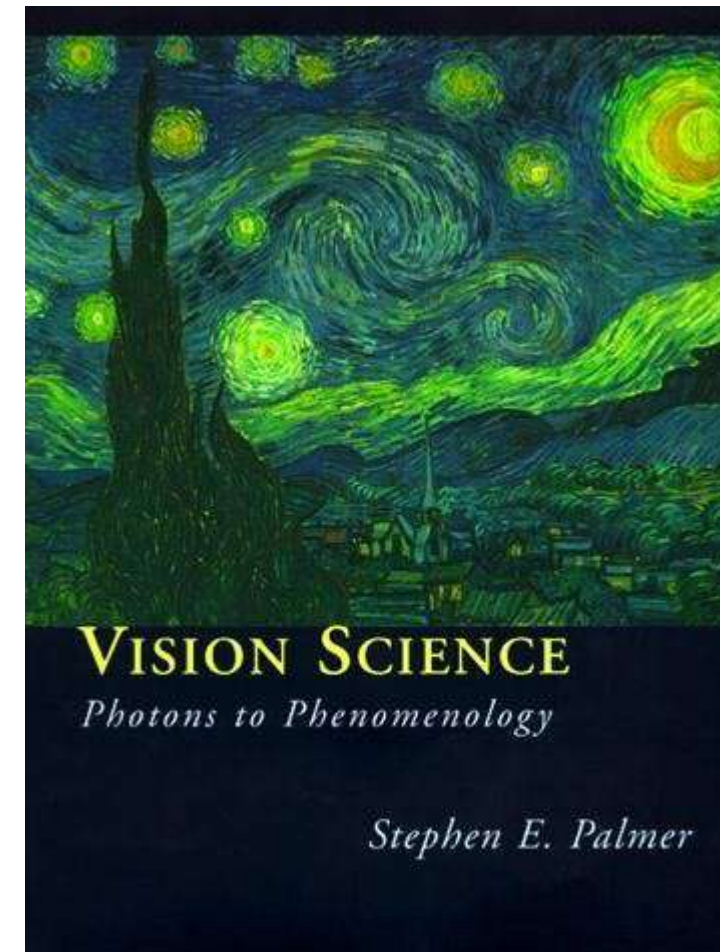
- Physics of color
- Human encoding of color
- Color spaces

[Slides based on material by Niebles 2017]



# What is color?

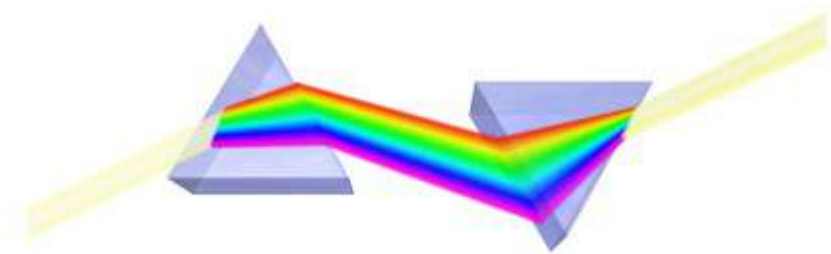
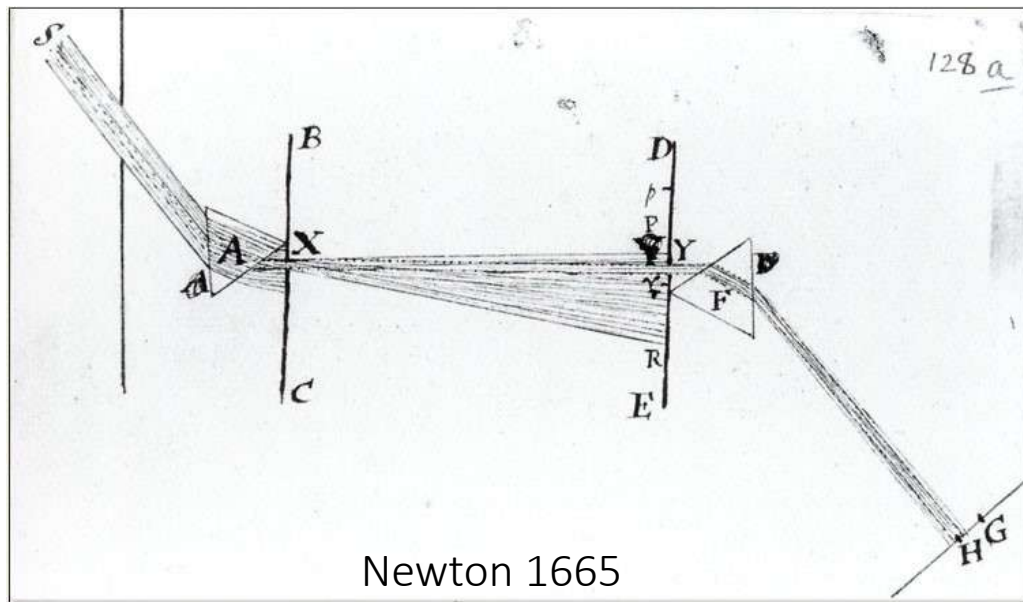
- The result of interaction between physical **light** in the environment and our **visual system**.
- A **psychological property** of our visual experiences when we look at objects and lights, **not a physical property** of those objects or lights.



[Slide credit: Lana Lazebnik]

# Color and light

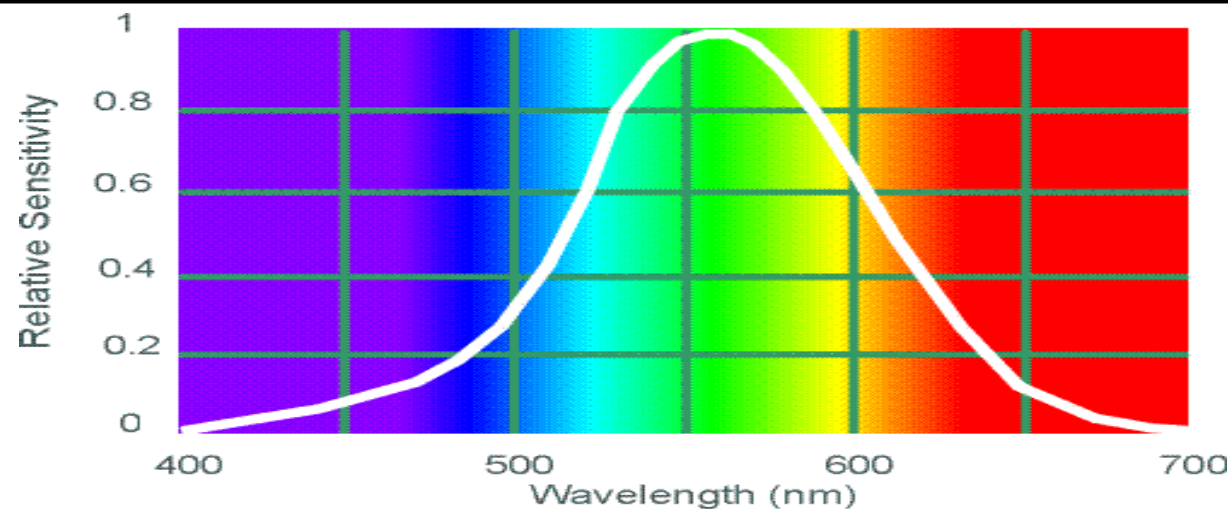
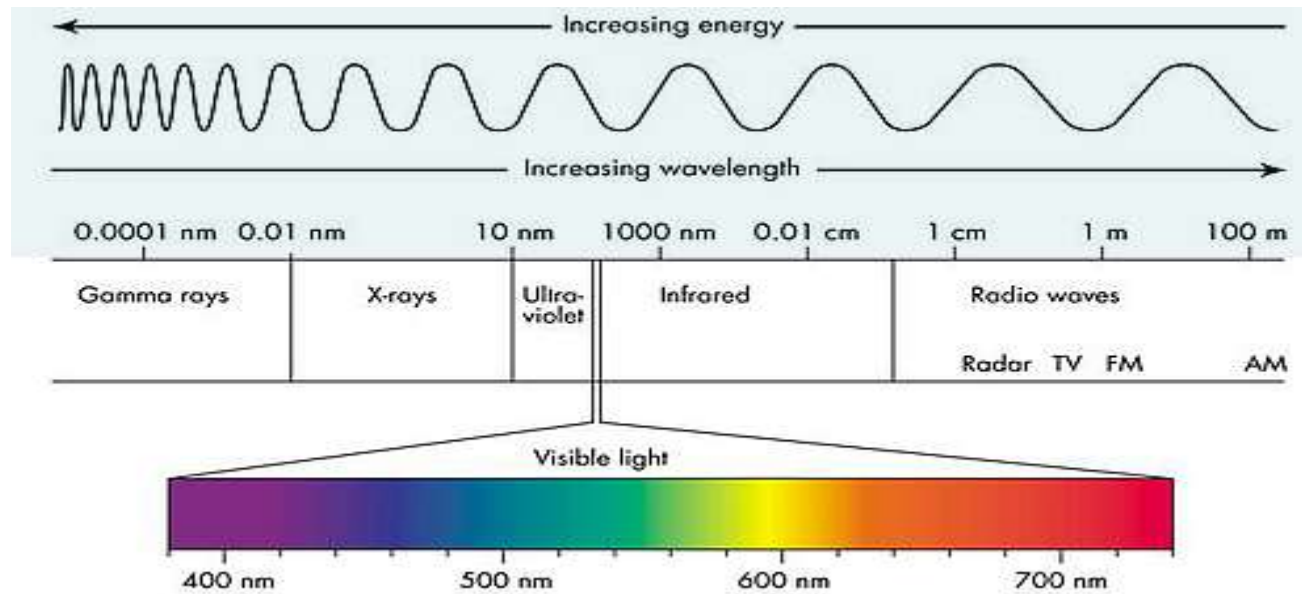
White light: composed of almost equal energy in all wavelengths of the visible spectrum



Try it:

<https://micro.magnet.fsu.edu/primer/java/scienceopticsu/newton/>

# Electromagnetic Spectrum

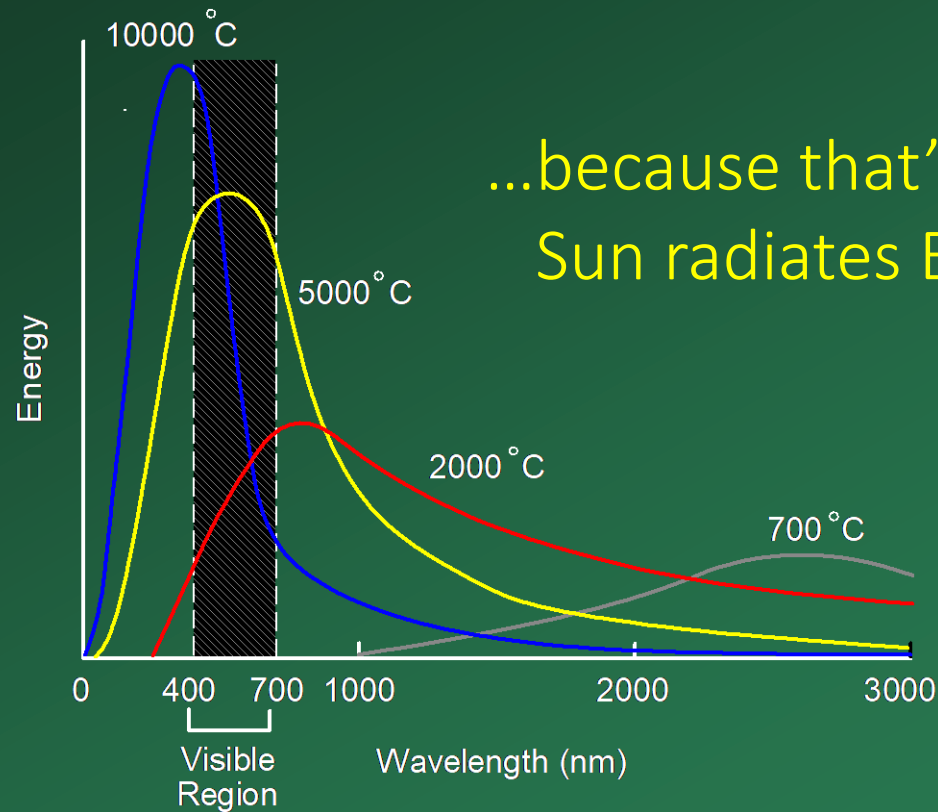


Human Luminance Sensitivity Function

<http://www.yorku.ca/eye/photopik.htm>

# Visible Light

Why do we see light of these wavelengths?

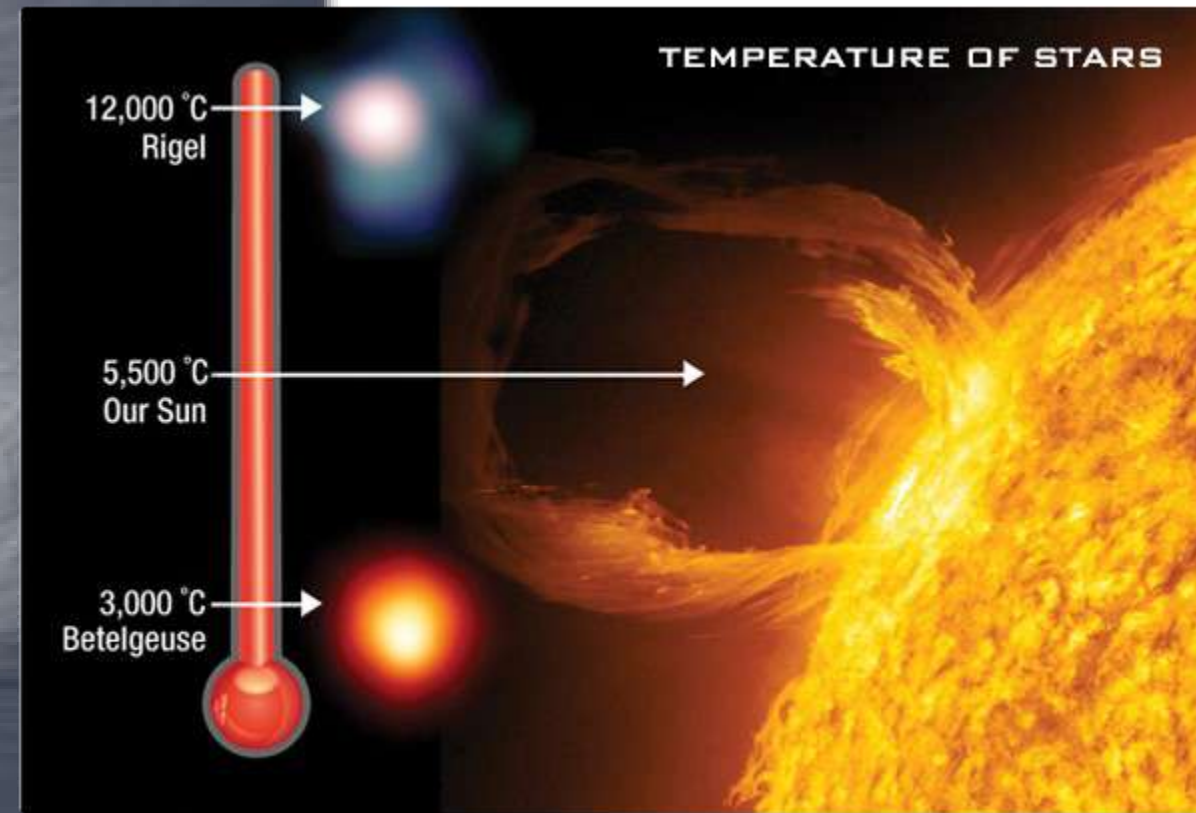
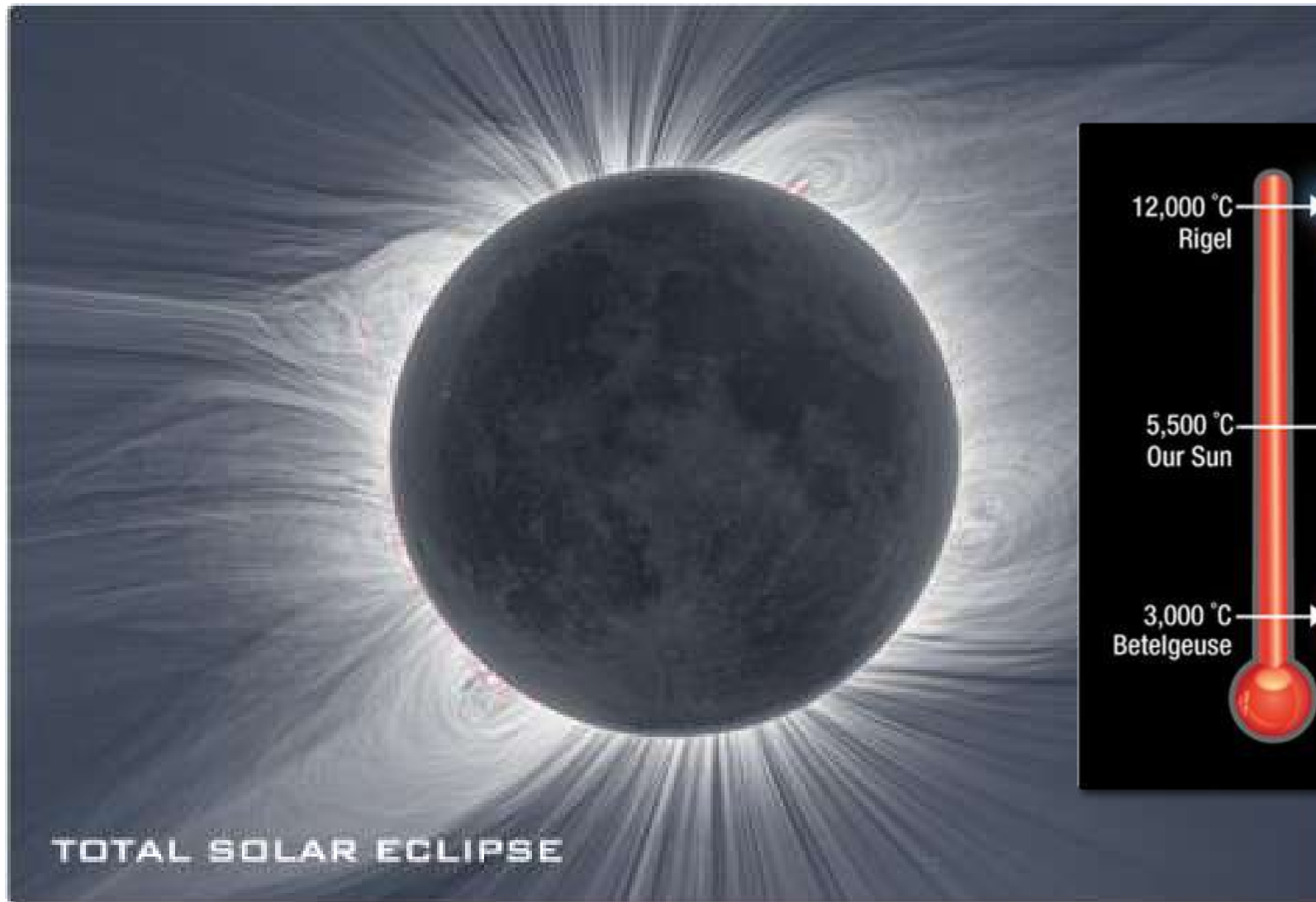


...because that's where the Sun radiates EM energy

© Stephen E. Palmer, 2002



Sun's temperature makes it emit yellow light more than any other color

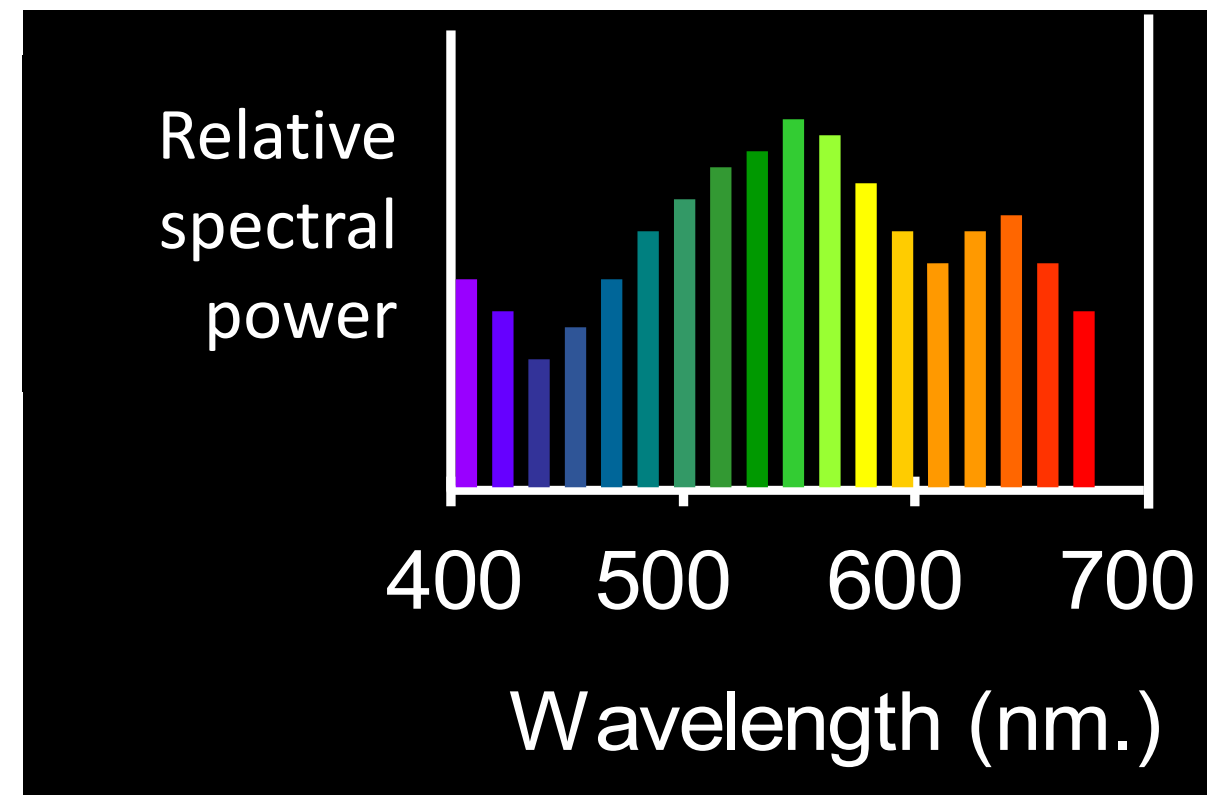


[https://science.nasa.gov/ems/09\\_visiblelight](https://science.nasa.gov/ems/09_visiblelight)

# The Physics of Light

Any source of visible light can be completely described physically by its **spectrum**:

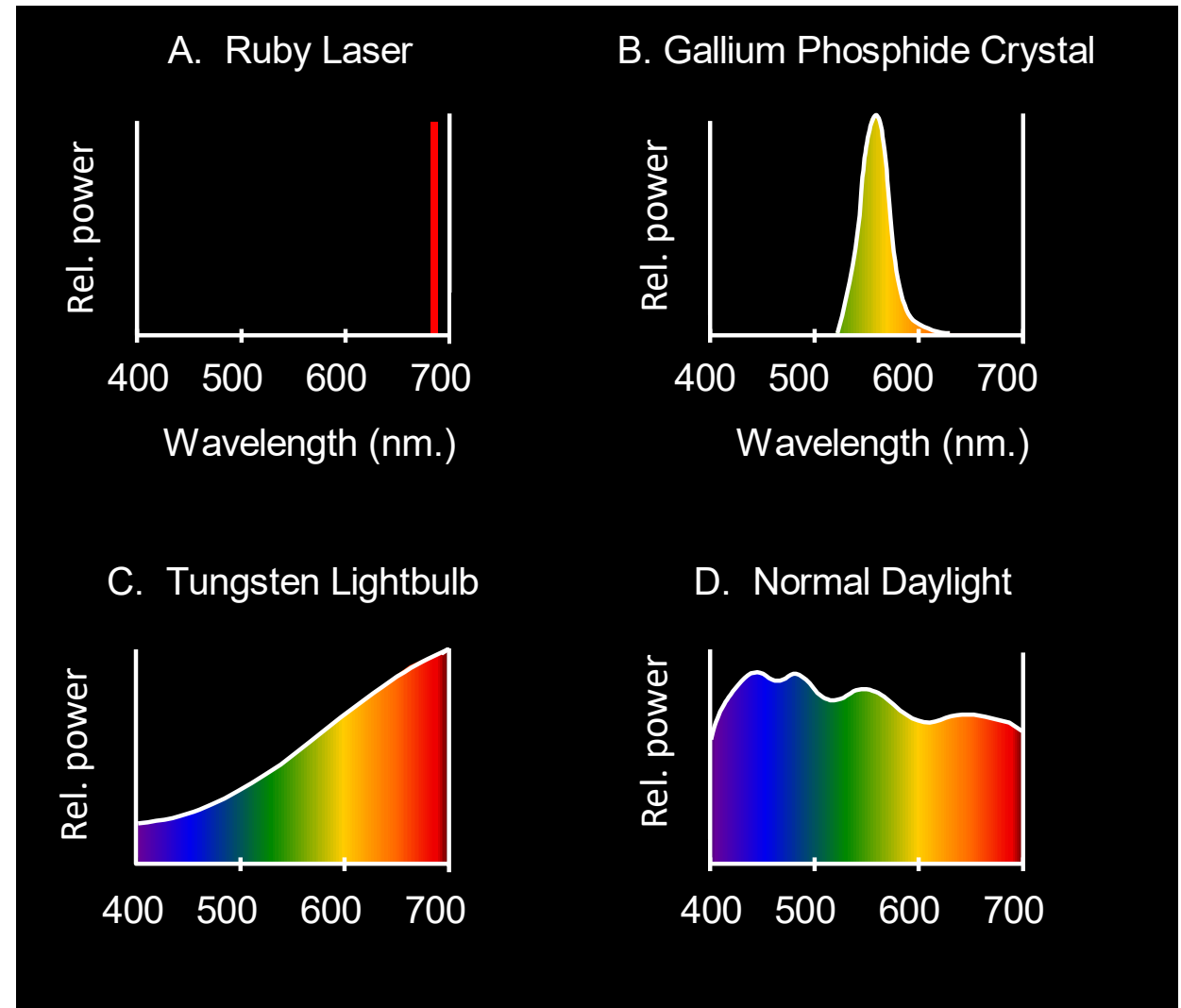
the amount of energy emitted (per time unit) at each wavelength 400 - 700 nm.



© Stephen E. Palmer, 2002

# The Physics of Light

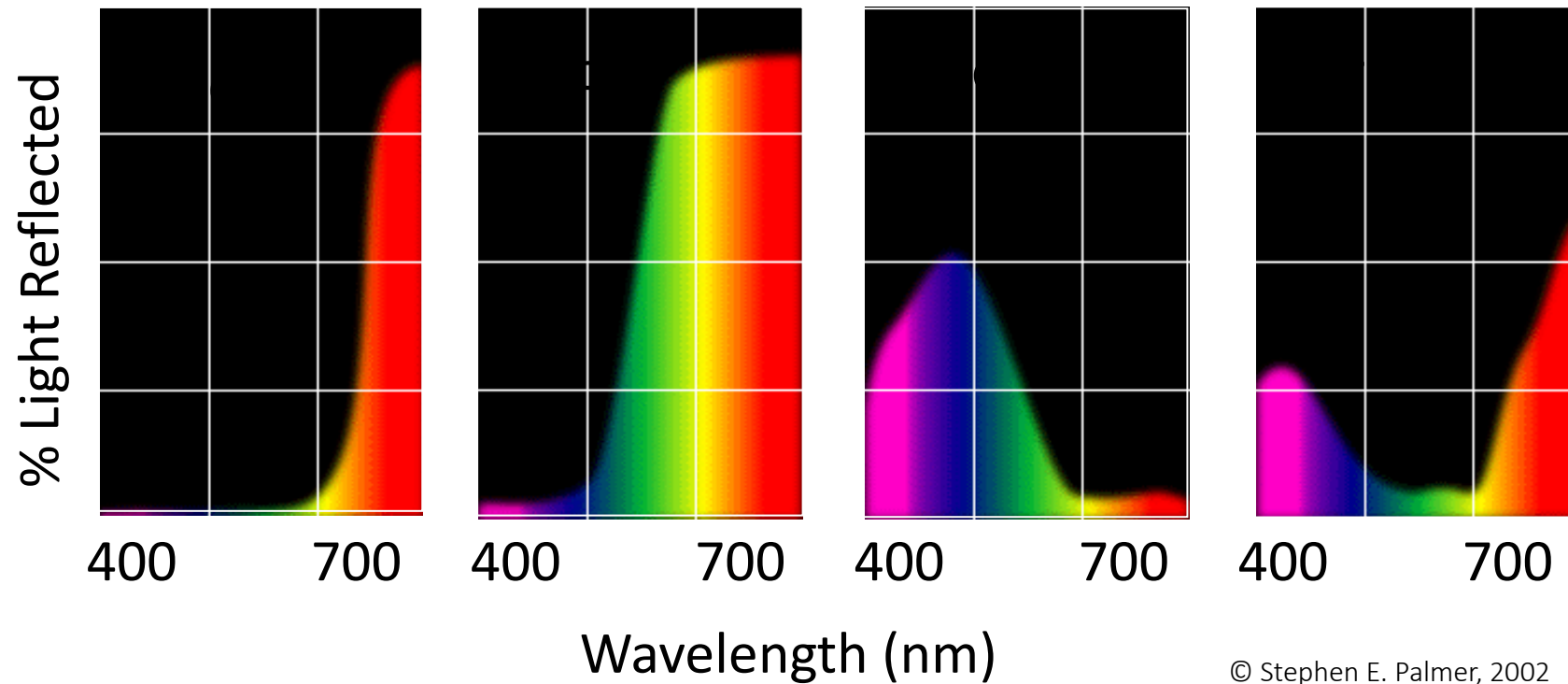
Some examples of the spectra of light sources



© Stephen E. Palmer, 2002

# The Physics of Light

Some examples of the **reflectance** spectra of surfaces

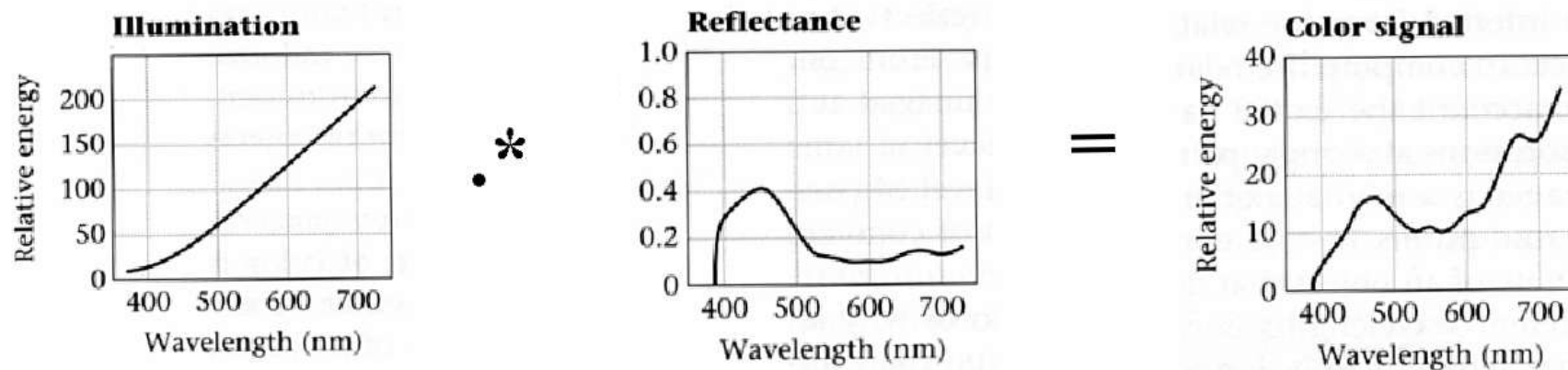


© Stephen E. Palmer, 2002



# Interaction of light and surfaces

Reflected color is the result of interaction of **light** source spectrum with **surface** reflectance



From Foundation of Vision by Brian Wandell, Sinauer Associates, 1995

# Interaction of light and surfaces

What is the observed color of any surface under monochromatic light?



[Olafur Eliasson, \*Room for one color\*](#)

[Slide by S. Lazebnik]

## Interaction of light and surfaces



[See here](#)

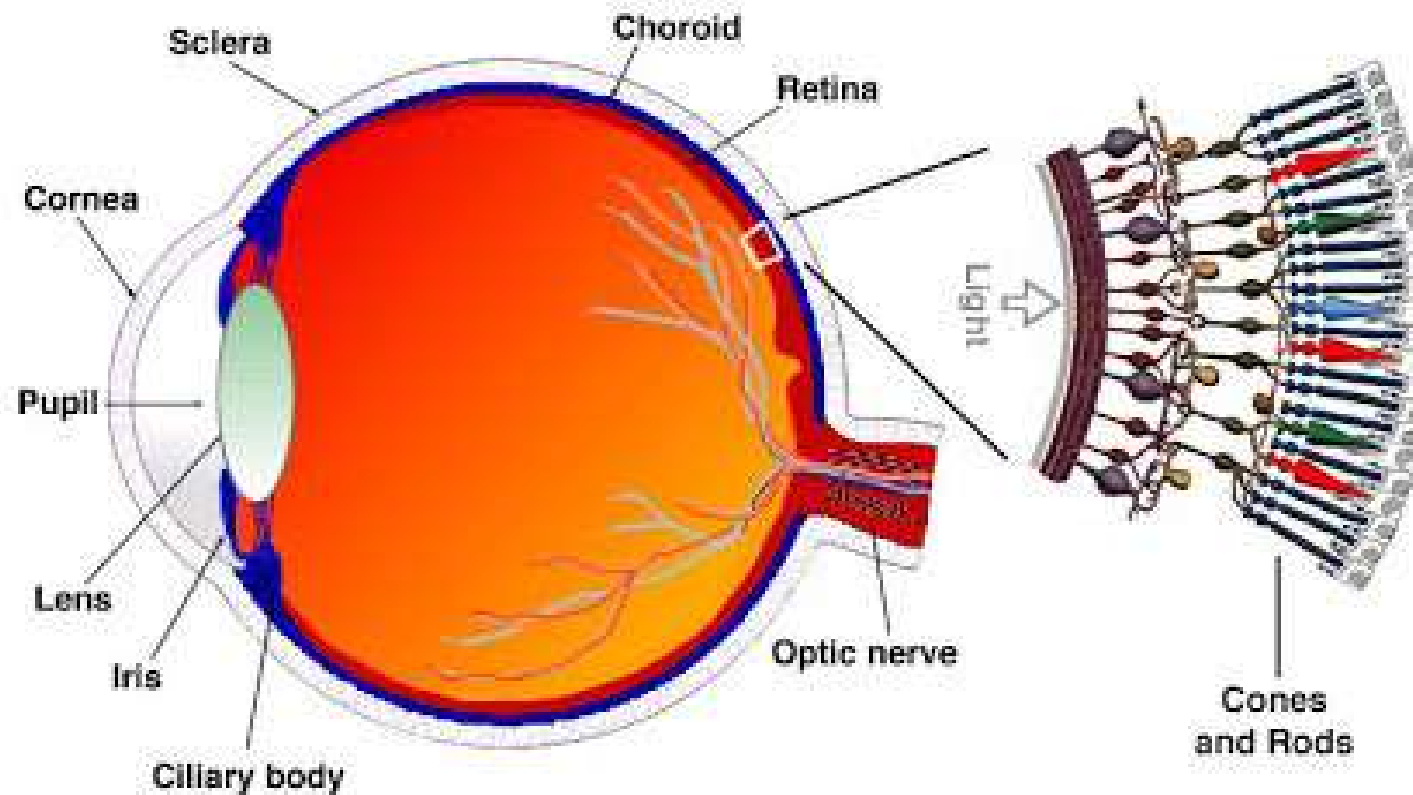
# Overview of Color

- Physics of color
- Human encoding of color
- Color spaces





# Two types of light-sensitive receptors



## Cones

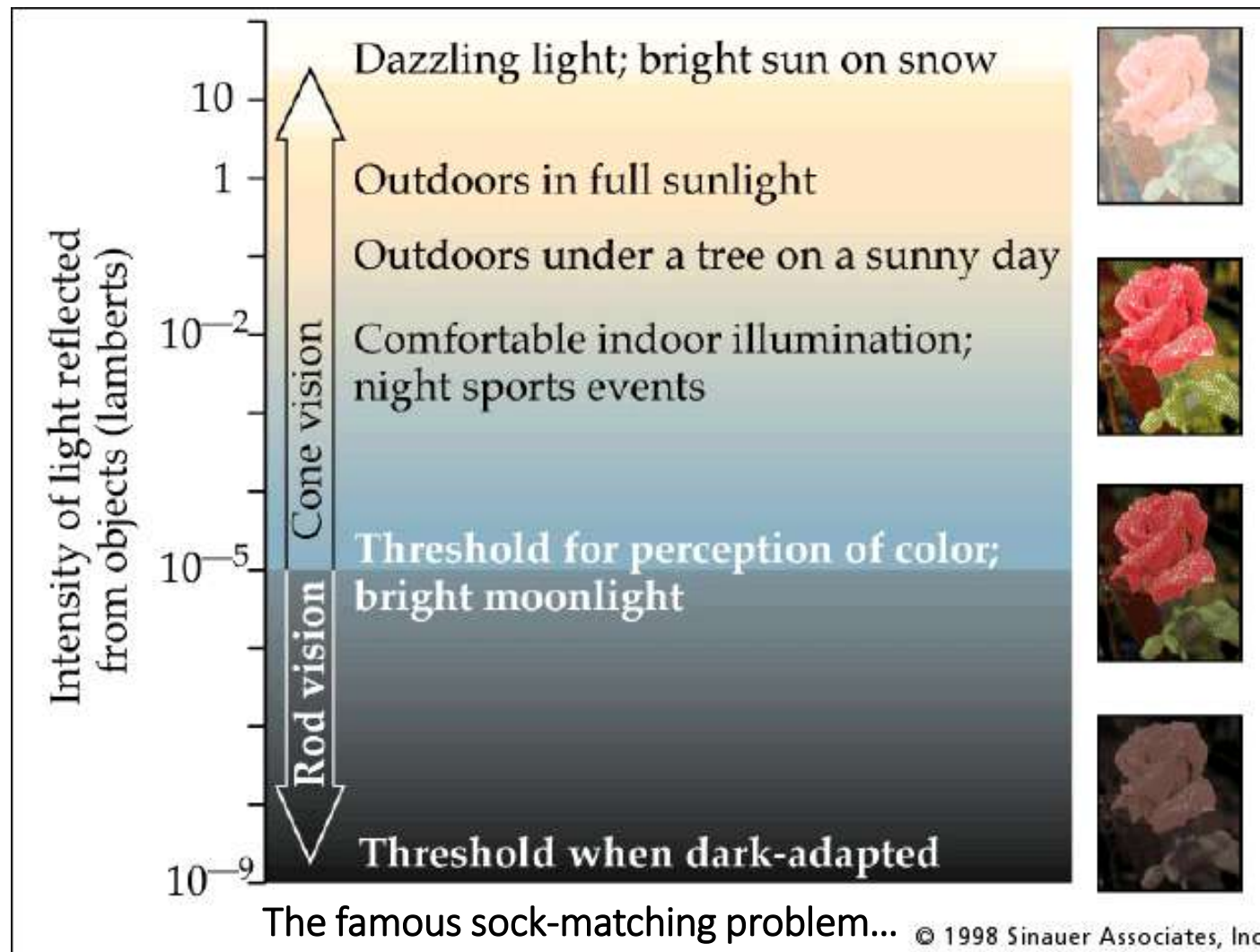
- cone-shaped
- less sensitive
- operate in high light
- color vision

## Rods

- rod-shaped
- highly sensitive
- operate at night
- gray-scale vision

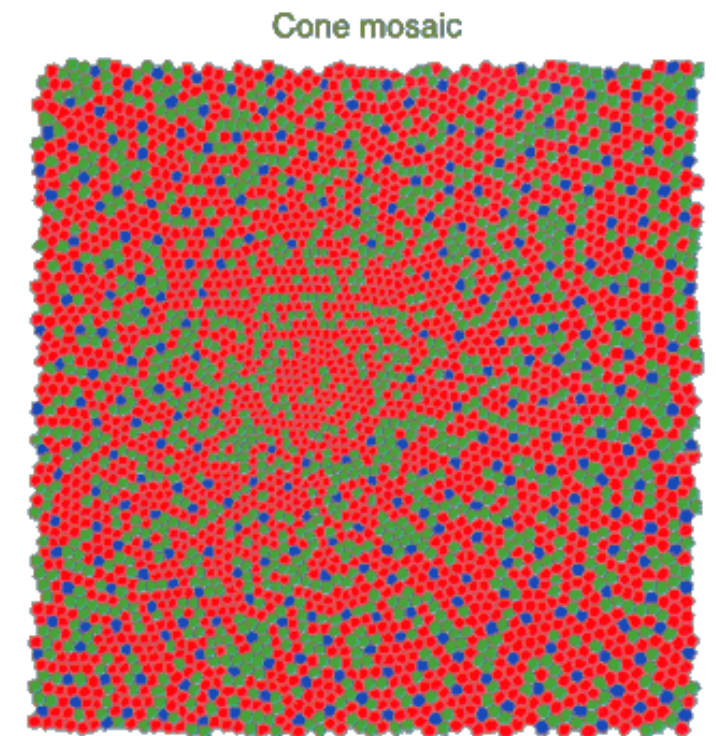
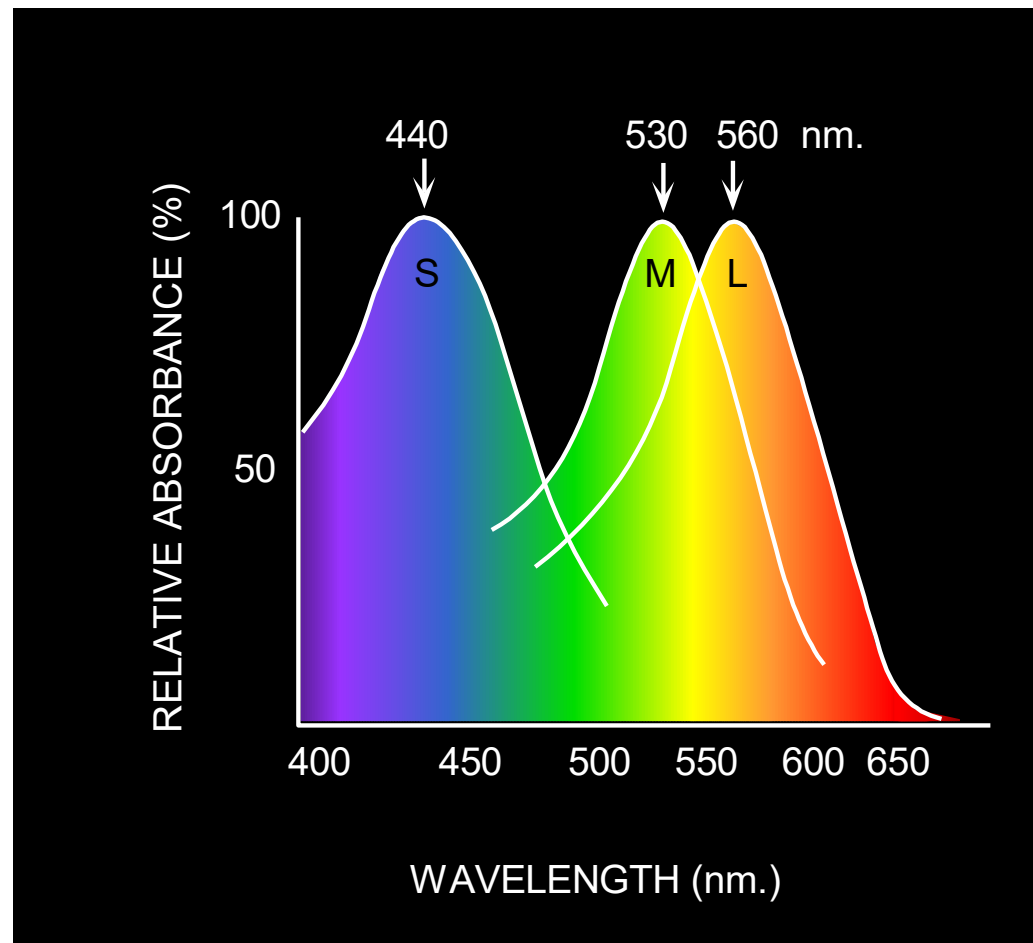
<http://www.blueconemonochromacy.org/how-the-eye-functions/>

# Rod / Cone sensitivity



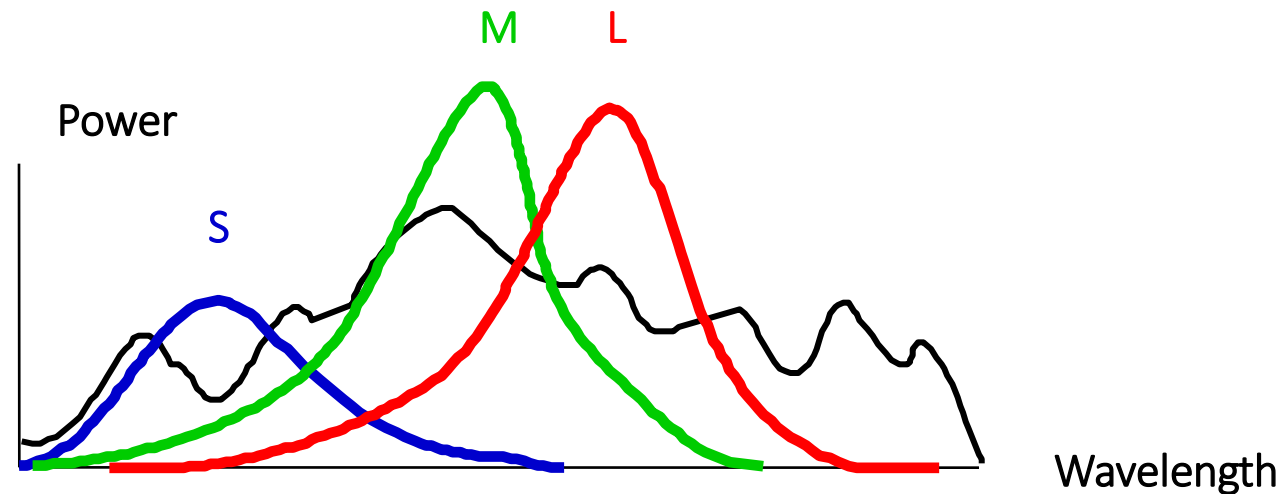
# Physiology of Color Vision

Three kinds of cones: Short (S), Medium (M), Long (L)



© Stephen E. Palmer, 2002

# Color perception



## Rods and cones act as filters on the spectrum

- To get the output of a filter, multiply its response curve by the spectrum, integrate over all wavelengths
  - Each cone yields **one** number
- Q: How can we represent an entire spectrum with 3 numbers?
- A: We can't! Most of the information is lost.
  - As a result, two different spectra may appear indistinguishable - such spectra are known as **metamers** ([see demo](#))

[Slide by Steve Seitz]

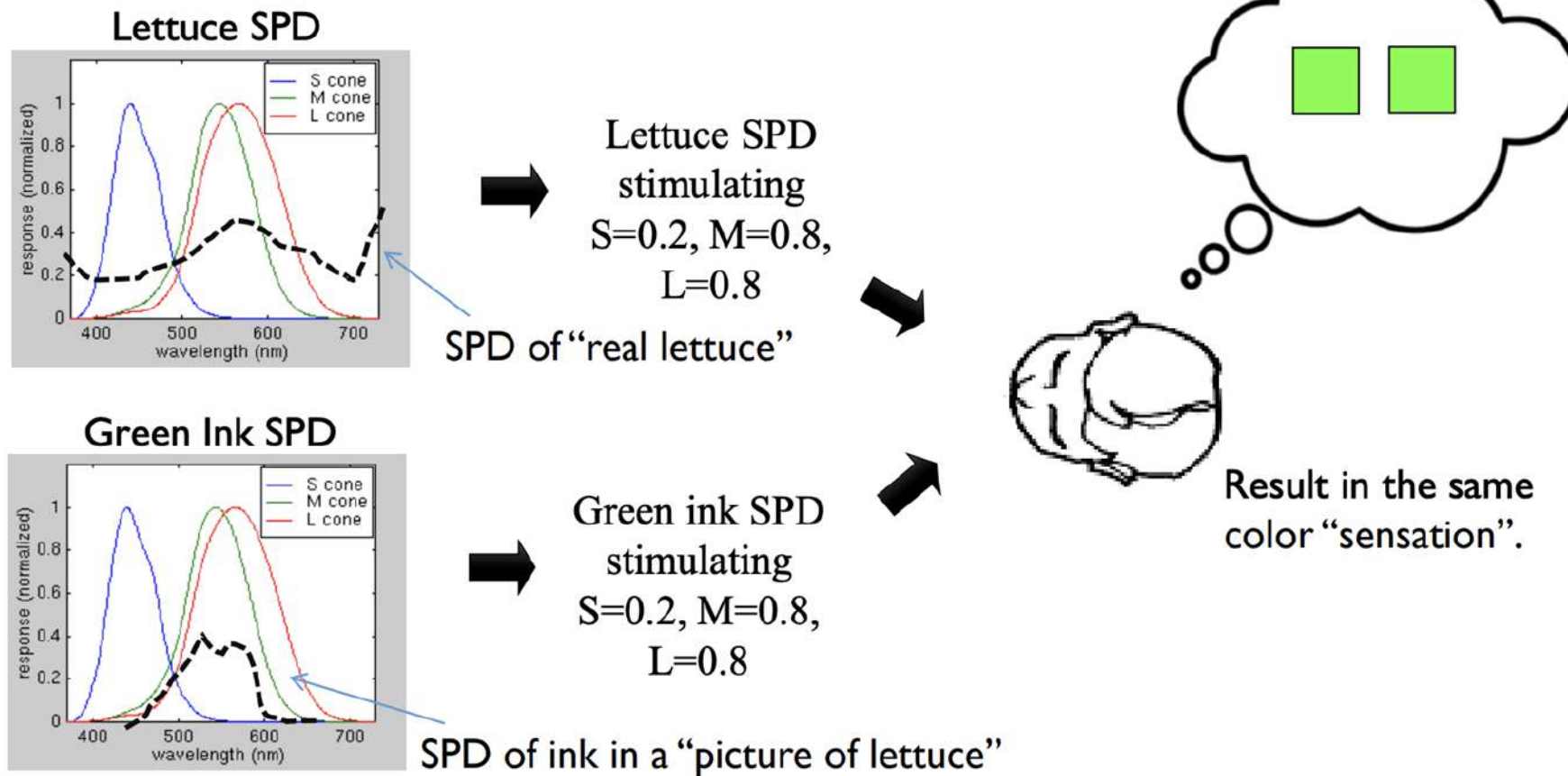


# We need metamers!



[Seitz]

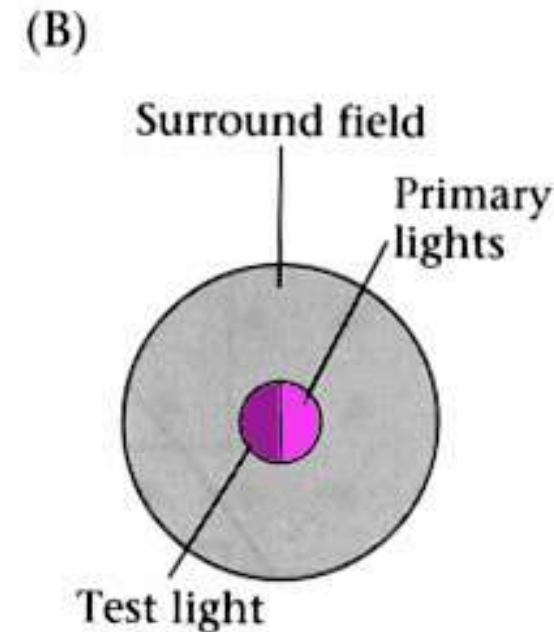
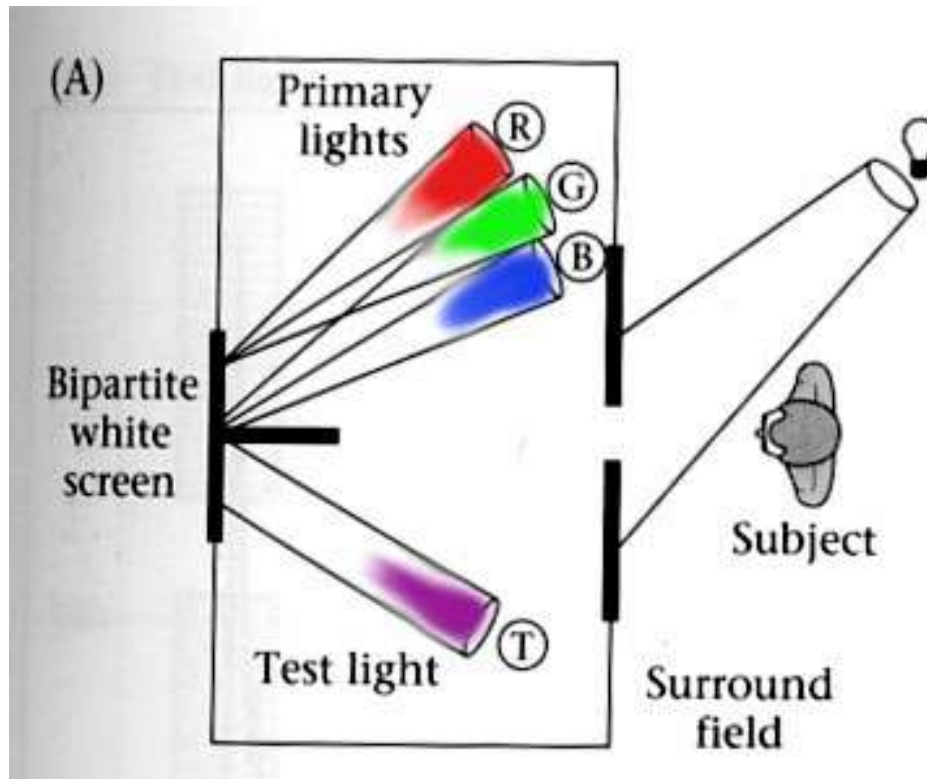
# We need metamers!



[Seitz]

# Standardizing color experience

- We would like to understand which spectra produce the same color sensation in people under similar viewing conditions
- Color matching experiments



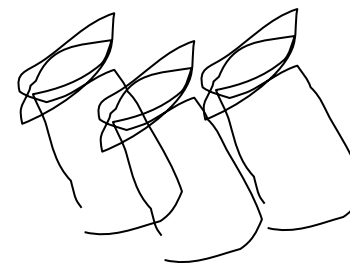
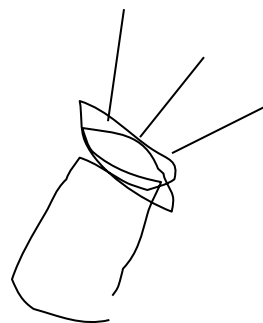
Foundations of Vision, by Brian Wandell, Sinauer Assoc., 1995

# Trichromacy

- In color matching experiments, most people can match any given light with three primaries
  - Primaries must be *independent*
- For the same light and same primaries, most people select the same weights
  - Exception: color blindness
- Trichromatic color theory
  - Three numbers seem to be sufficient for encoding color
  - Dates back to 18<sup>th</sup> century (Thomas Young)

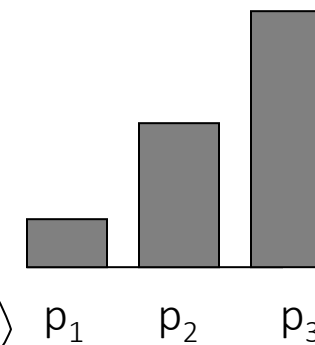
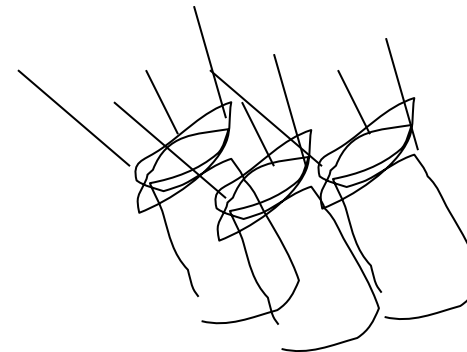
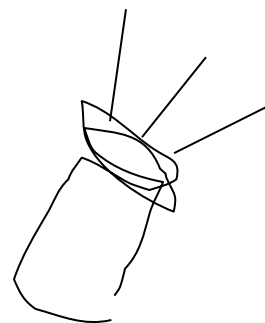
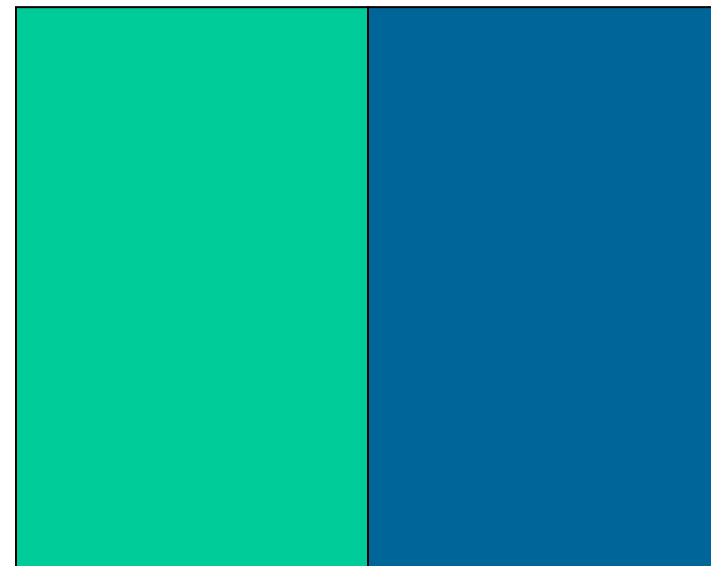


# Color matching experiment 1



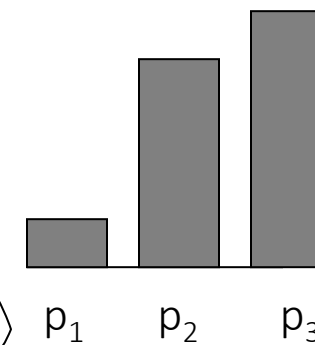
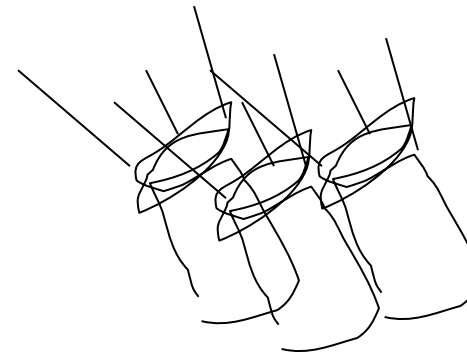
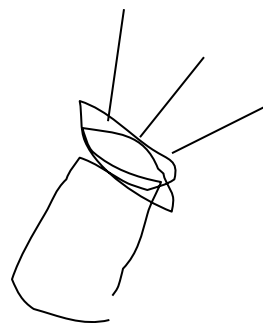
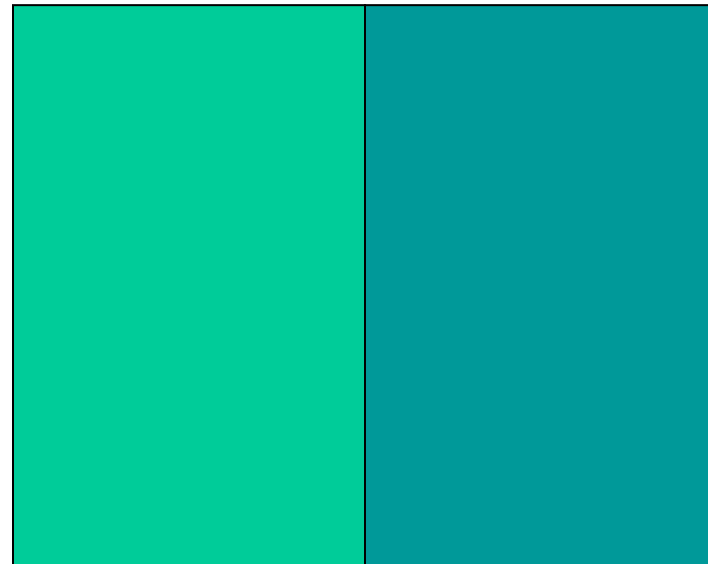
Source: W. Freeman

# Color matching experiment 1



Source: W. Freeman

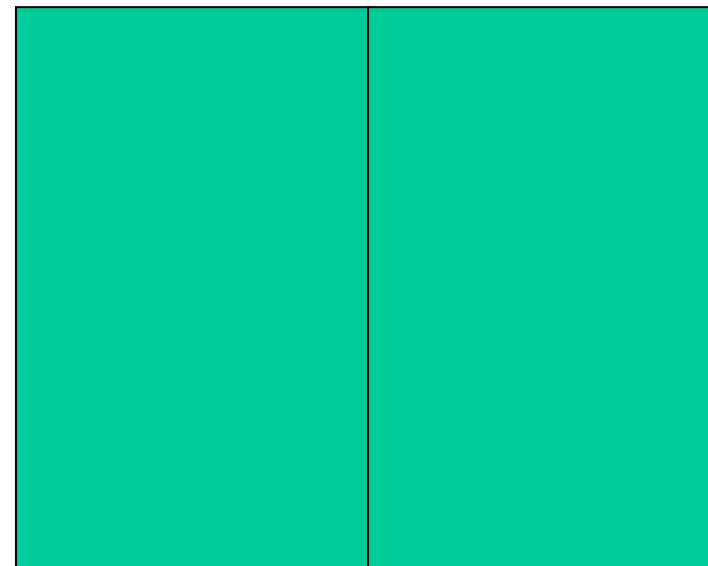
# Color matching experiment 1



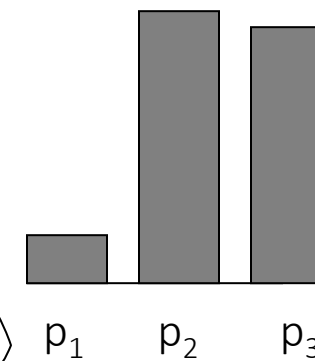
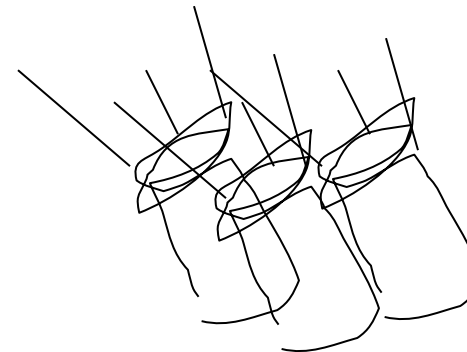
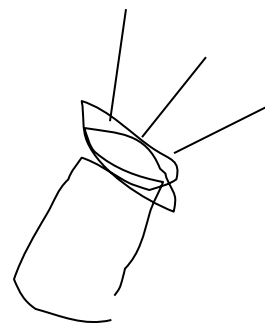
Source: W. Freeman



# Color matching experiment 1



The primary color amounts needed for a match

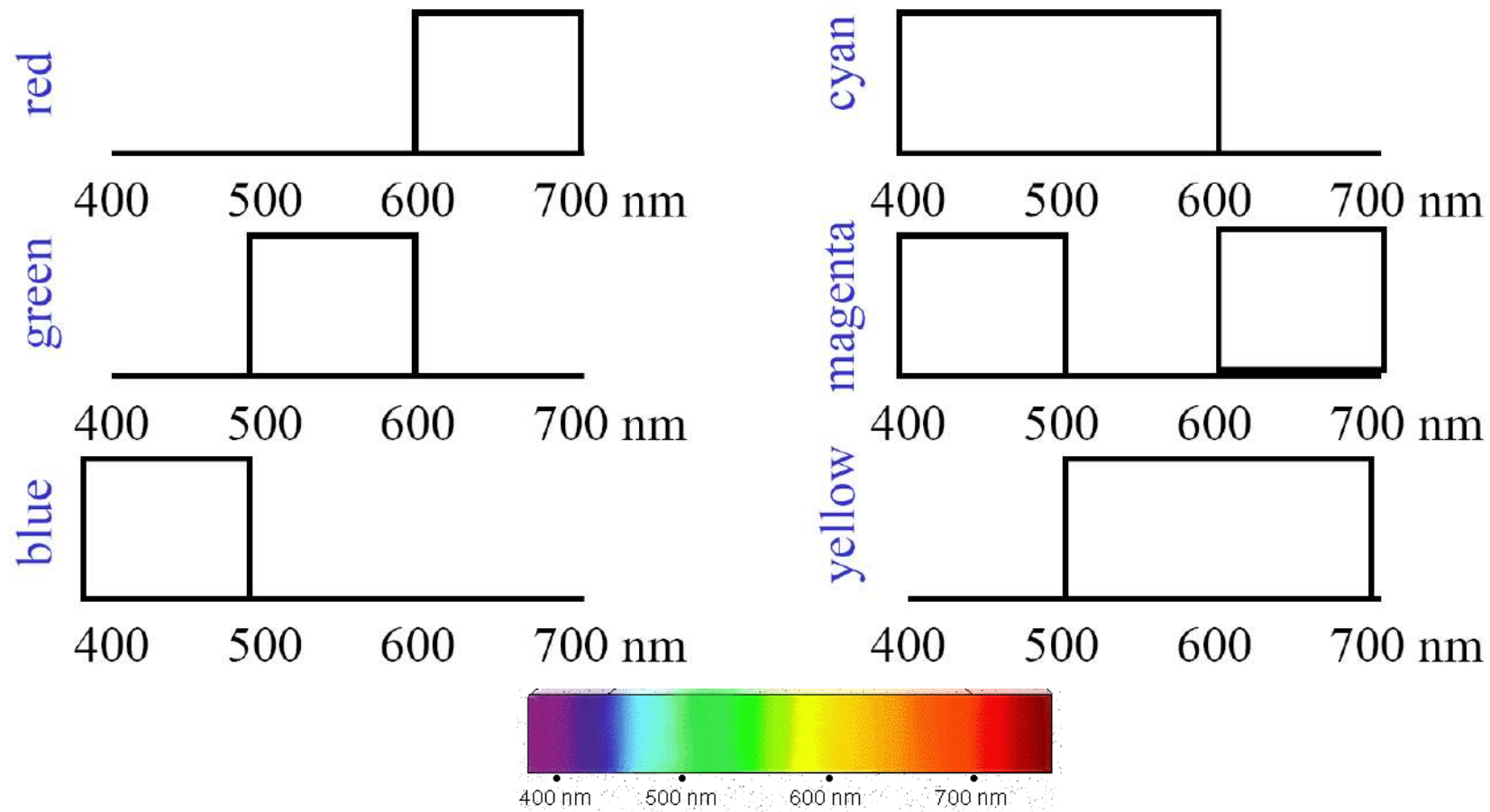


Source: W. Freeman



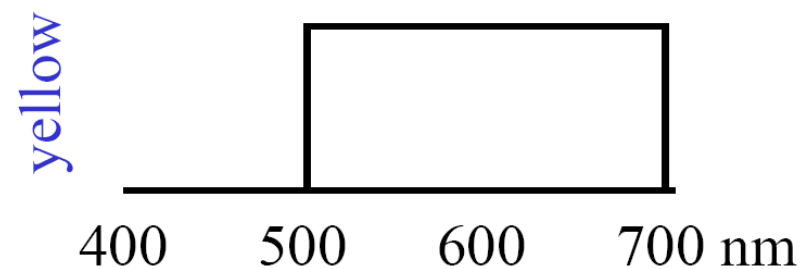
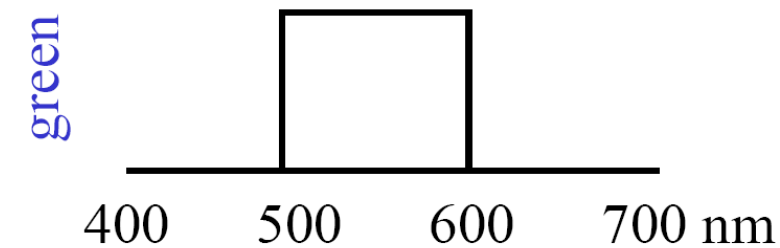
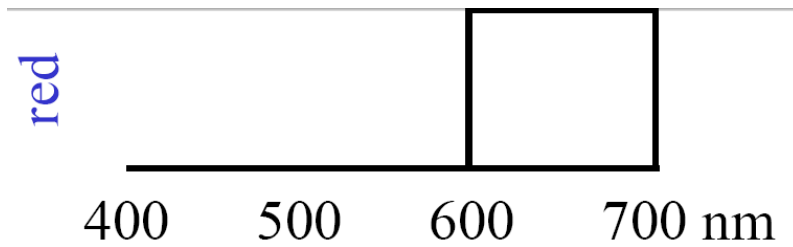


# Color mixing

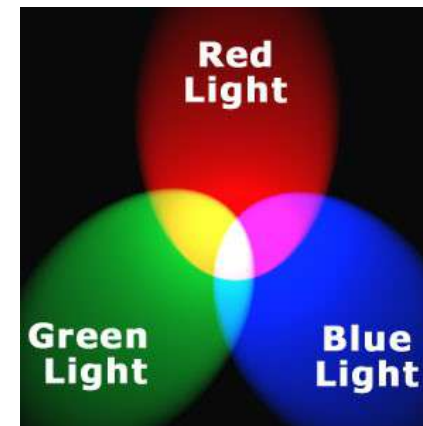


Source: W. Freeman

# Additive color mixing



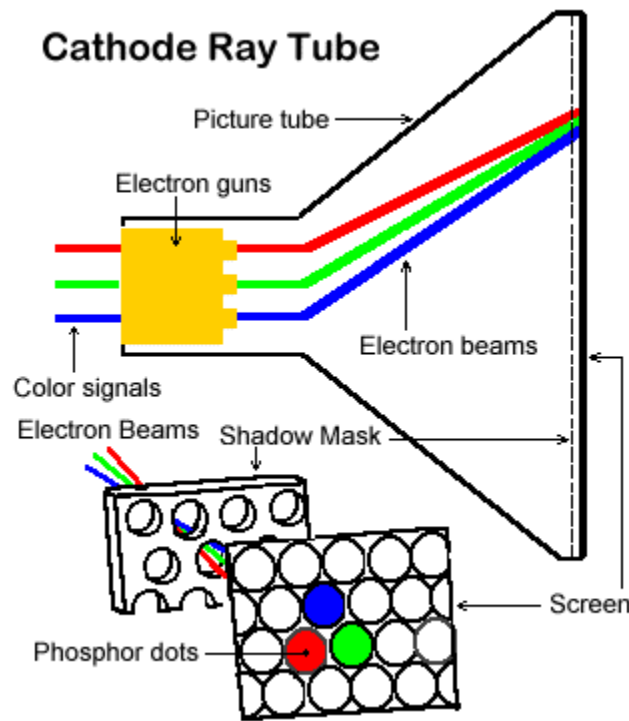
Colors combine by *adding* color spectra



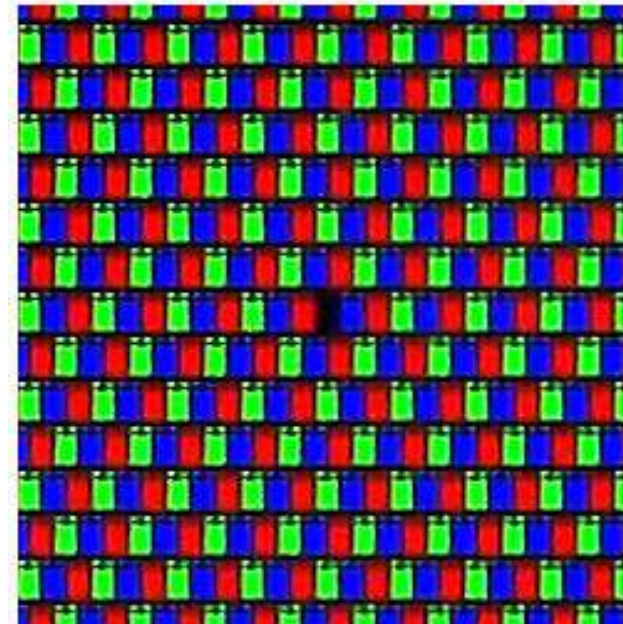
Light *adds* to existing black.

Source: W. Freeman

# Examples of additive color systems

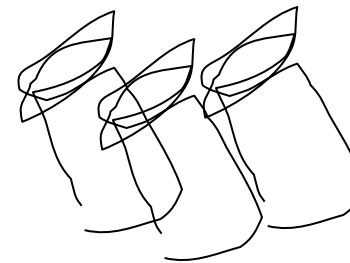
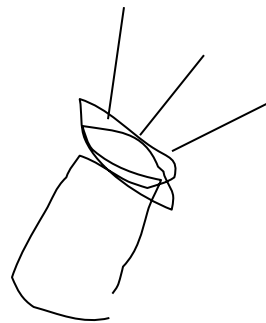


CRT phosphors



Multiple projectors

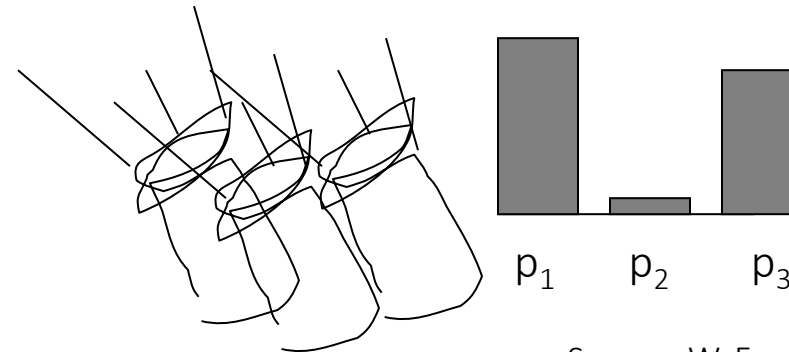
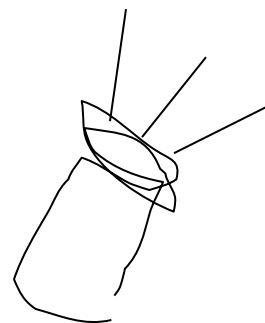
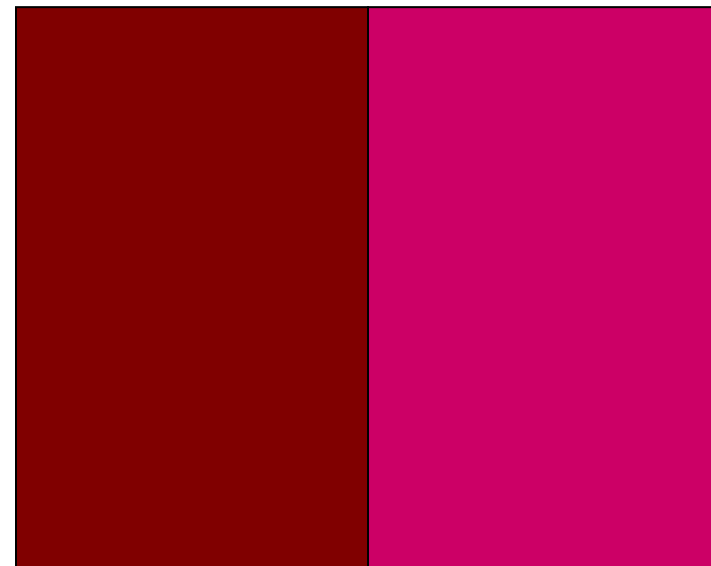
## Color matching experiment 2



Source: W. Freeman

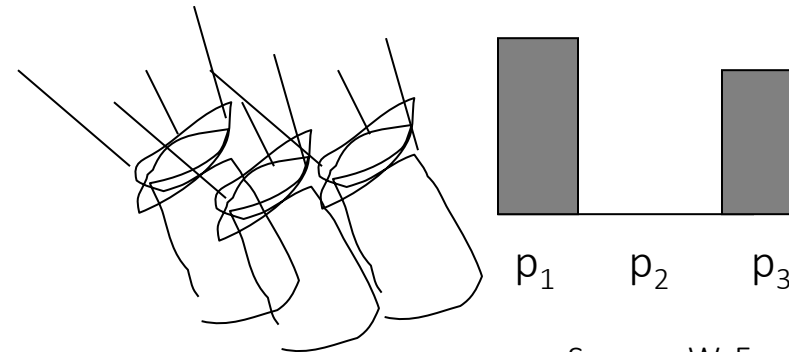
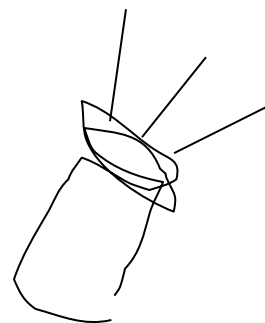
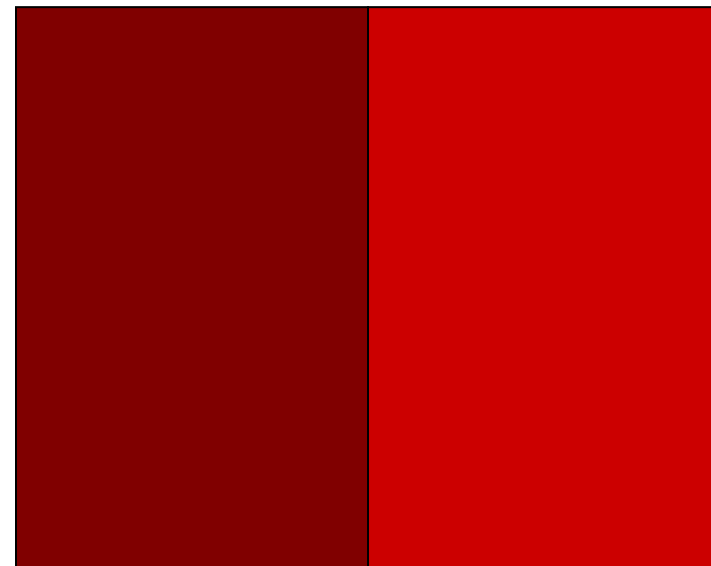


# Color matching experiment 2



Source: W. Freeman

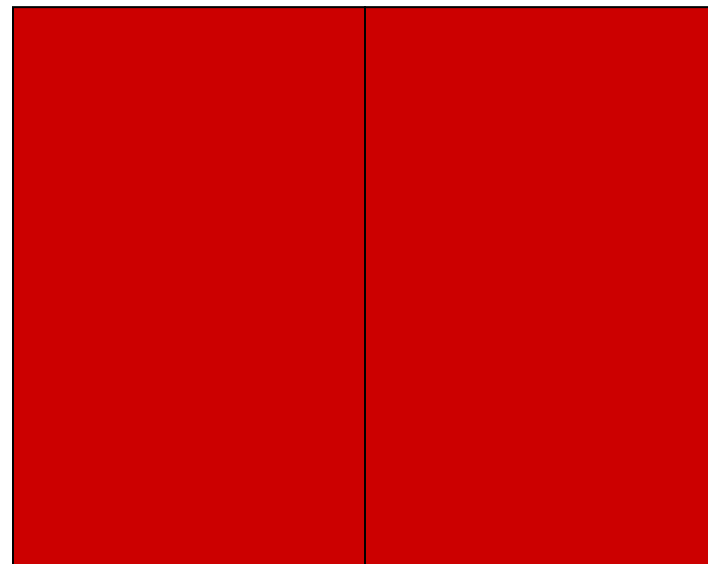
# Color matching experiment 2



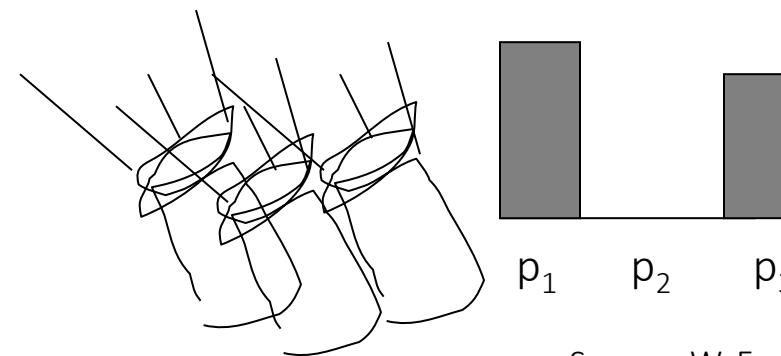
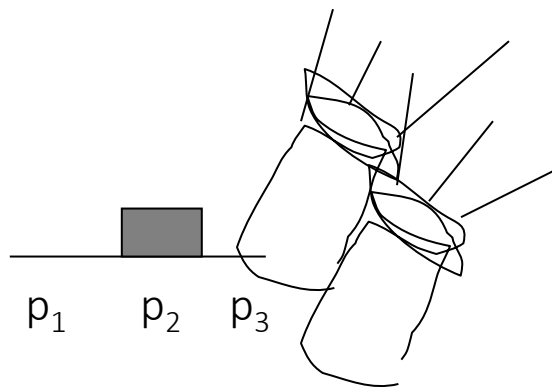
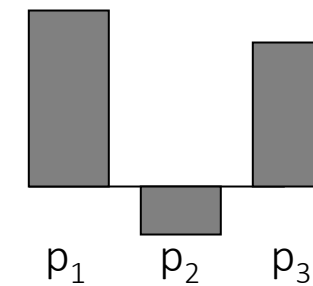
Source: W. Freeman

# Color matching experiment 2

We say a “negative” amount of  $p_2$  was needed to make the match, because we added it to the test color’s side.

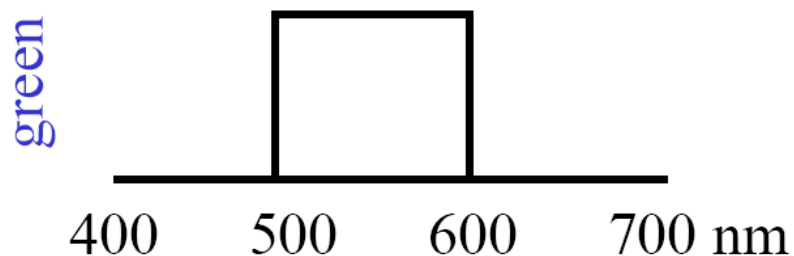
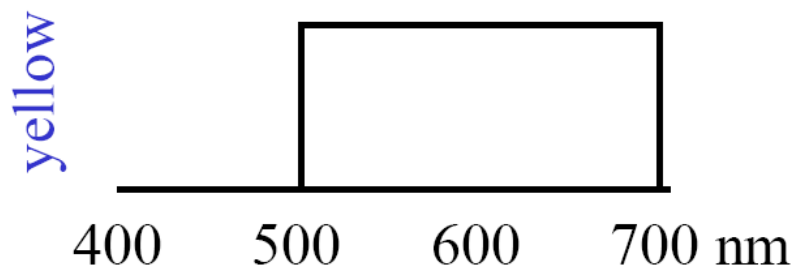
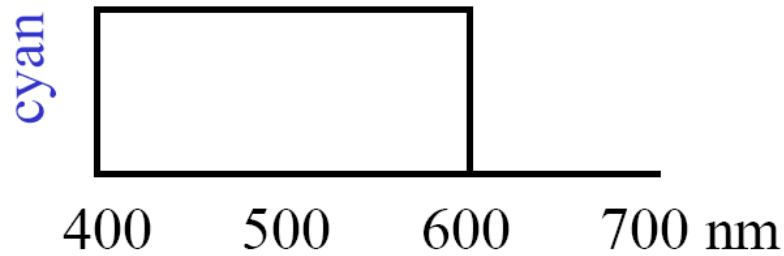


The primary color amounts needed for a match:

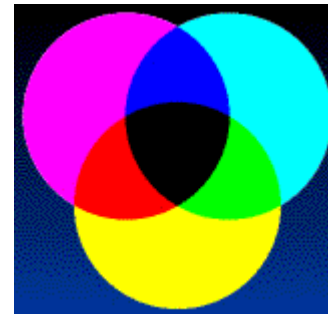


Source: W. Freeman

# Subtractive color mixing



Colors combine by *multiplying* color spectra.



Pigments *remove* color from incident light (white).

Source: W. Freeman



# Examples of subtractive color systems

- Printing on paper
- Crayons
- Photographic film

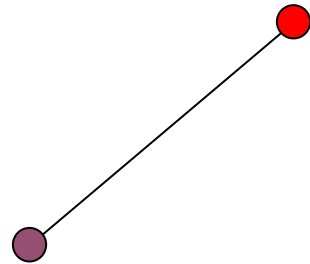


# Overview of Color

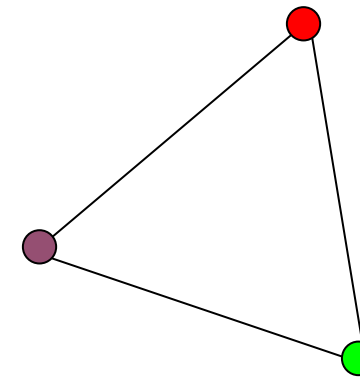
- Physics of color
- Human encoding of color
- Color spaces

# Linear color spaces

- Defined by a choice of three *primaries*
- The coordinates of a color are given by the weights of the primaries used to match it



mixing two lights produces  
colors that lie along a straight  
line in color space



mixing three lights produces  
colors that lie within the triangle  
they define in color space

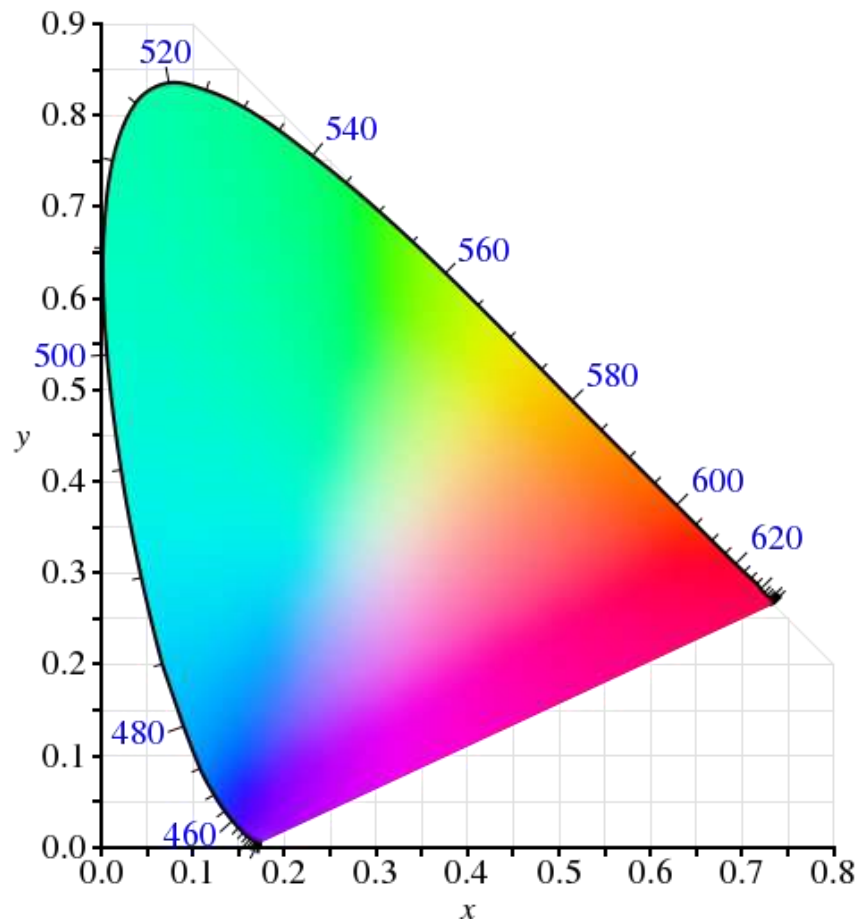
# Linear color spaces

- Pick some primaries
- Can mix those primaries to match any color inside the triangle
- There is a commission that studies color!
  - Commission internationale de l'éclairage (CIE) is a 100-year-old organization that creates international standards related to light and color.

[Shapiro]



# CIE 1931 XYZ color space

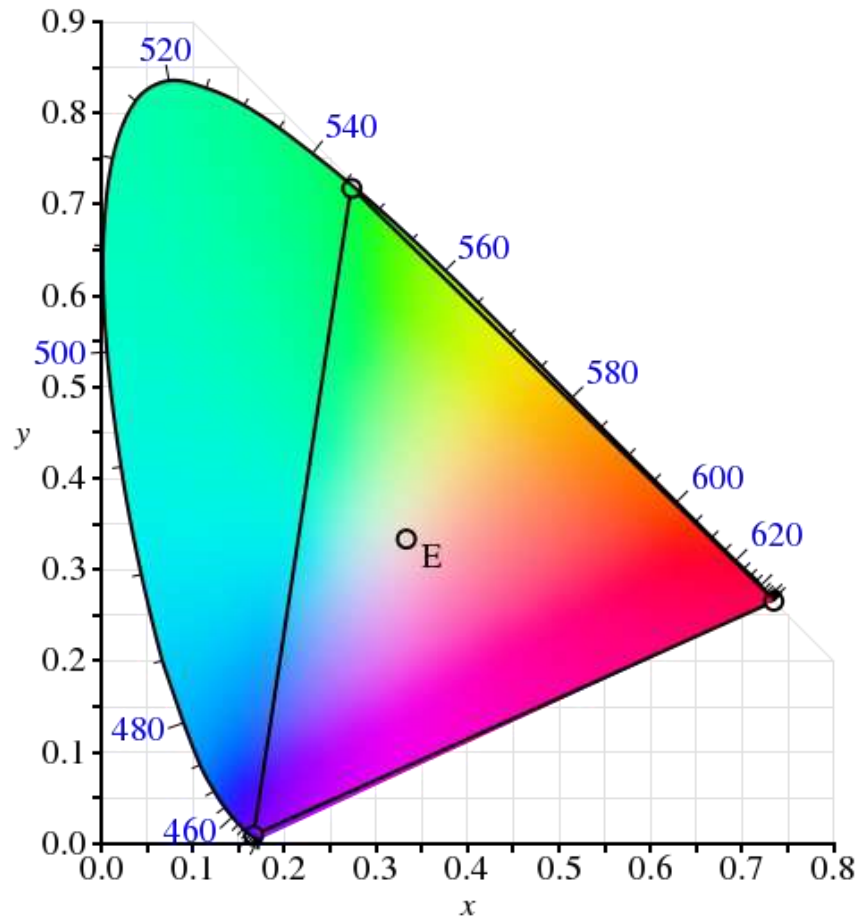


- First attempt to produce a color space based on measurements of human color perception
- Y corresponds to brightness or *luminance* of a color. Z is like blue light stimulation, and X is a mix (a linear combination) of cone response curves chosen to be nonnegative.
- 2D visualization - draw  $(x,y)$ , where:  
 $x = X/(X+Y+Z)$ ,  $y = Y/(X+Y+Z)$

[http://en.wikipedia.org/wiki/CIE\\_1931\\_color\\_space](http://en.wikipedia.org/wiki/CIE_1931_color_space)

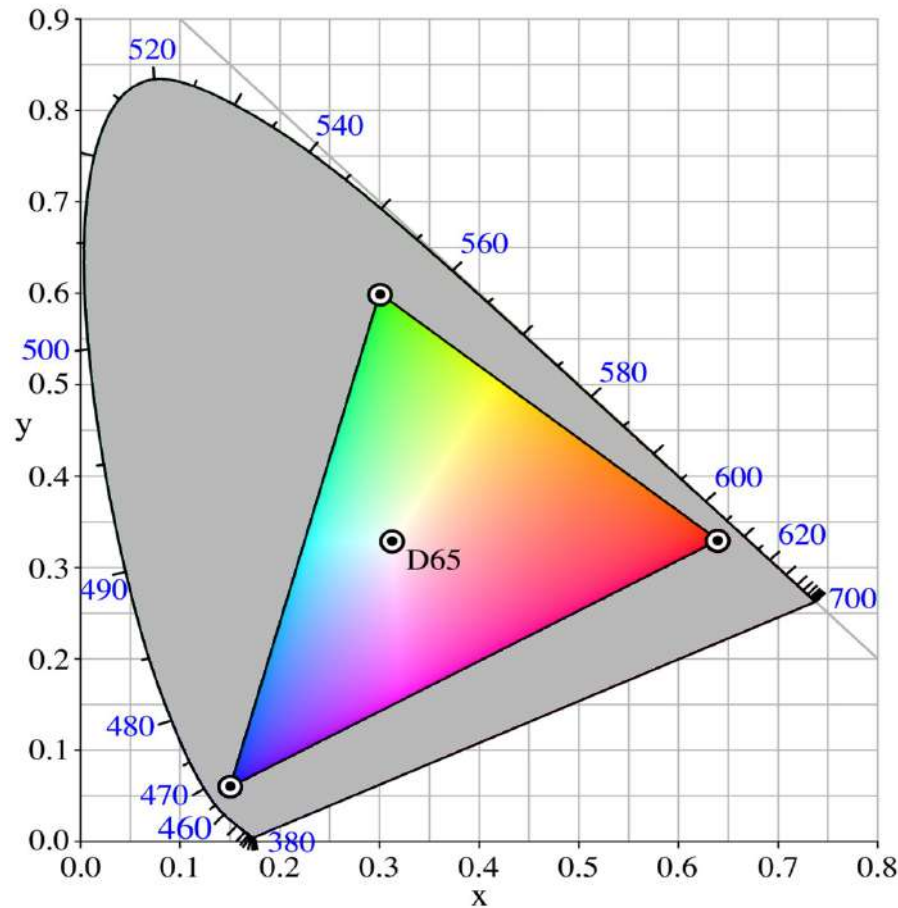
[Shapiro]

# 'Theoretical' CIE RGB primaries



[Shapiro]

# Practical sRGB primaries, MSFT 1996



- sRGB (standard Red Green Blue) is an RGB color space that HP and Microsoft created cooperatively in 1996 to use on monitors, printers, and the Internet
- It was subsequently standardized by the IEC (International Electrotechnical Commission) as IEC 61966-2-1:1999

[Shapiro]

# What does this mean for computers?

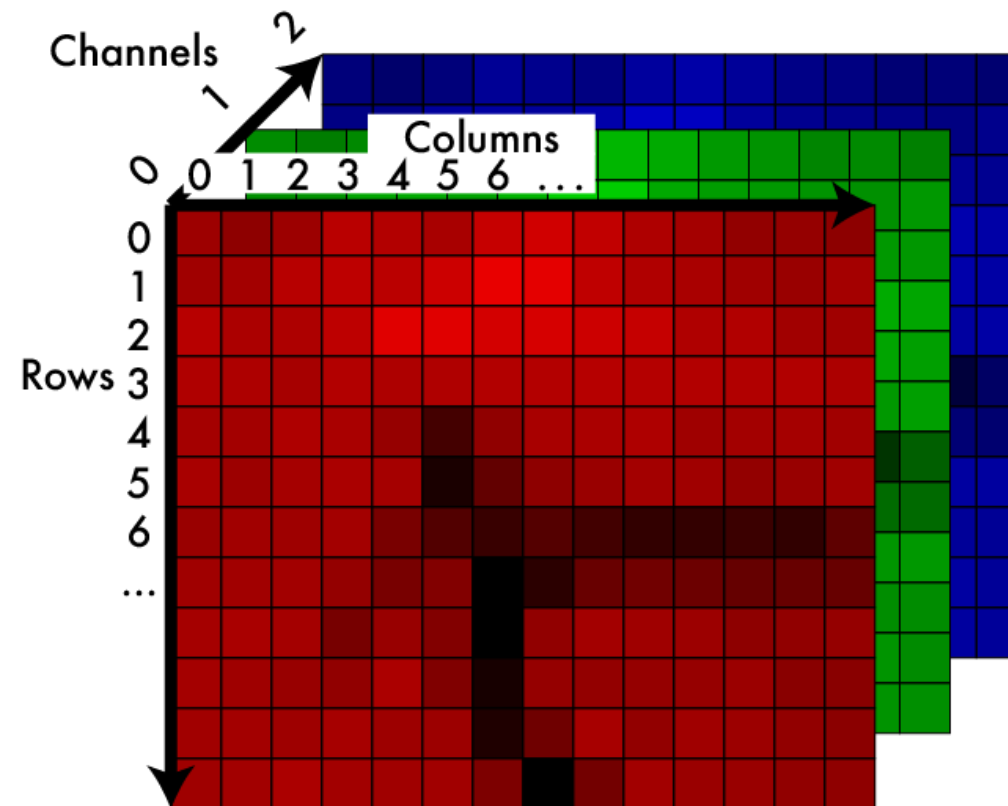
- We represent images as grids of pixels
- Each pixel has a color, 3 components: RGB
- Not every color can be represented in RGB!
  - Have to go out in the real world sometimes
- RGB is not fully accurate
- We can represent a color with 3 numbers
  - #ff00ff; (1.0, 0.0, 1.0); 255,0,255
  - What color is this?

[Shapiro]



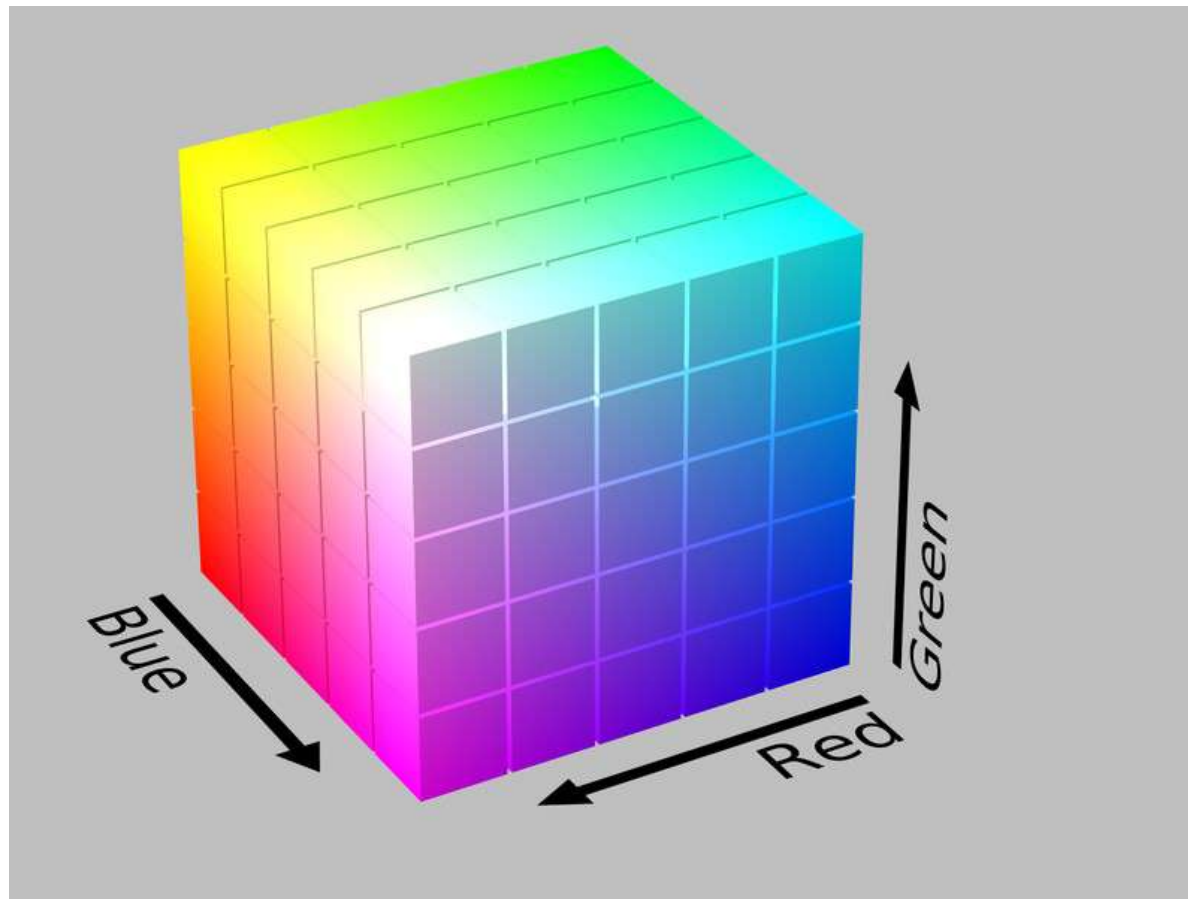
# Image: 2d array of color

- Some range
  - [0,255]
  - [0.0,1.0]
- We'll talk more about this later.



[Shapiro]

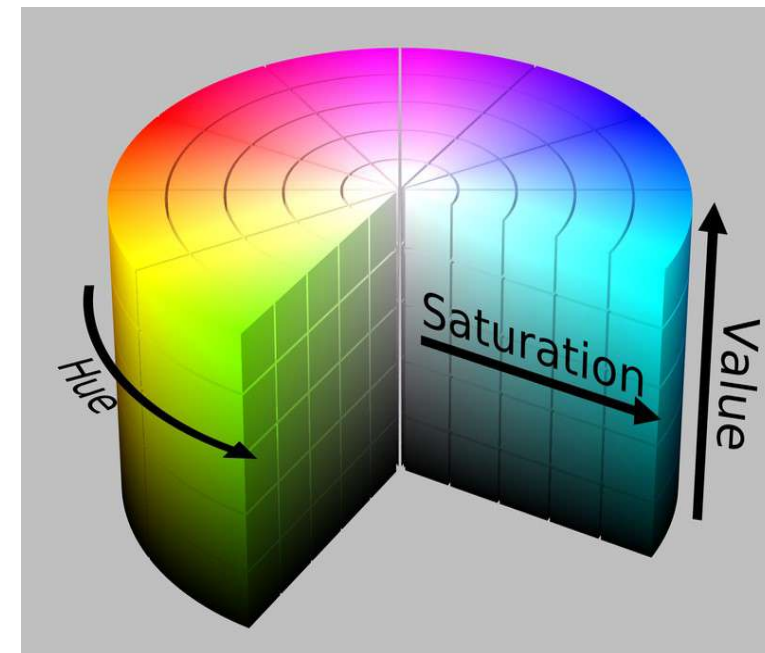
## RGB is a cube...



[Shapiro]

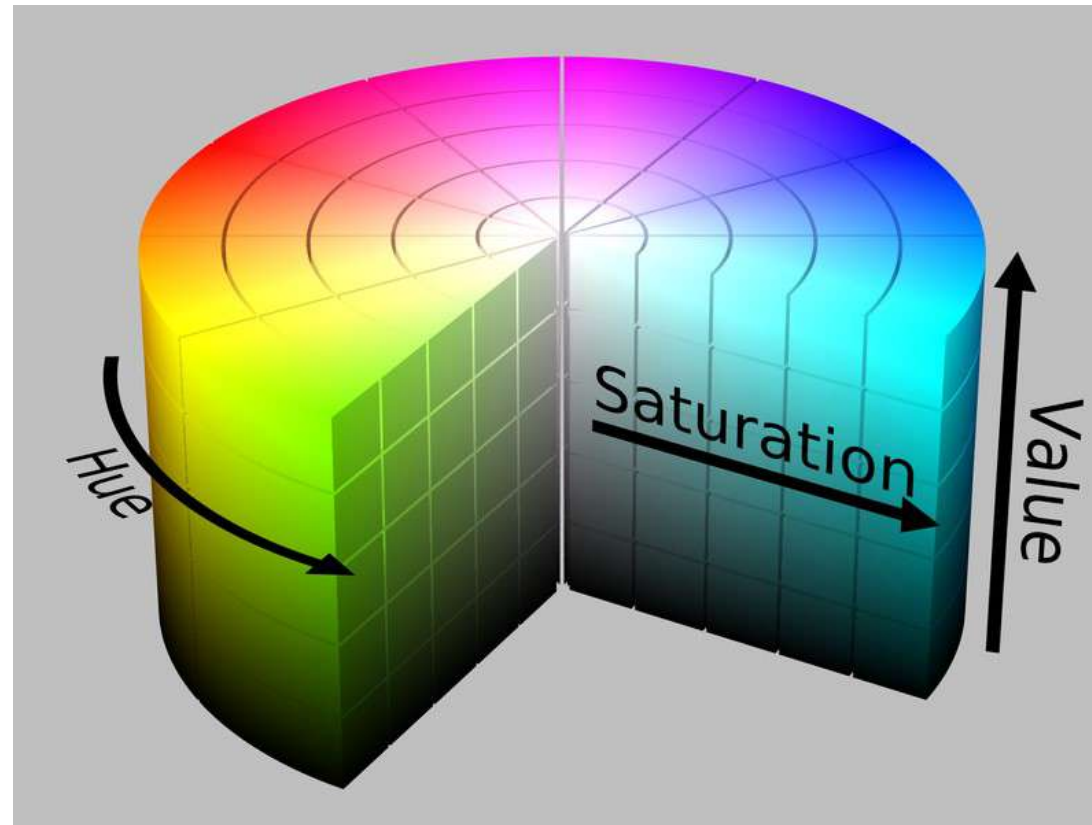
# Non-linear color spaces: Hue, Saturation, Value

- Different model based on perception of light
- **Hue**: what color
- **Saturation**: how much color
- **Value**: how bright
- Allows easy image transforms
  - Shift the hue
  - Increase saturation



[Shapiro]

# Hue, Saturation, Value: cylinder!



[Shapiro]

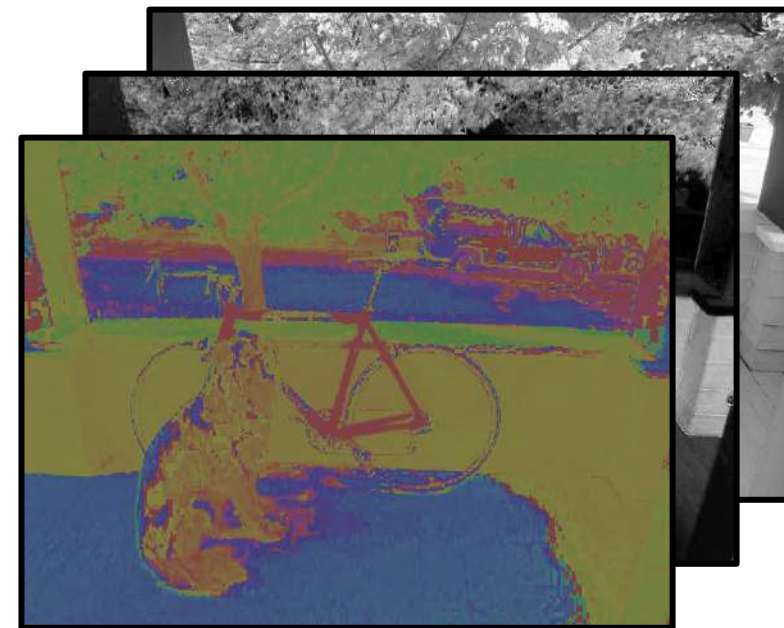
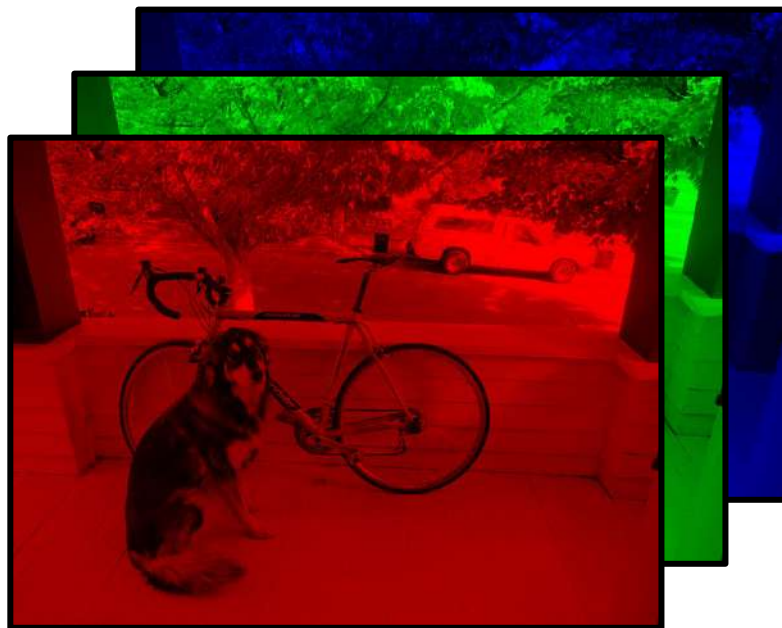


## An RGB Image



[Shapiro]

## Still 3d tensor, different info

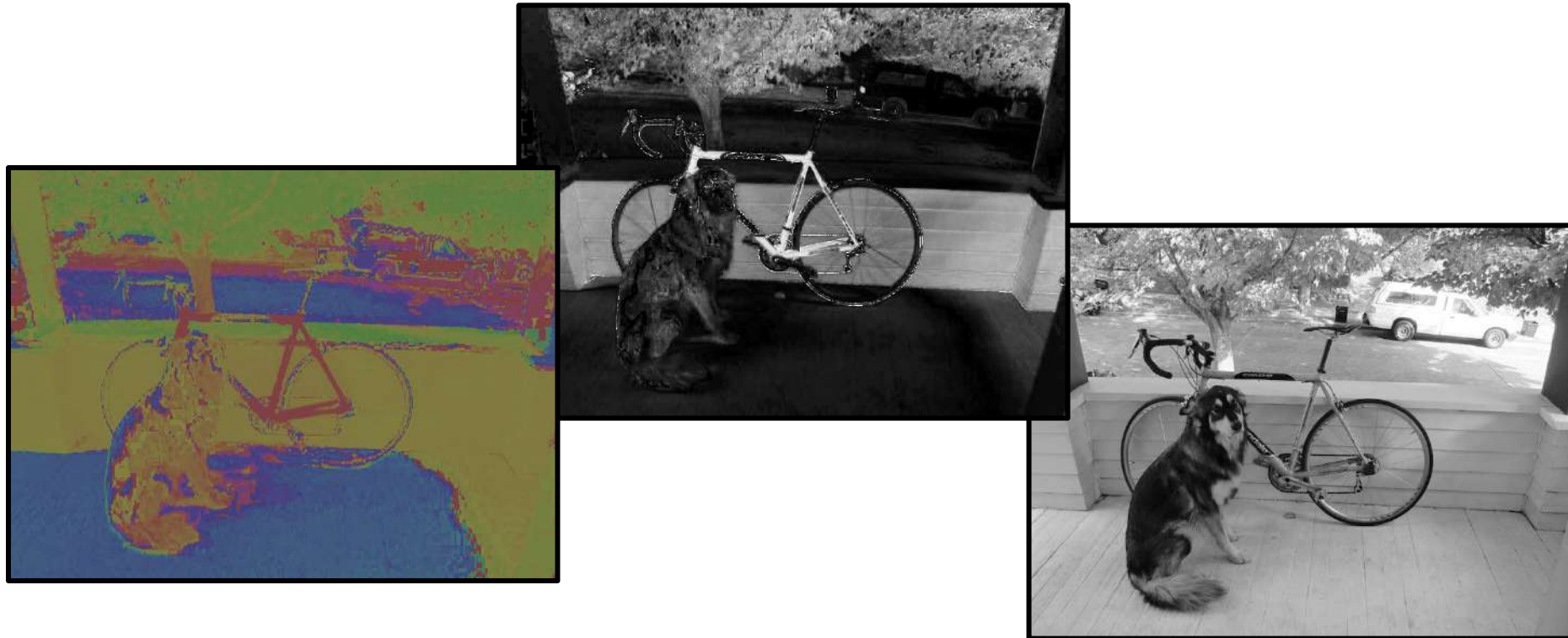


[Shapiro]

## Hue

## Saturation

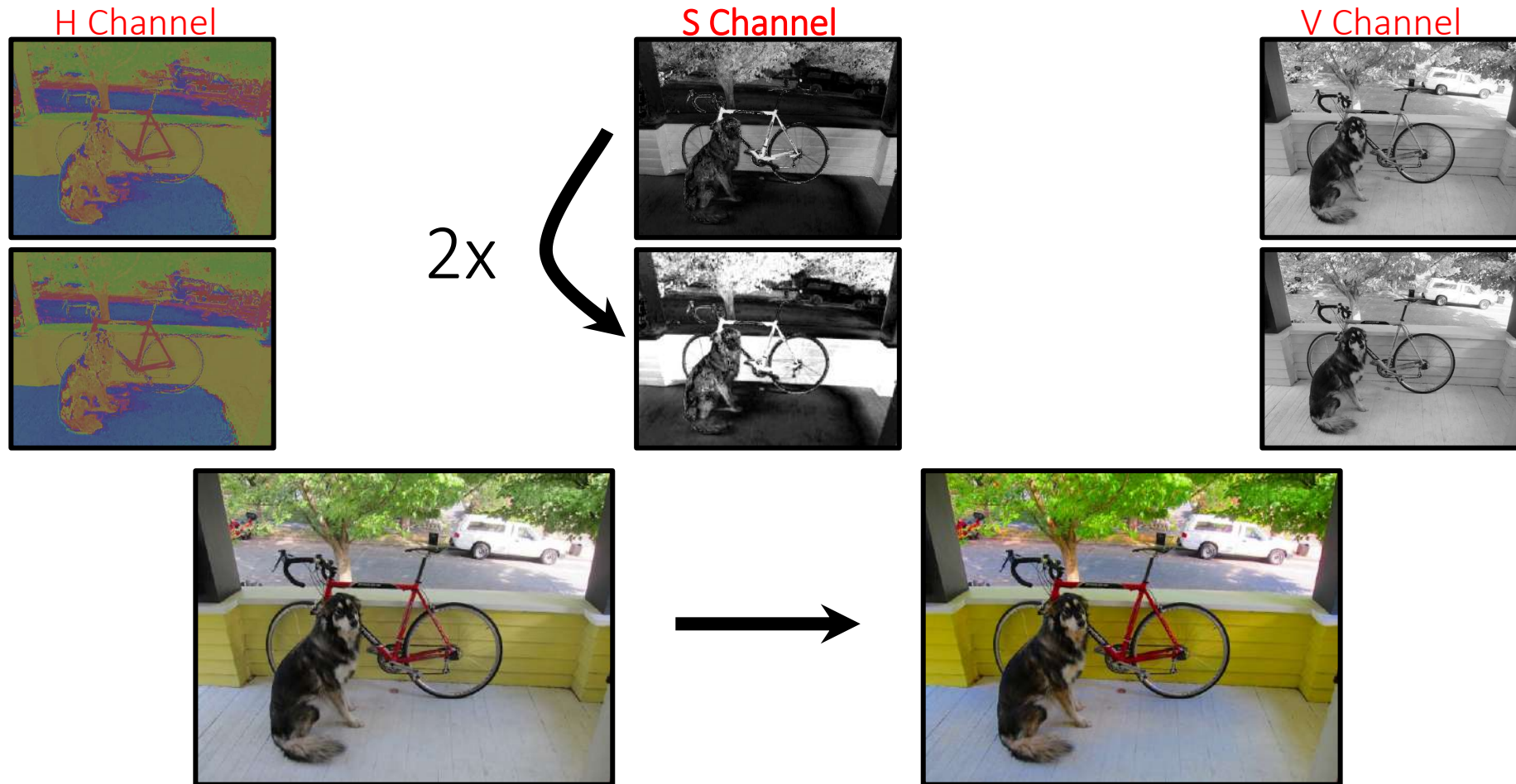
## Value



[Shapiro]



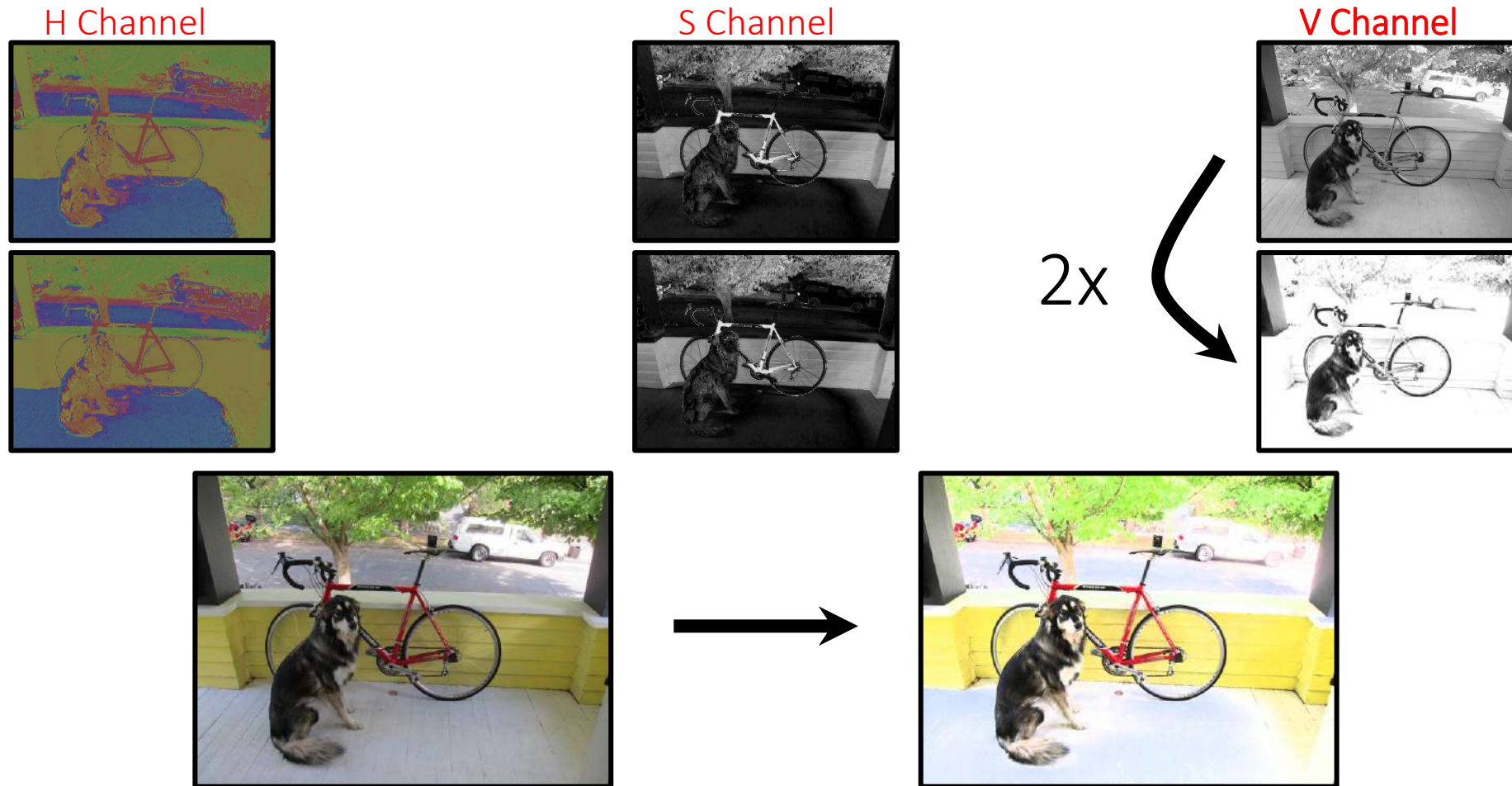
# More saturation = intense colors



[Shapiro]

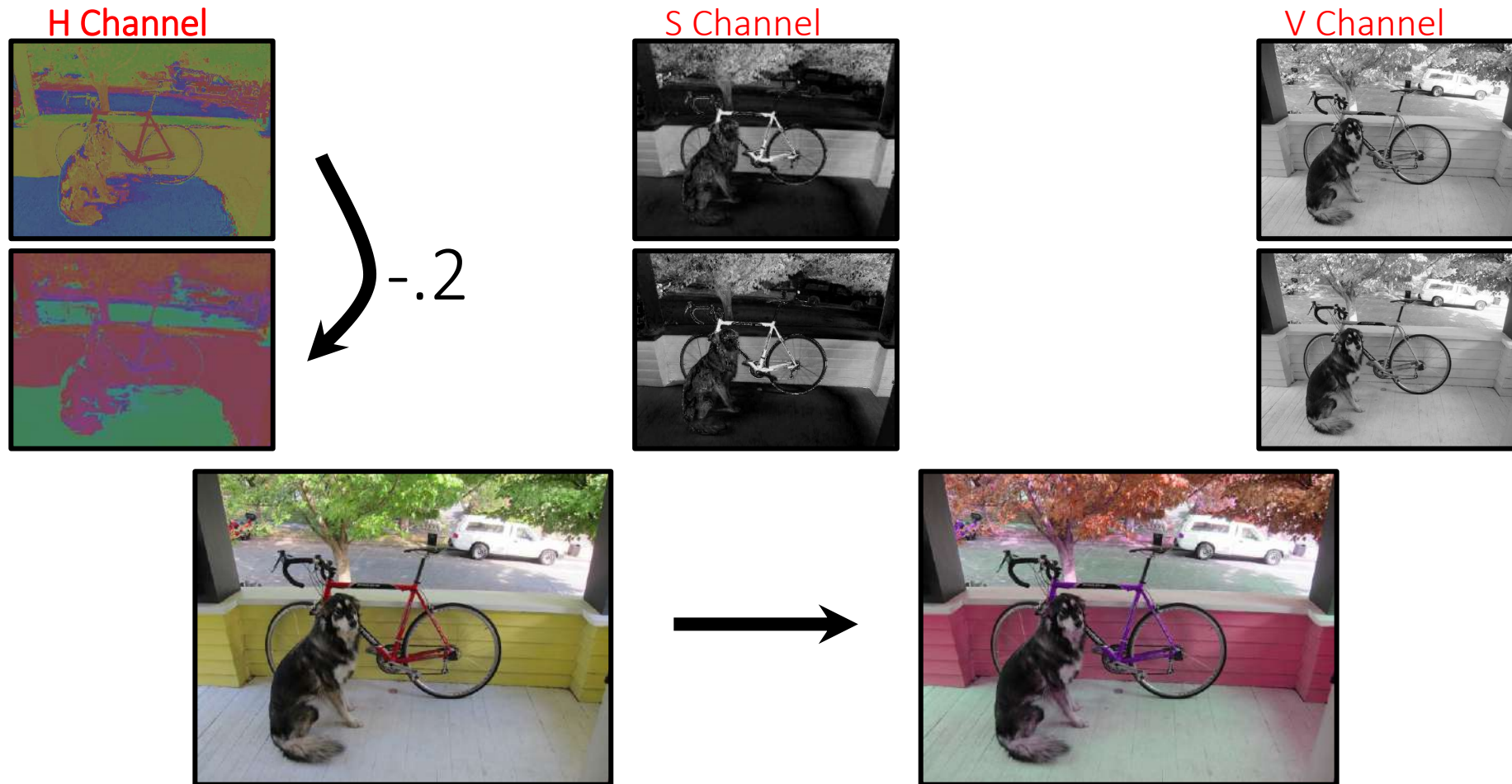


# More value = lighter image



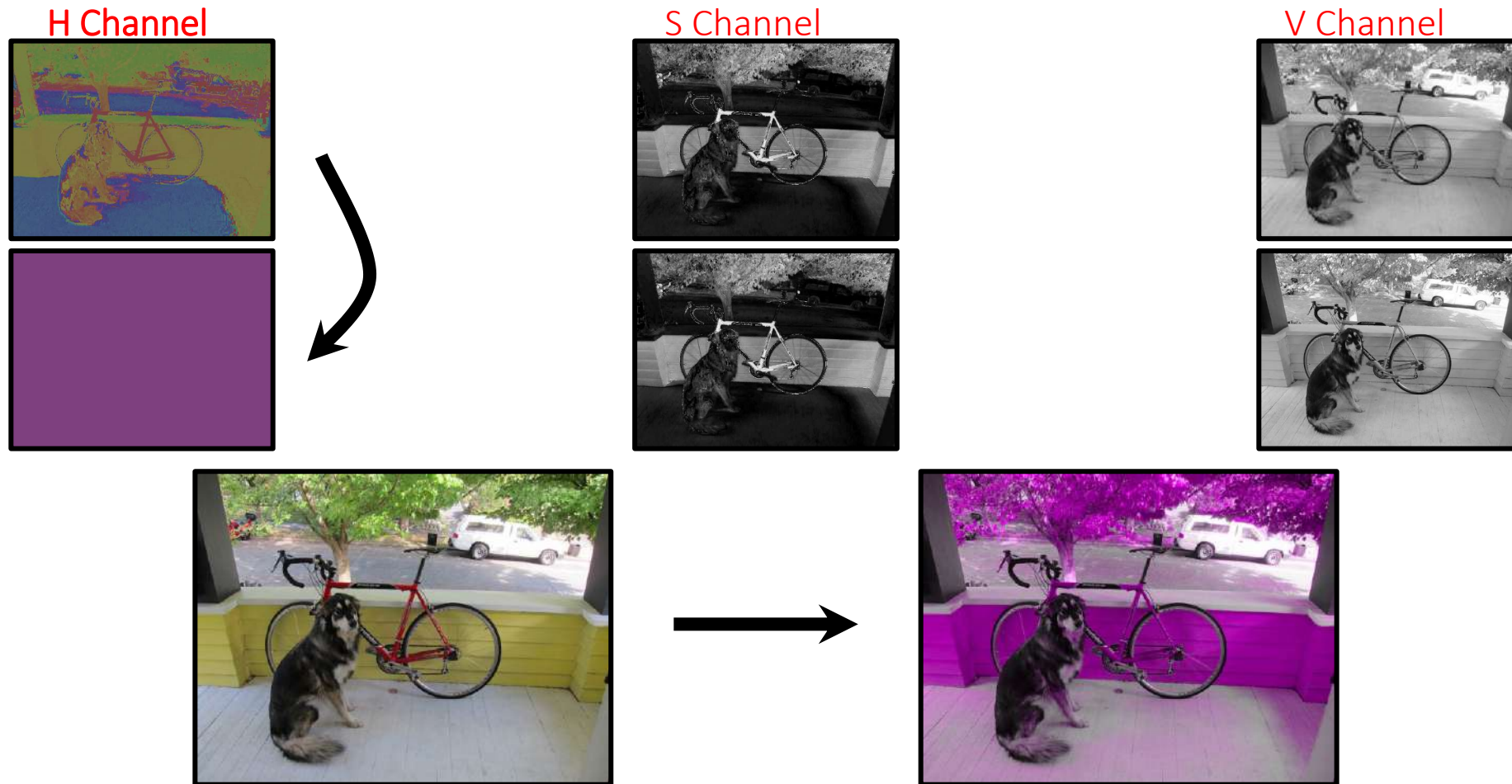
[Shapiro]

# Shift hue = shift colors



[Shapiro]

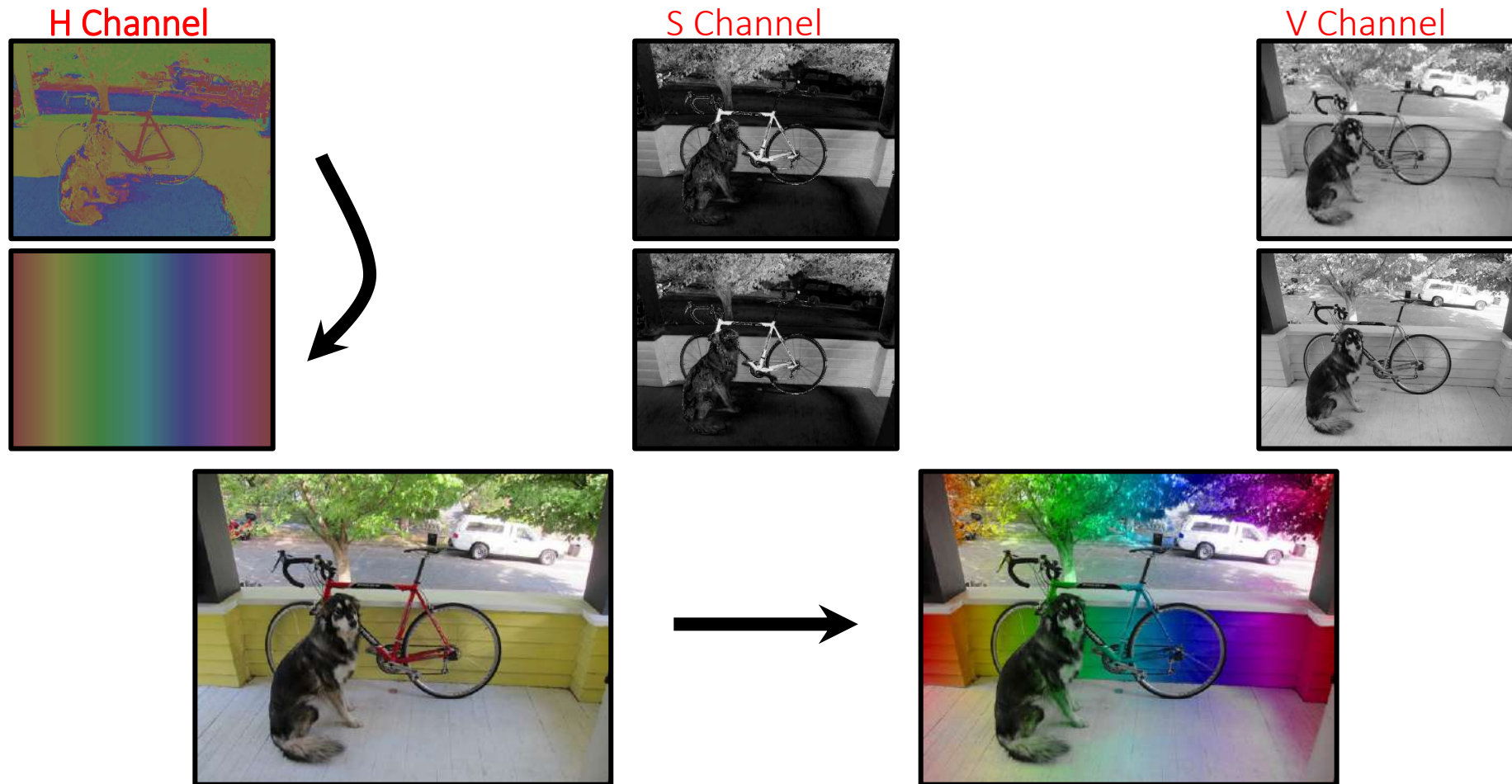
# Set hue to your favorite color!



[Shapiro]



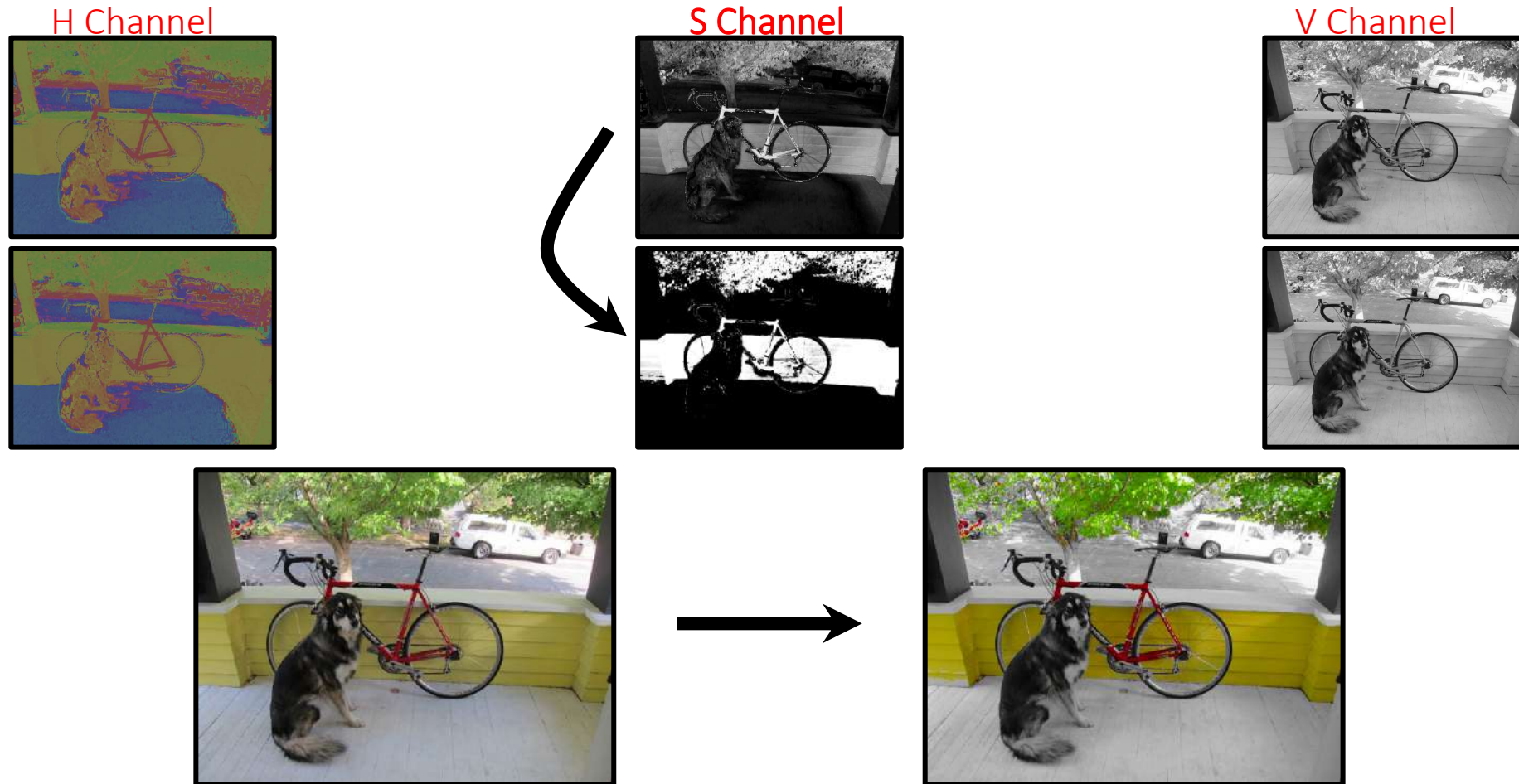
## Or pattern...



[Shapiro]



## Increase and threshold saturation



[Shapiro]



[Shapiro]

**MAI4CAREU**

Master programmes in Artificial  
Intelligence 4 Careers in Europe



# Thank you.





University of Cyprus – MSc Artificial Intelligence

# MAI644 – COMPUTER VISION

## Lecture 3: Fundamentals - Cameras

**Melinos Averkiou**

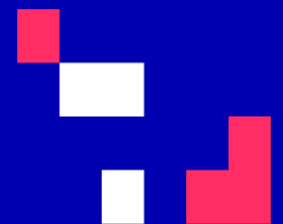
CYENS Centre of Excellence

University of Cyprus - Department of Computer Science

[m.averkiou@cyens.org.cy](mailto:m.averkiou@cyens.org.cy)



**CYENS**  
CENTRE OF EXCELLENCE





## Last time

- Physics of color
- Human encoding of color
- Color spaces

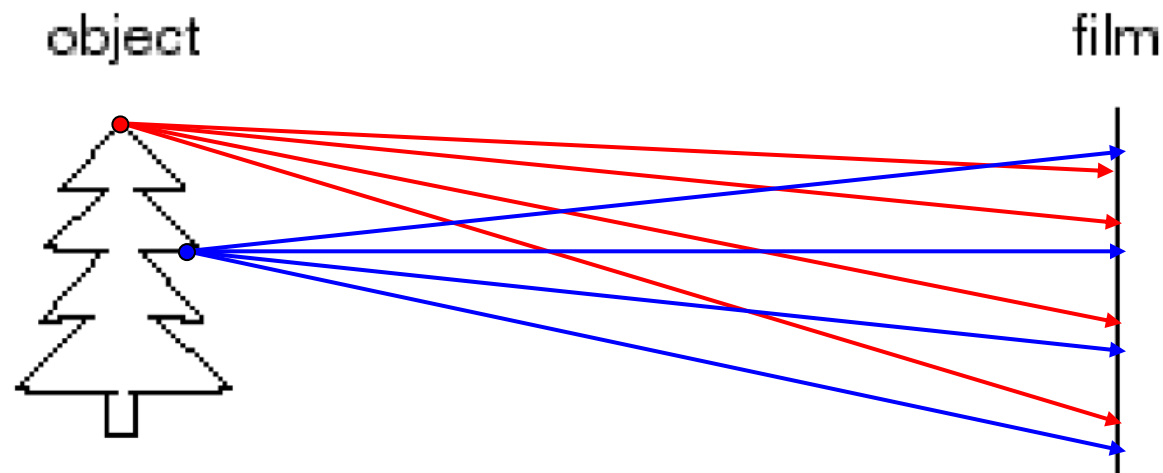
# Today's Agenda - Overview of Cameras

- Pinhole Camera model
  - Aperture
  - Camera Obscura
- Cameras with lenses
  - Thin lens equation
  - Depth of field
  - Field of view
- Digital cameras
  - Bayer filters
  - Debayering

# Today's Agenda - Overview of Cameras

- Pinhole Camera model
  - Aperture
  - Camera Obscura
- Cameras with lenses
  - Thin lens equation
  - Depth of field
  - Field of view
- Digital cameras
  - Bayer filters
  - Debayering

# Let's design a camera

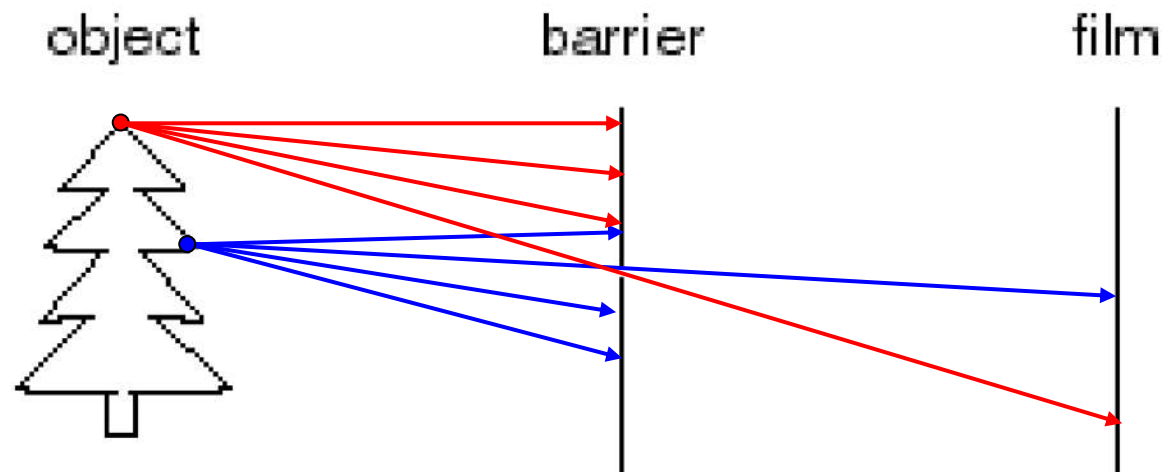


- Idea 1: put a piece of film in front of an object
- Do we get a reasonable image?

[Slide by Steve Seitz]



# Pinhole camera



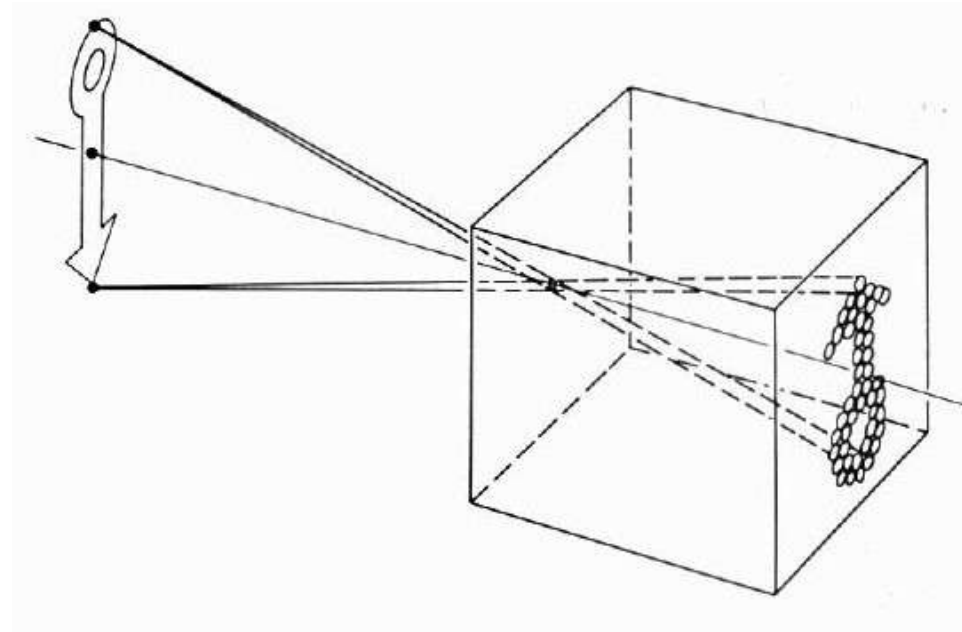
Add a barrier to block off most of the rays

- This reduces blurring
- The opening known as the **aperture**
- How does this transform the image?

[Slide by Steve Seitz]

# Pinhole Camera model

- Captures **pencil of rays** – all rays through a single point
- The point is called **Center of Projection**
- The image is formed on the **Image Plane**



[Slide by Steve Seitz]

# Pinhole cameras everywhere



Sun “shadows” during a solar eclipse  
by Henrik von Wendt <http://www.flickr.com/photos/hvw/2724969199/>

[Slide by Steve Seitz]

# Pinhole cameras everywhere



Sun “shadows” during a solar eclipse

<http://www.flickr.com/photos/73860948@N08/6678331997/>

[Slide by Steve Seitz]



# Pinhole cameras everywhere



Tree shadow during a solar eclipse

photo credit: Nils van der Burg <http://www.physicstogo.org/index.cfm>

[Slide by Steve Seitz]

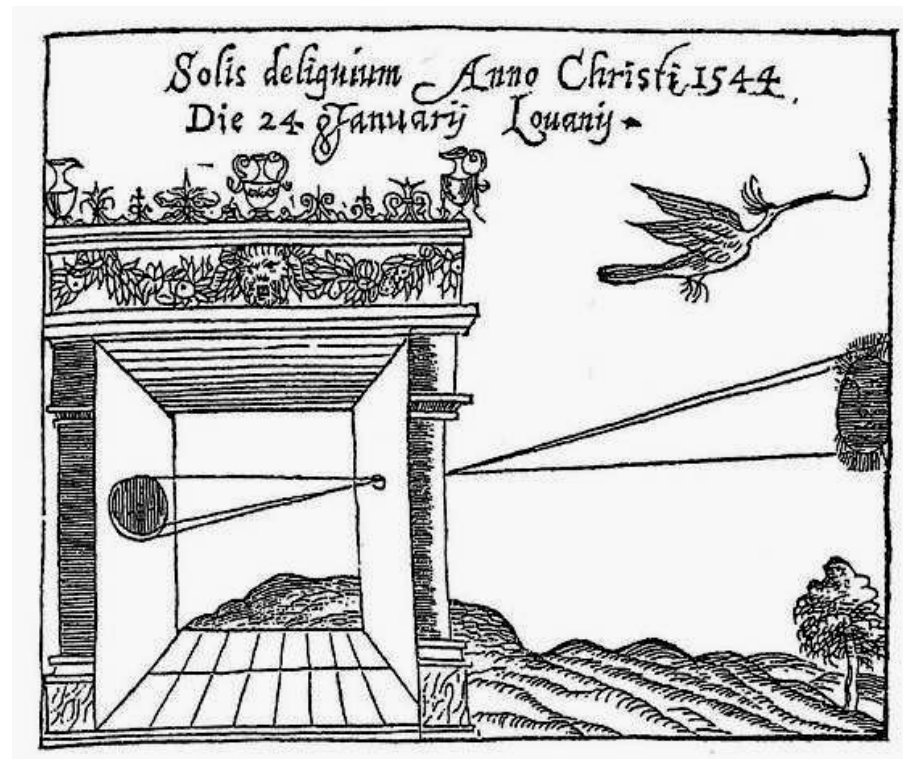
# Today's Agenda - Overview of Cameras

- Pinhole Camera model
  - Aperture
  - Camera Obscura
- Cameras with lenses
  - Thin lens equation
  - Depth of field
  - Field of view
- Digital cameras
  - Bayer filters
  - Debayering

# Pinhole Camera - Camera Obscura

## The first camera

- Known to Aristotle



Gemma Frisius, 1558

[Slide by Steve Seitz]

# How to Turn a Room into a Camera Obscura



[https://www.youtube.com/watch?v=hsXo4gD7iWI&ab\\_channel=GeorgeEastmanMuseum](https://www.youtube.com/watch?v=hsXo4gD7iWI&ab_channel=GeorgeEastmanMuseum)



## Home-made pinhole camera



Why so blurry?



<http://www.debevec.org/Pinhole/>

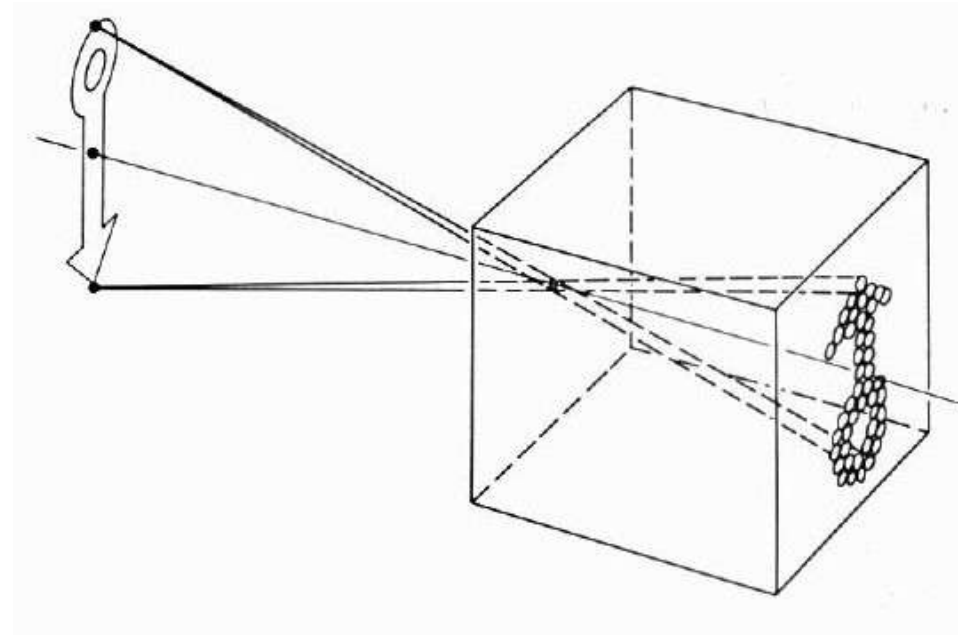
[Slide by A. Efros]



# Pinhole Camera - Camera Obscura

## The first camera

- Known to Aristotle
- How does the aperture size affect the image?

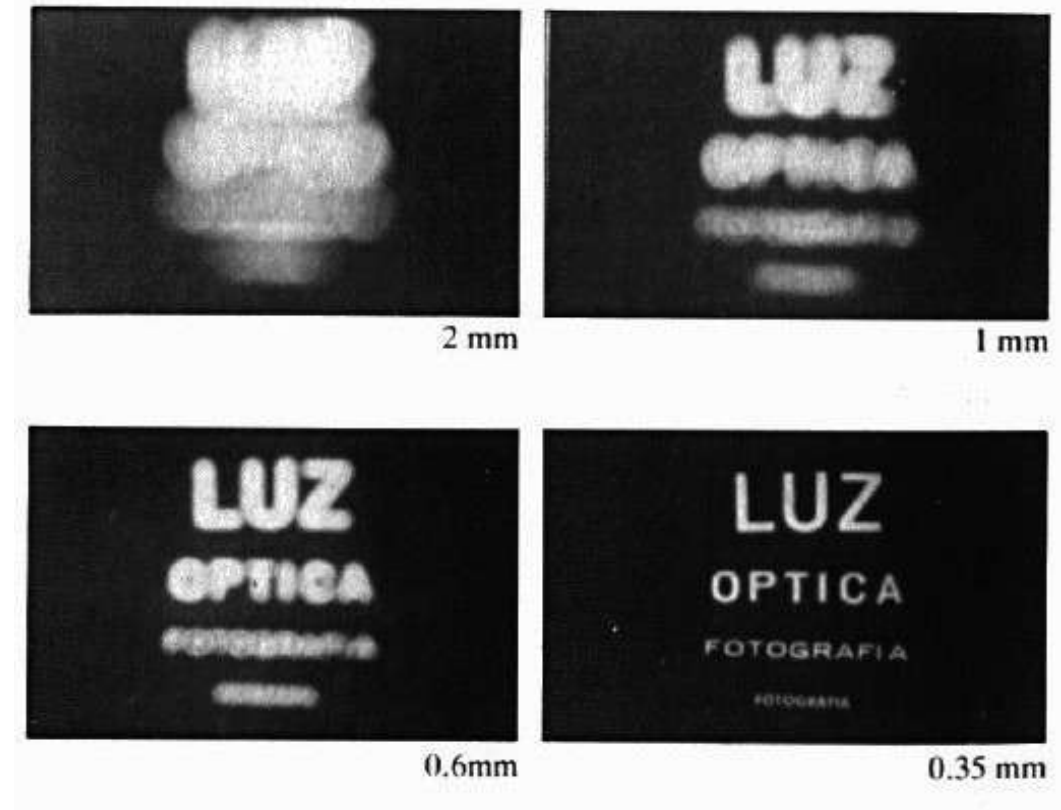


[Slide by Steve Seitz]



# Shrinking the aperture

Why not make the aperture as small as possible?

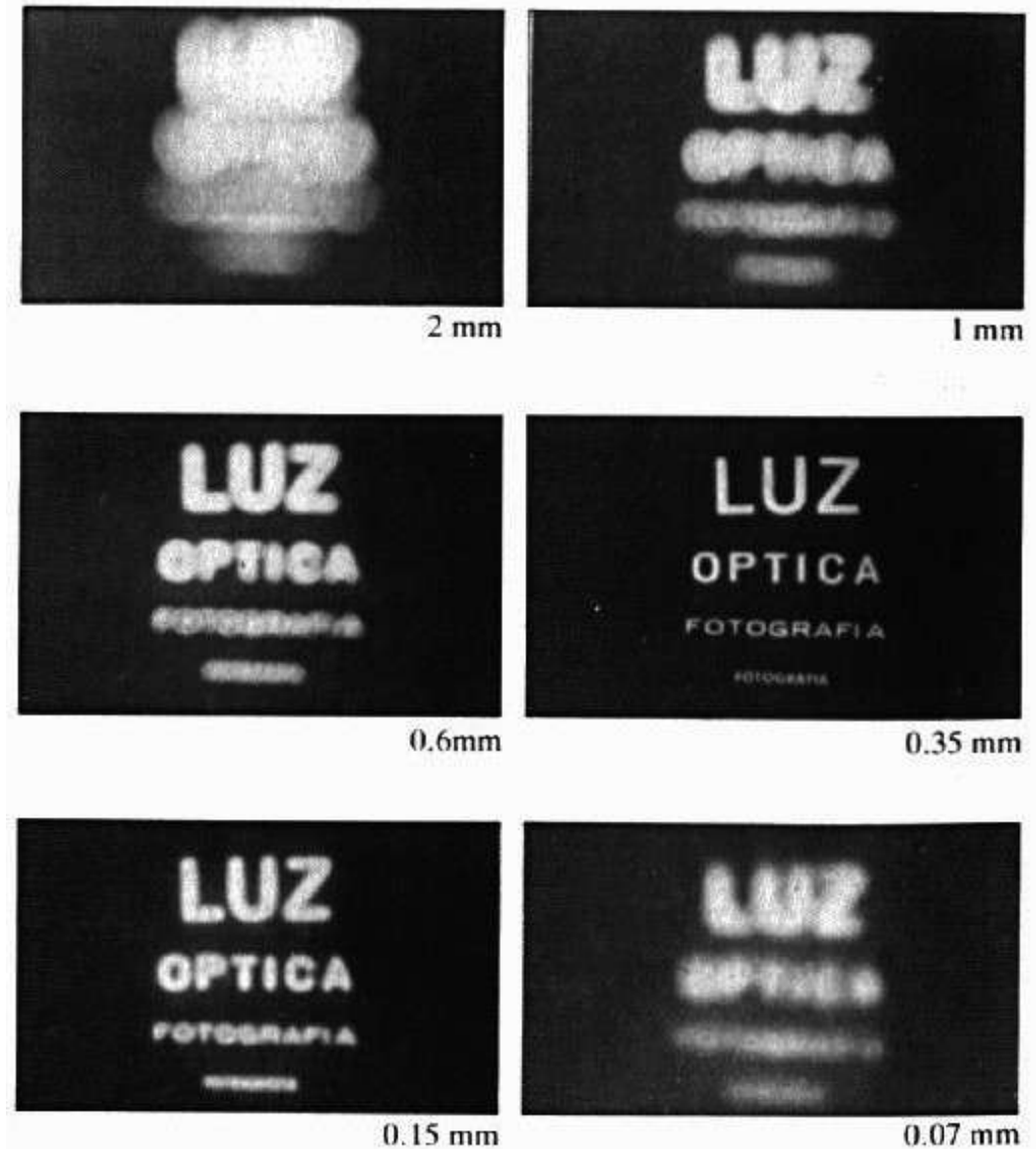
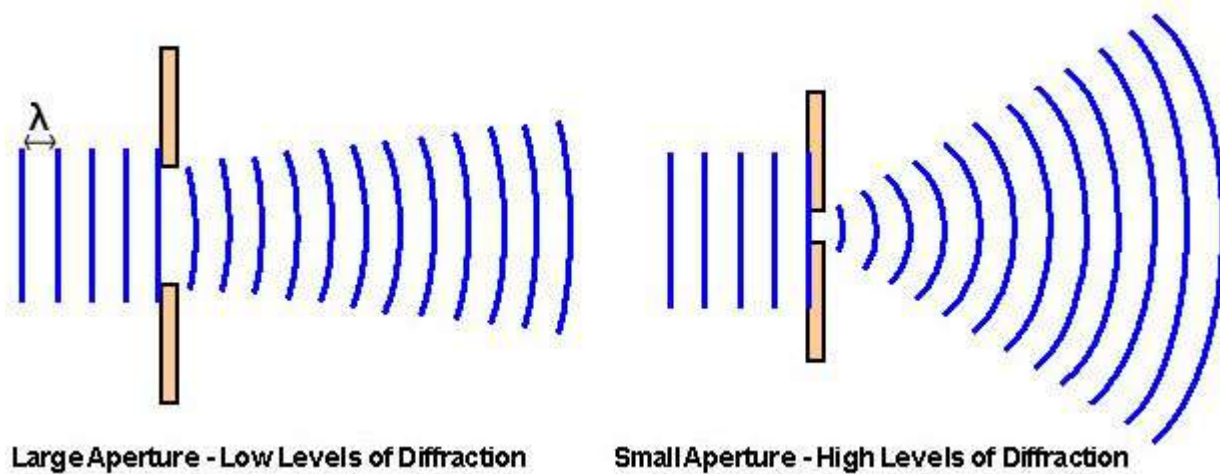


[Slide by Steve Seitz]

# Shrinking the aperture

Why not make the aperture as small as possible?

- Less light gets through
- Diffraction effects



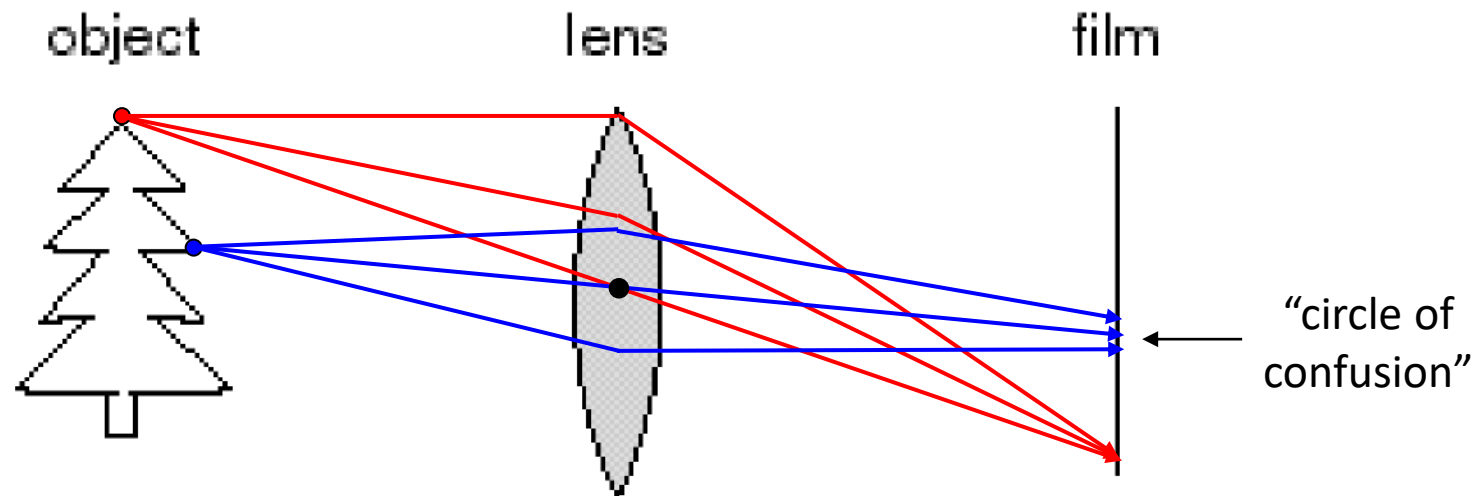
[Slide by Steve Seitz]



# Today's Agenda - Overview of Cameras

- Pinhole Camera model
  - Aperture
  - Camera Obscura
- Cameras with lenses
  - Thin lens equation
  - Depth of field
  - Field of view
- Digital cameras
  - Bayer filters
  - Debayering

# Adding a lens

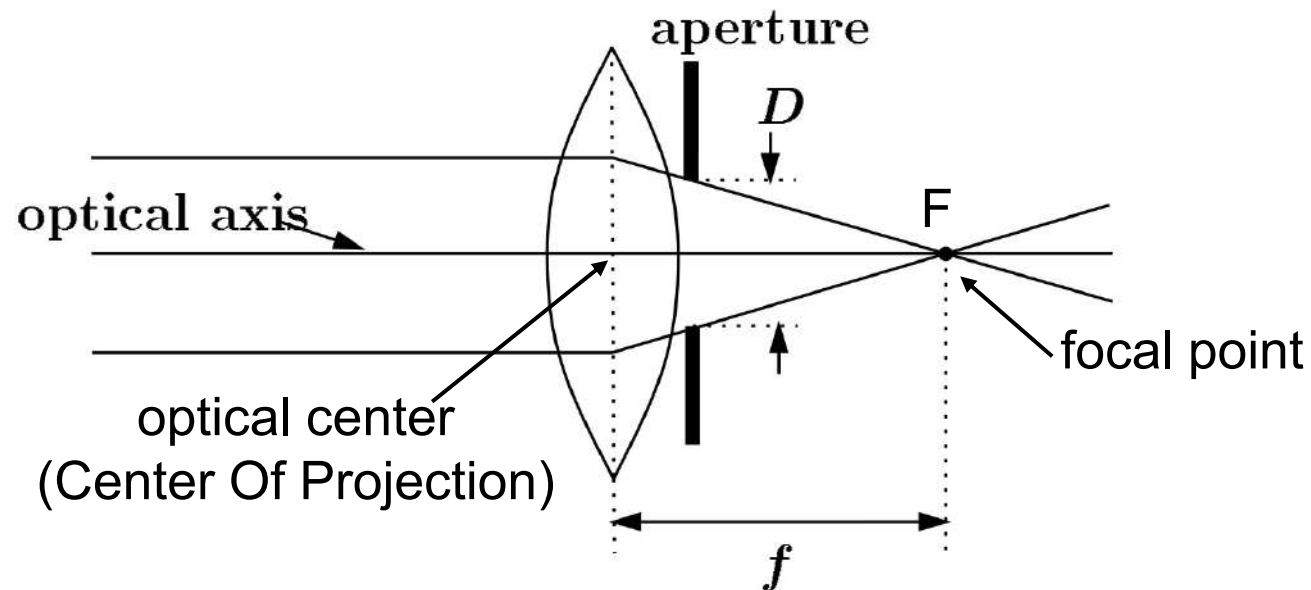


A lens focuses light onto the film

- Rays passing through the center are not deviated
- There is a specific distance at which objects are “in focus”
  - other points project to a “circle of confusion” in the image

[Slide by Steve Seitz]

# Lenses



A lens focuses rays parallel to its axis onto a single **focal point**

- Focal point is on a plane located at a distance  $f$  (**focal length**) beyond the plane of the lens
  - $f$  is a function of the shape and index of refraction of the lens
- Aperture of diameter  $D$  restricts the range of rays
  - aperture may be on either side of the lens
- Lenses are typically spherical (easier to produce)

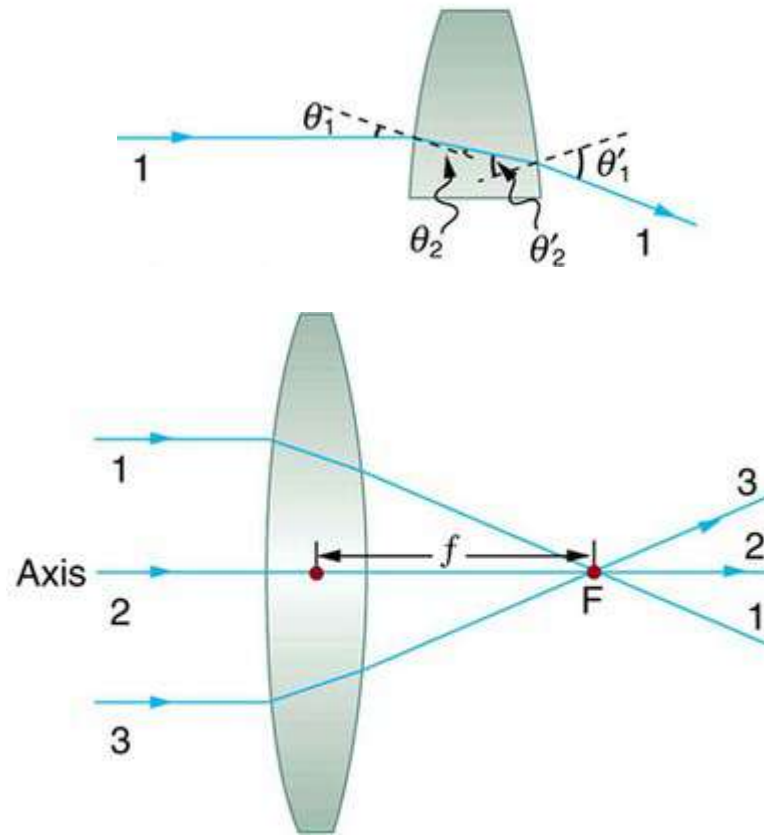
[Slide by Steve Seitz]

# Today's Agenda - Overview of Cameras

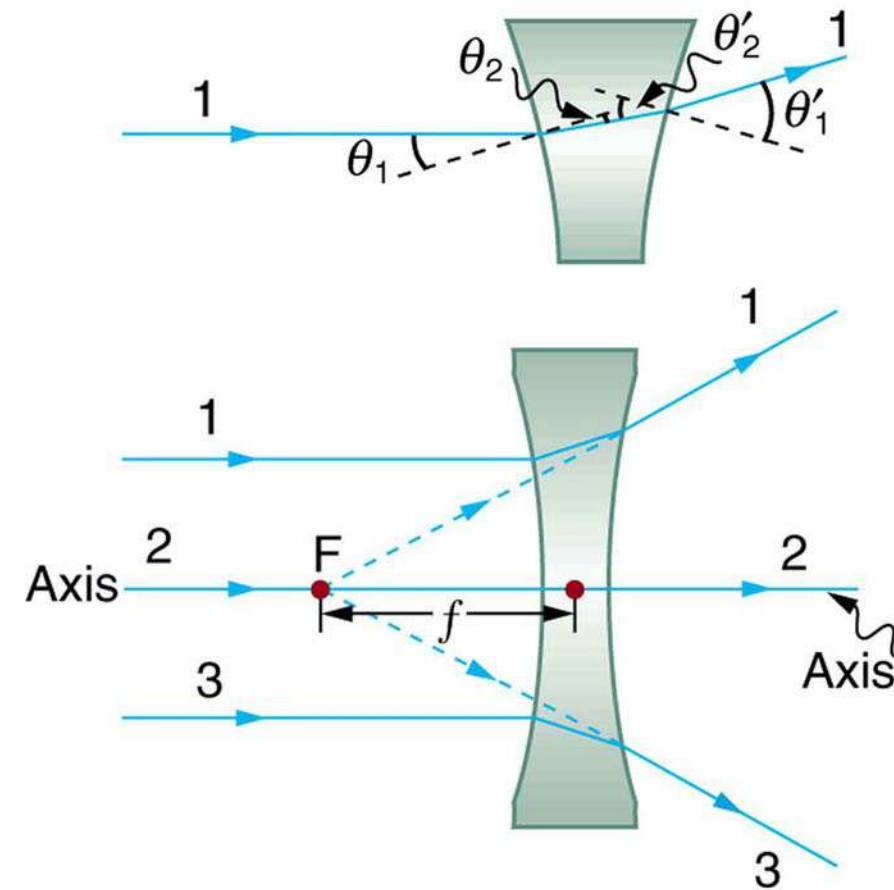
- Pinhole Camera model
  - Aperture
  - Camera Obscura
- Cameras with lenses
  - Thin lens equation
  - Depth of field
  - Field of view
- Digital cameras
  - Bayer filters
  - Debayering



# Thin lenses



Converging – Convex lens



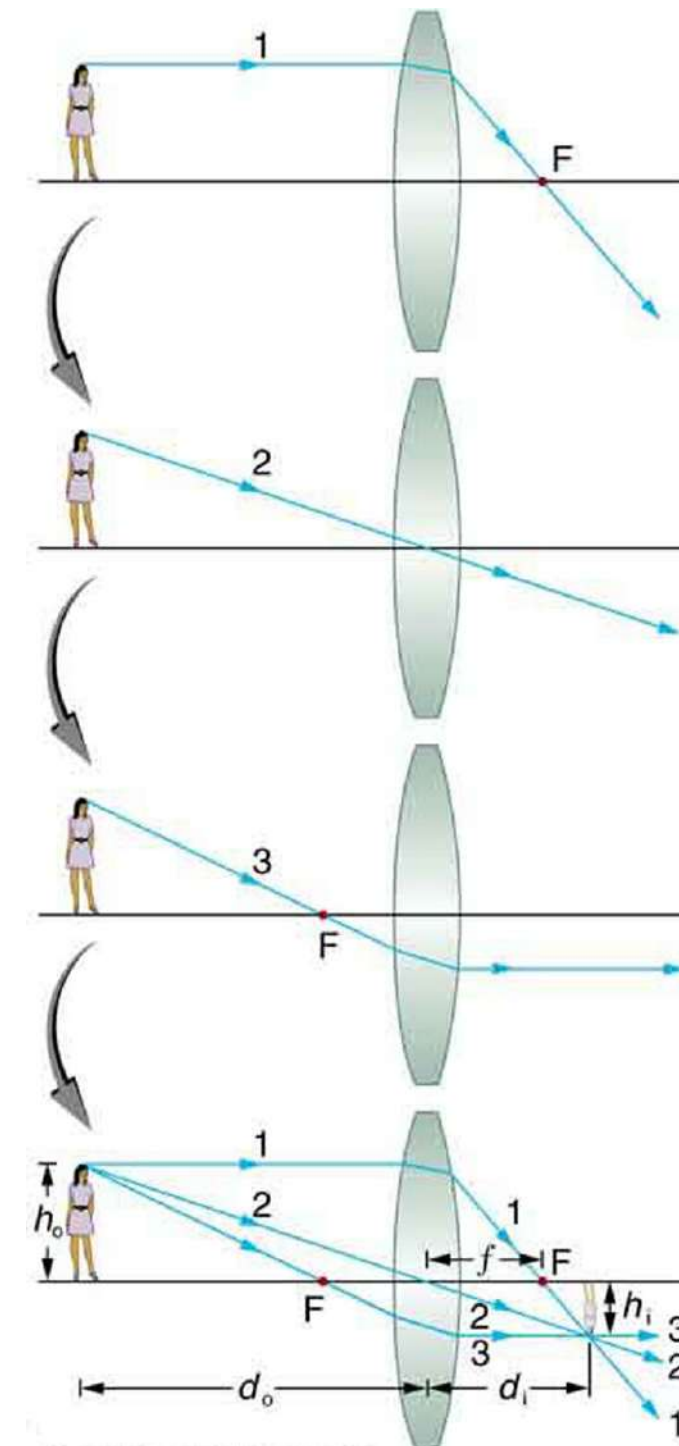
Diverging – Concave lens

[See here](#)

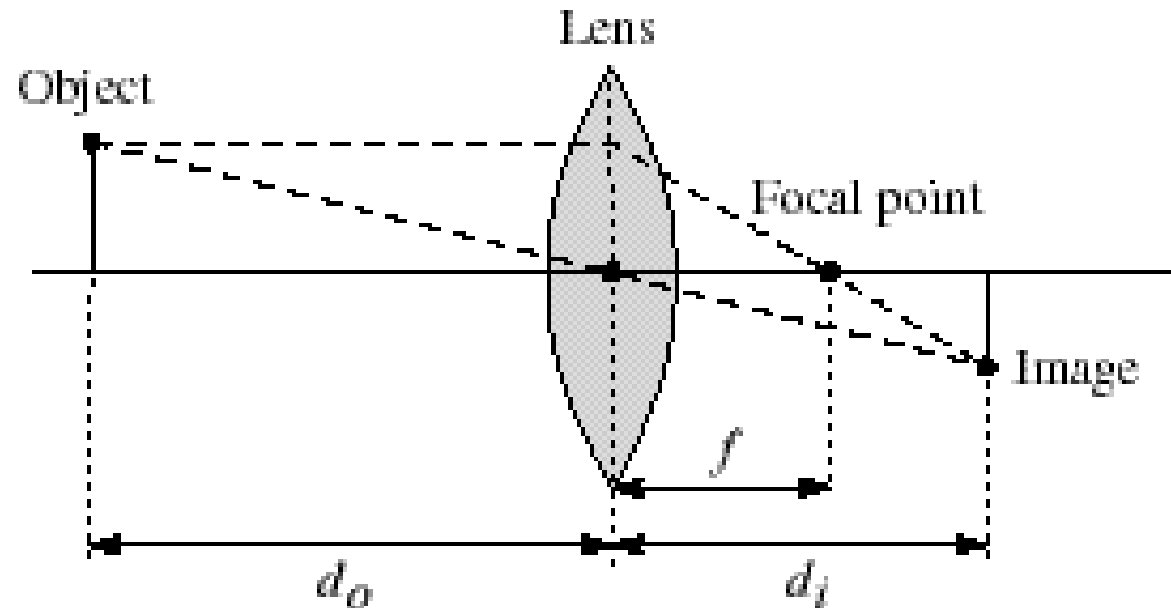
# Thin lenses – ray tracing

- Five rules for image formation with thin lenses
  1. A ray entering a converging lens parallel to its axis passes through the focal point  $F$  of the lens on the other side.
  2. A ray entering a diverging lens parallel to its axis seems to come from the focal point  $F$ .
  3. A ray passing through the center of either a converging or a diverging lens does not change direction.
  4. A ray entering a converging lens through its focal point exits parallel to its axis.
  5. A ray that enters a diverging lens by heading toward the focal point on the opposite side exits parallel to the axis.

[See here](#)



# Thin lens equation

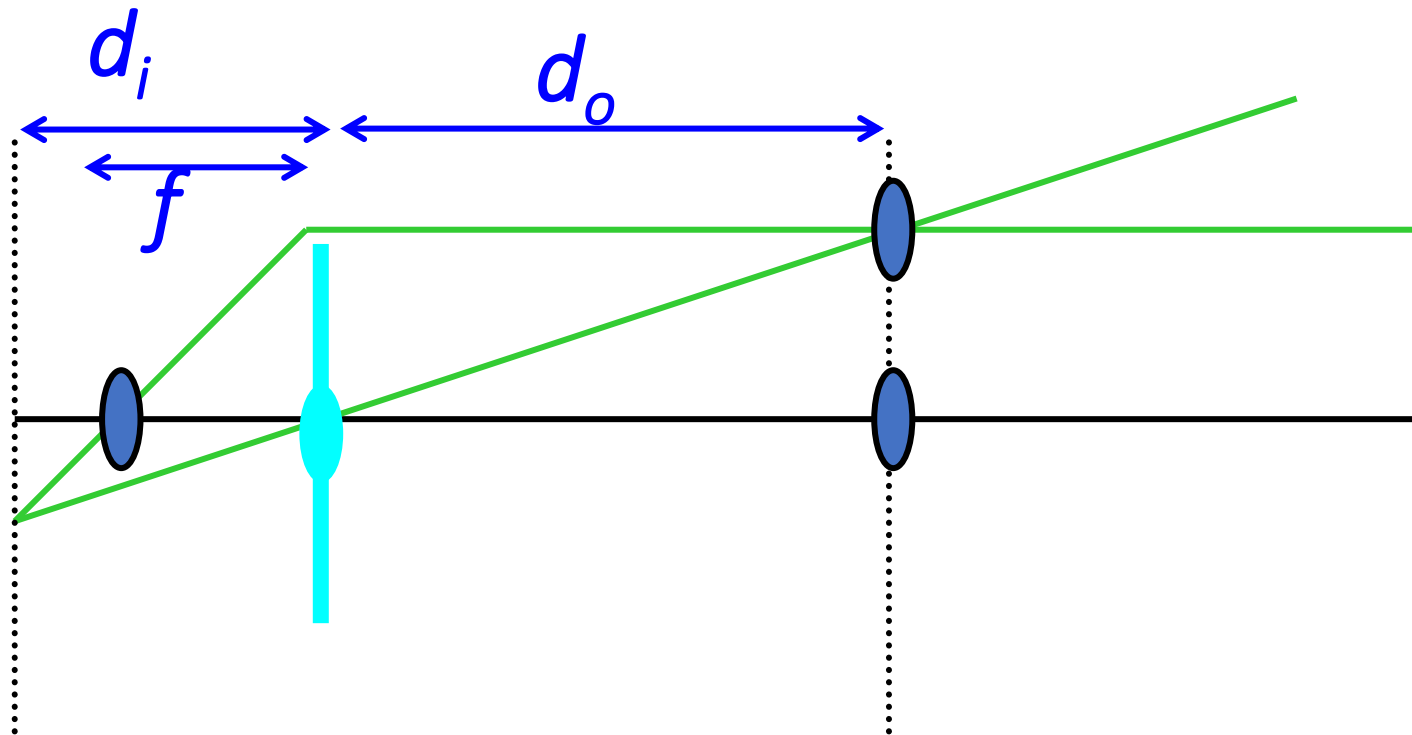


Thin lens equation: 
$$\frac{1}{d_o} + \frac{1}{d_i} = \frac{1}{f}$$

- Any object point satisfying this equation is in focus
- Thin lens applet (needs java player):  
[http://www.phy.ntnu.edu.tw/java/Lens/lens\\_e.html](http://www.phy.ntnu.edu.tw/java/Lens/lens_e.html) (by Fu-Kwun Hwang)

[Slide by Steve Seitz]

# Thin lens formula



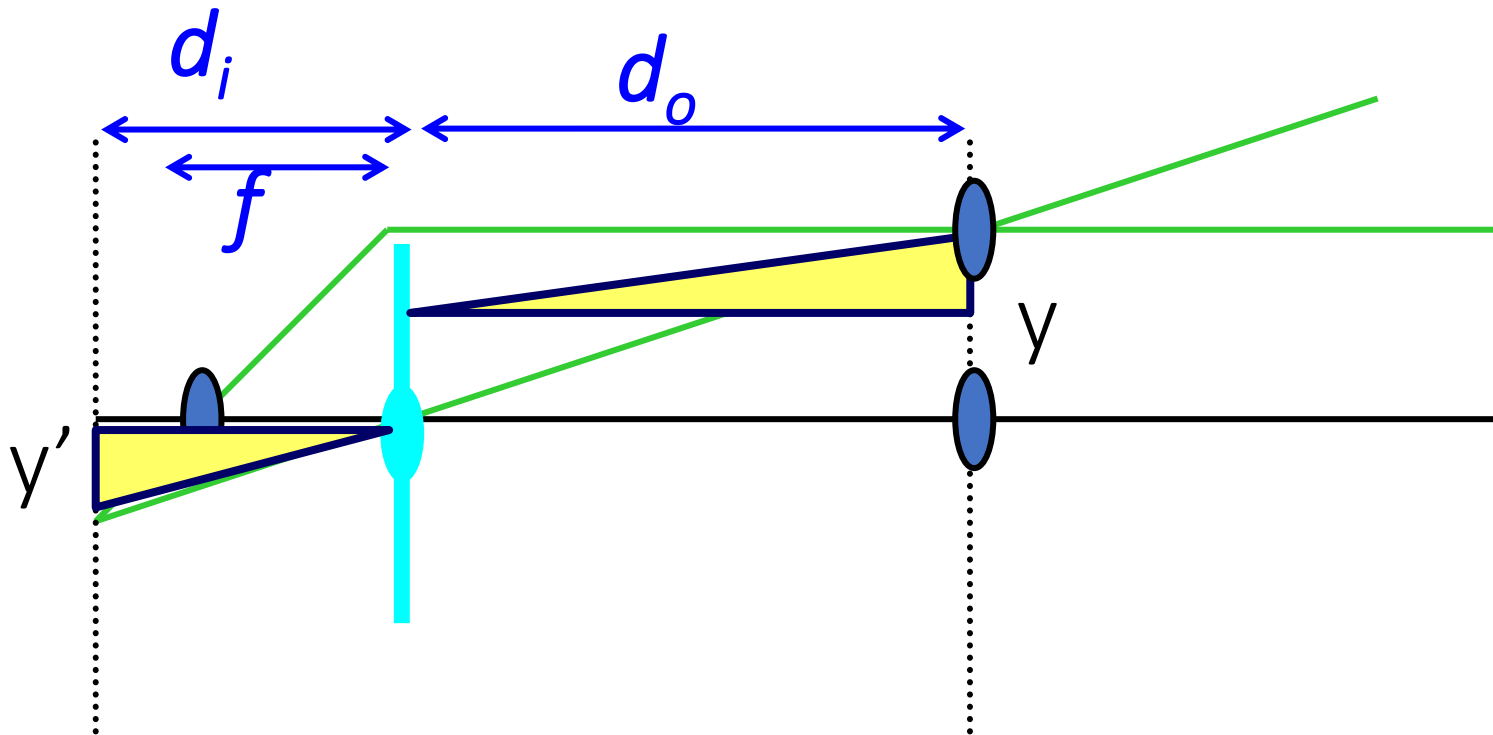
[Slide by Fredo Durand]





# Thin lens formula

$$y'/y = d_i/d_o$$

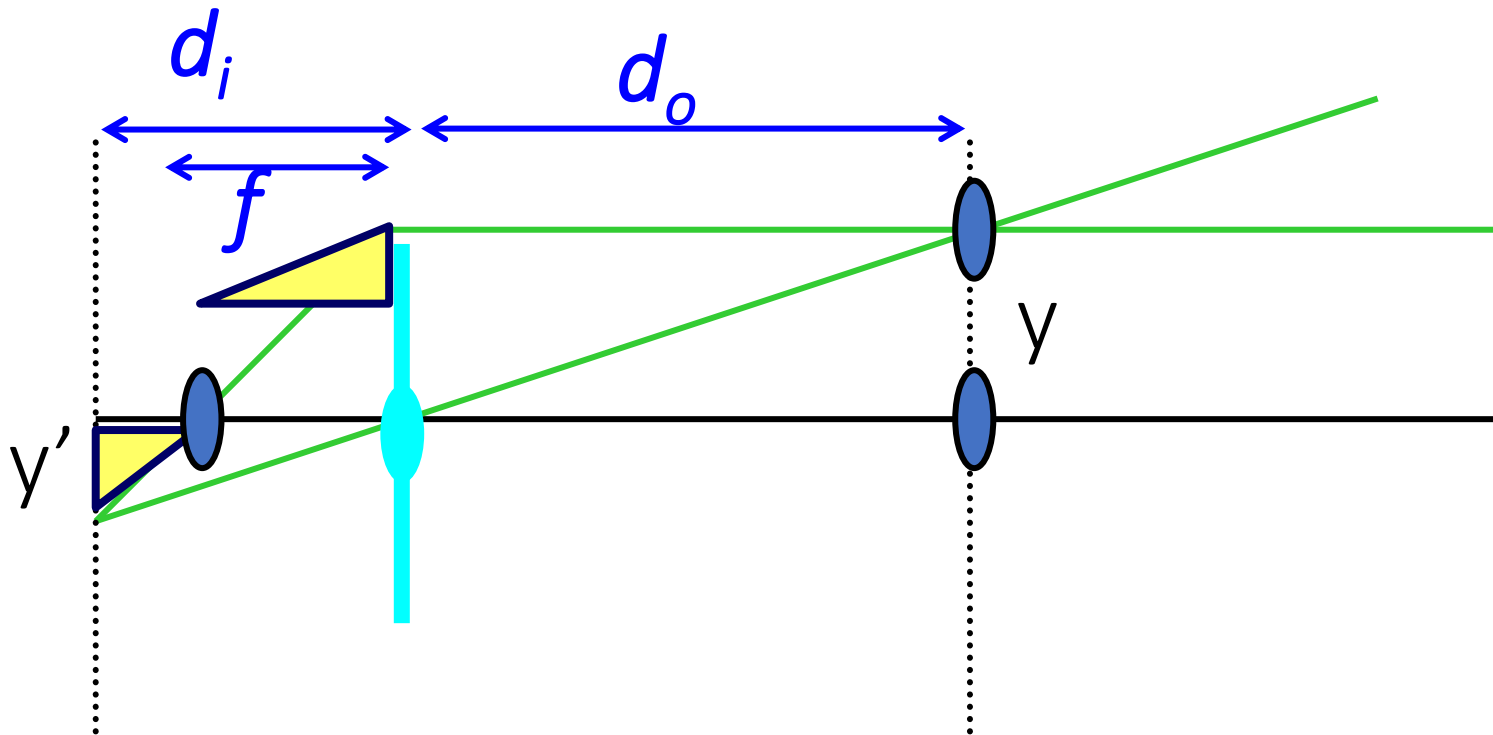


[Slide by Fredo Durand]

# Thin lens formula

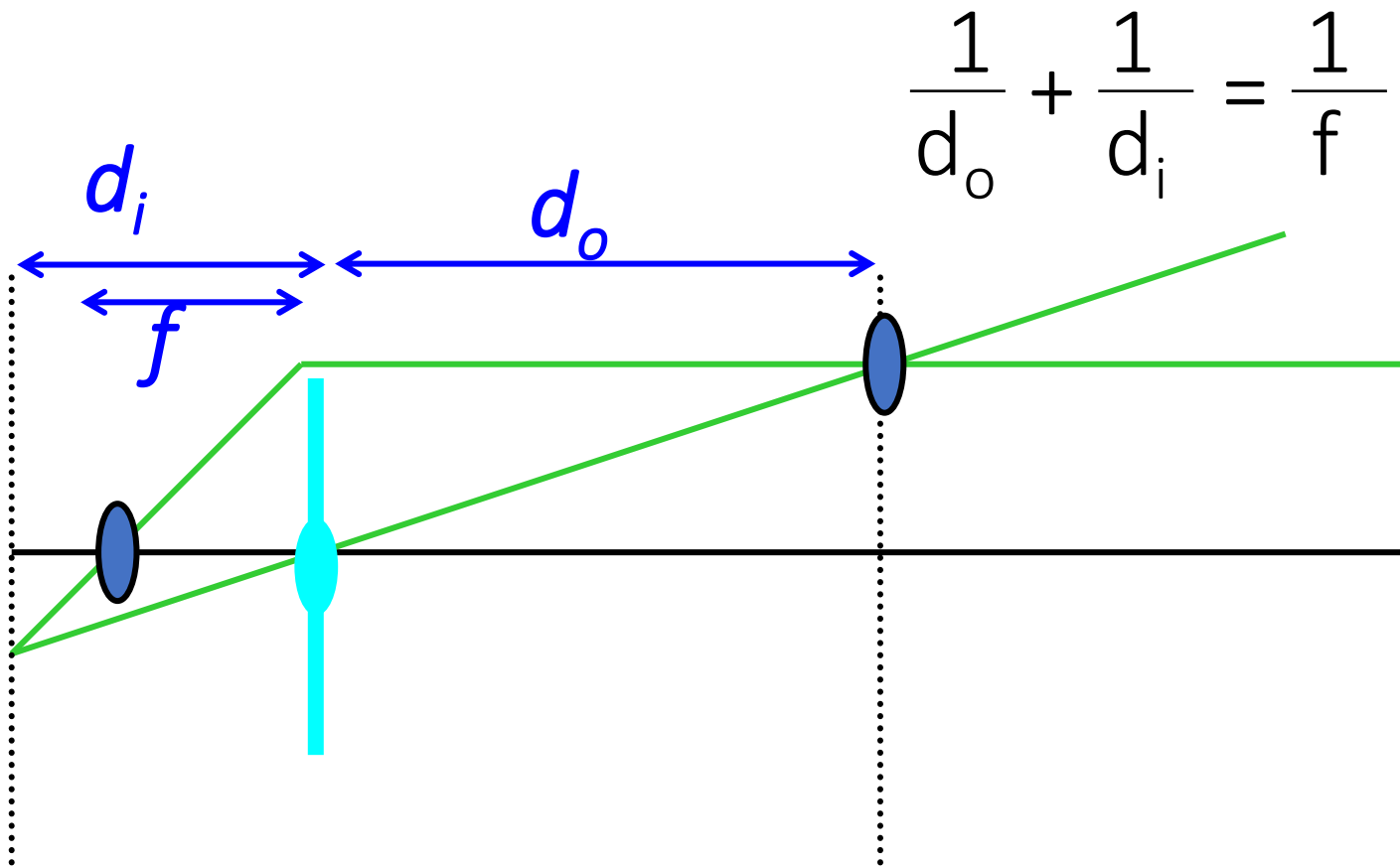
$$y'/y = d_i/d_o$$

$$y'/y = (d_i - f)/f$$



[Slide by Fredo Durand]

# Thin lens formula



Any point satisfying the thin lens equation is in focus.

[Slide by Fredo Durand]



# Today's Agenda - Overview of Cameras

- Pinhole Camera model
  - Aperture
  - Camera Obscura
- Cameras with lenses
  - Thin lens equation
  - Depth of field
  - Field of view
- Digital cameras
  - Bayer filters
  - Debayering

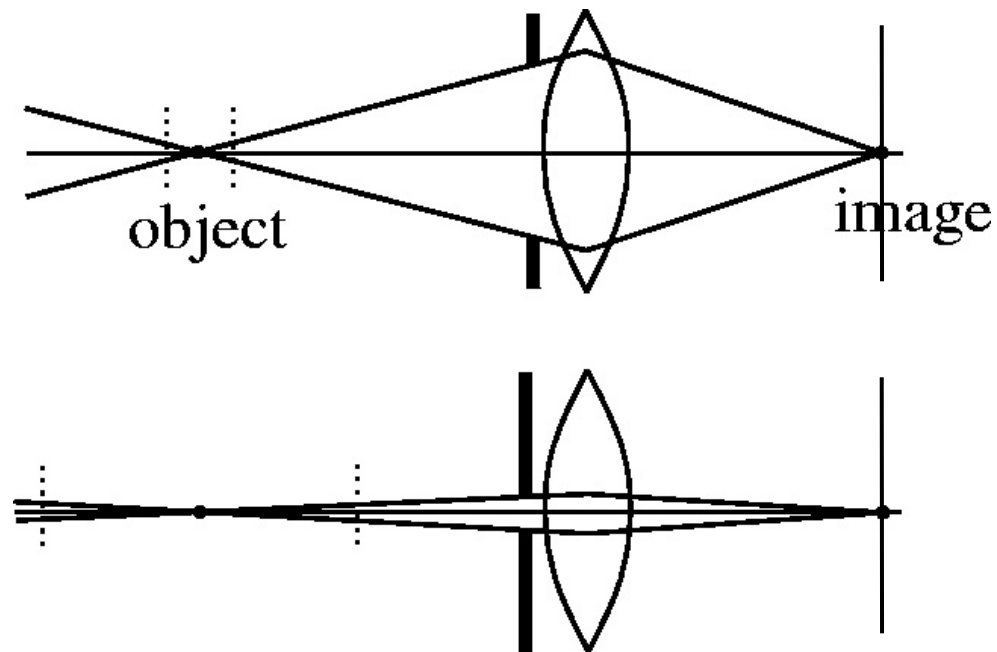
## Depth of field



<http://www.cambridgeincolour.com/tutorials/depth-of-field.htm>

[Slide by A. Efros]

# Depth of field



f / 5.6



f / 32

Changing the aperture size affects depth of field

- A smaller aperture increases the range in which the object is approximately in focus
- But small aperture reduces amount of light – need to increase exposure

Flower images from Wikipedia [http://en.wikipedia.org/wiki/Depth\\_of\\_field](http://en.wikipedia.org/wiki/Depth_of_field)

## Depth of field



Large aperture = small DOF



Small aperture = large DOF

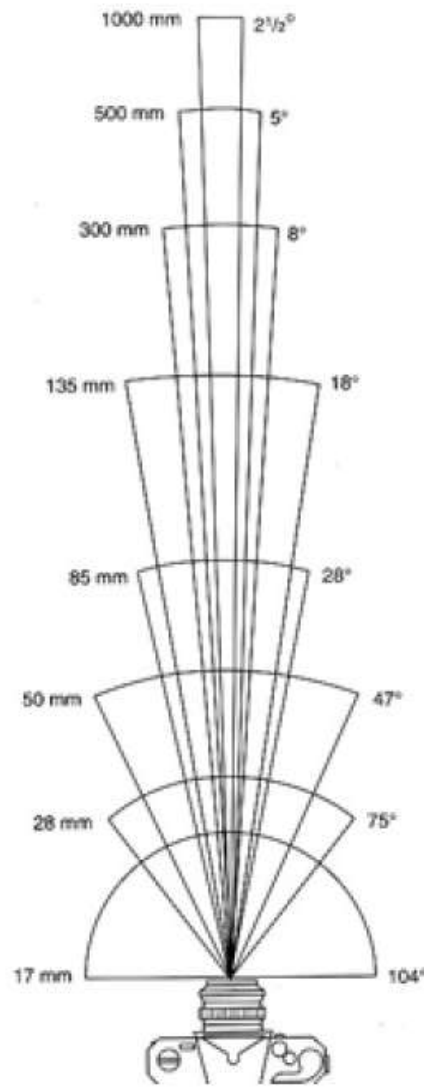
[Slide by A. Efros]



# Today's Agenda - Overview of Cameras

- Pinhole Camera model
  - Aperture
  - Camera Obscura
- Cameras with lenses
  - Thin lens equation
  - Depth of field
  - Field of view
- Digital cameras
  - Bayer filters
  - Debayering

## Field of view



17mm



28mm



50mm

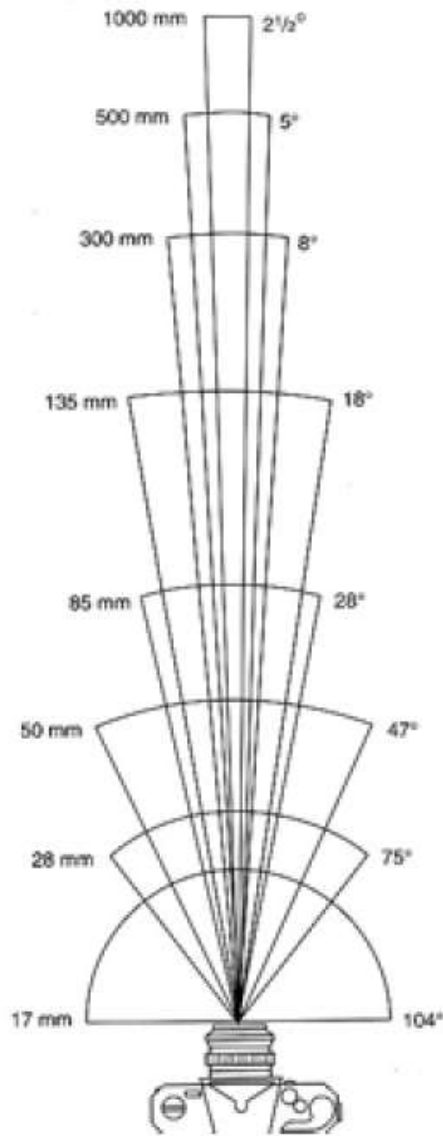


85mm

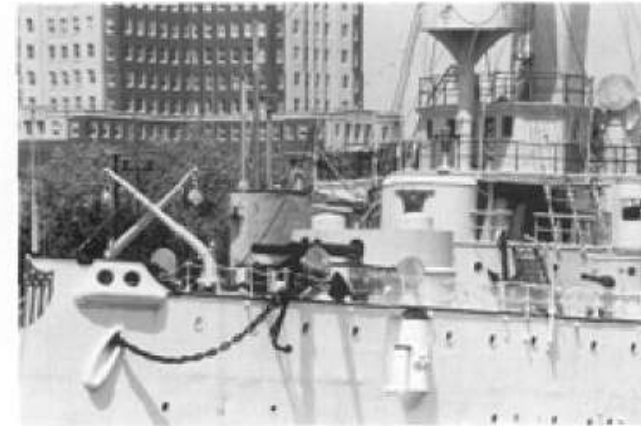
### From London and Upton

[Slide by A. Efros]

## Field of view



135mm



300mm



500mm

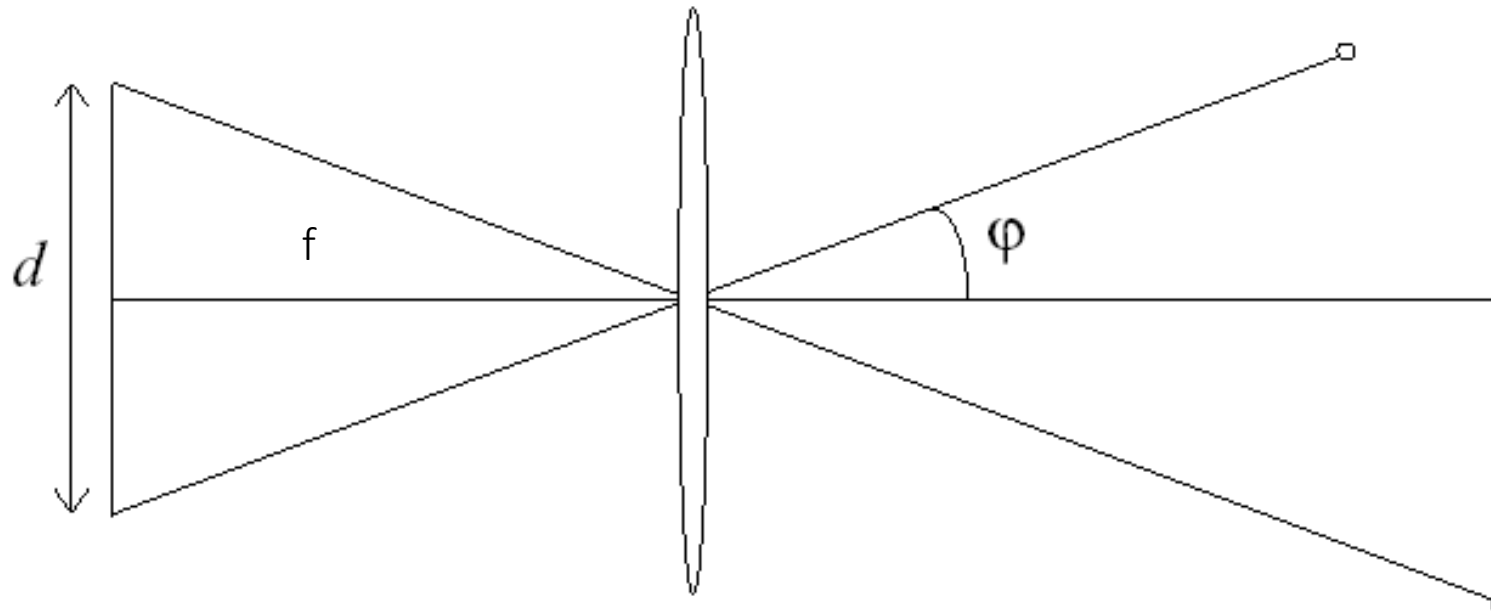


1000mm

### From London and Upton

[Slide by A. Efros]

# Field of view



FOV depends on focal length  $f$  and size of the camera sensor size  $d$

$$\varphi = \tan^{-1}\left(\frac{d}{2f}\right)$$

Smaller FOV = larger Focal Length

[Slide by A. Efros]

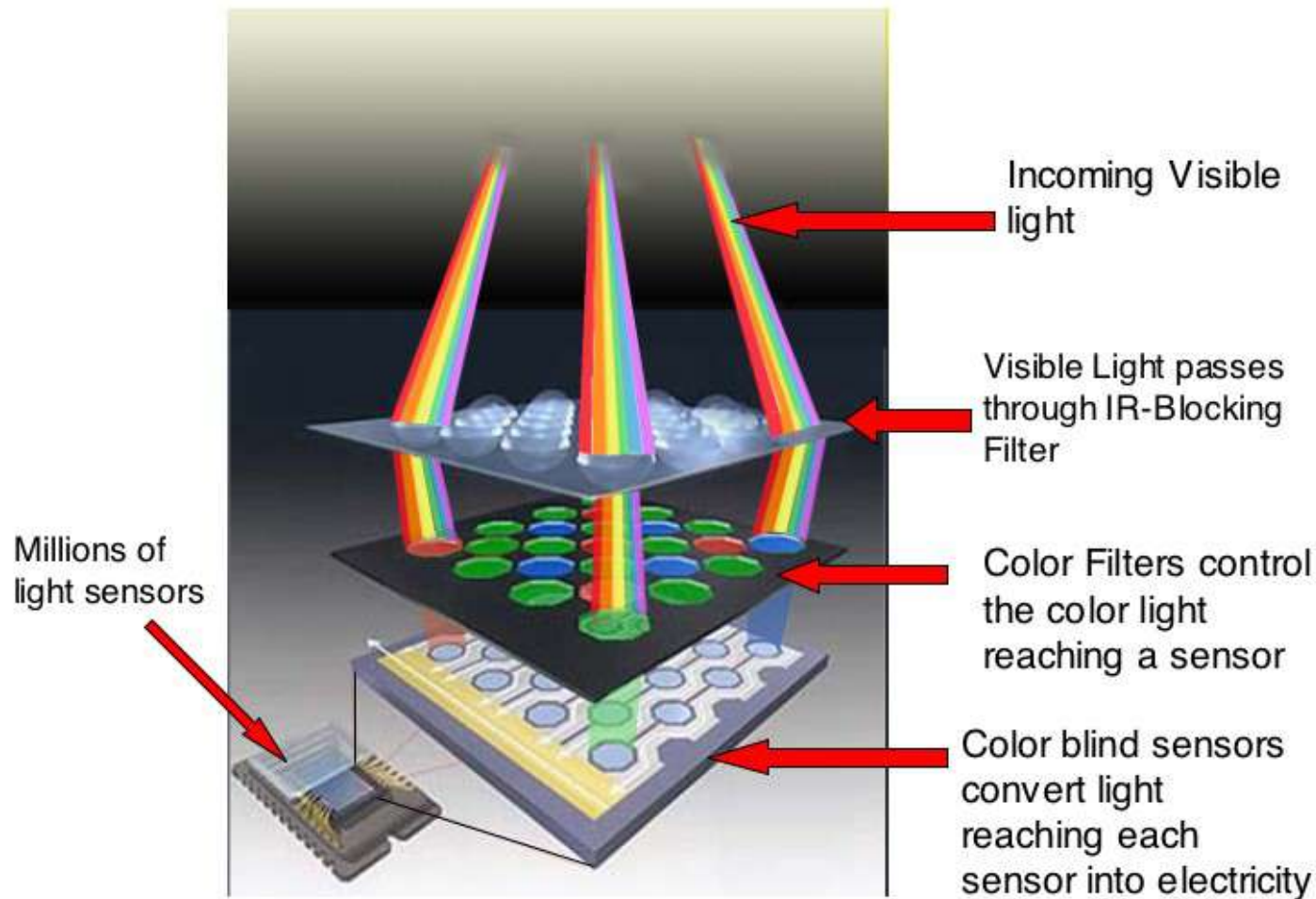


# Today's Agenda - Overview of Cameras

- Pinhole Camera model
  - Aperture
  - Camera Obscura
- Cameras with lenses
  - Thin lens equation
  - Depth of field
  - Field of view
- Digital cameras
  - Bayer filters
  - Debayering

# From light to pixels

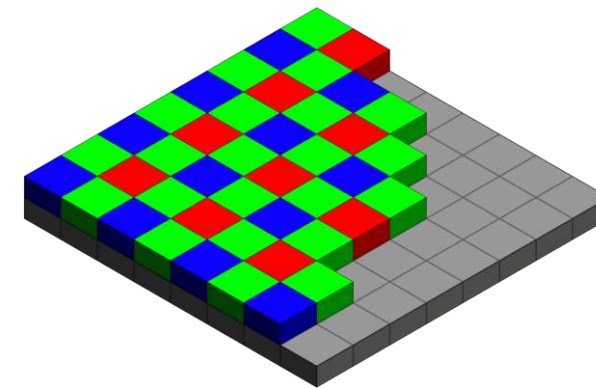
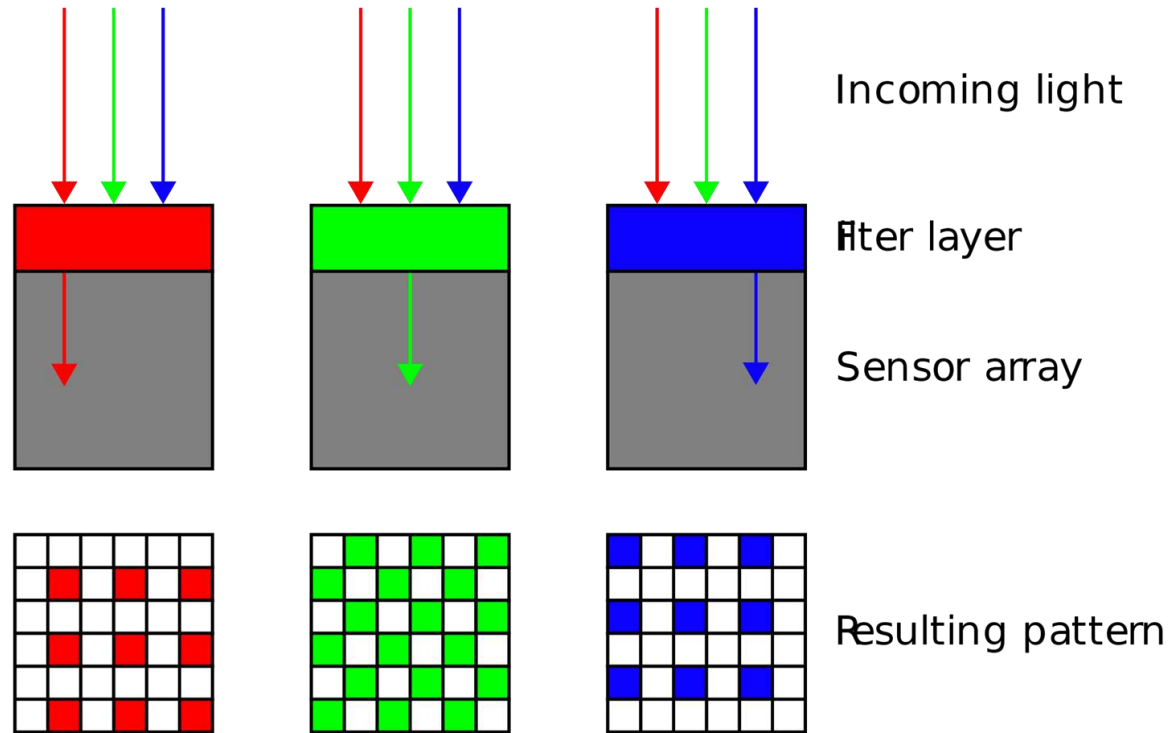
## RGB Inside the Camera



<https://sites.google.com/a/globalsystemscience.org/digital-earth-watch/tools/digital-cameras/light-entering-a-camera>

[Slide by Steve Seitz]

# Bayer filters



[https://en.wikipedia.org/wiki/Bayer\\_filter](https://en.wikipedia.org/wiki/Bayer_filter)

1/4 of pixels see red light (e.g.)

- Q: how do you get red at every pixel?
- A: Need to interpolate -- called *debayering*

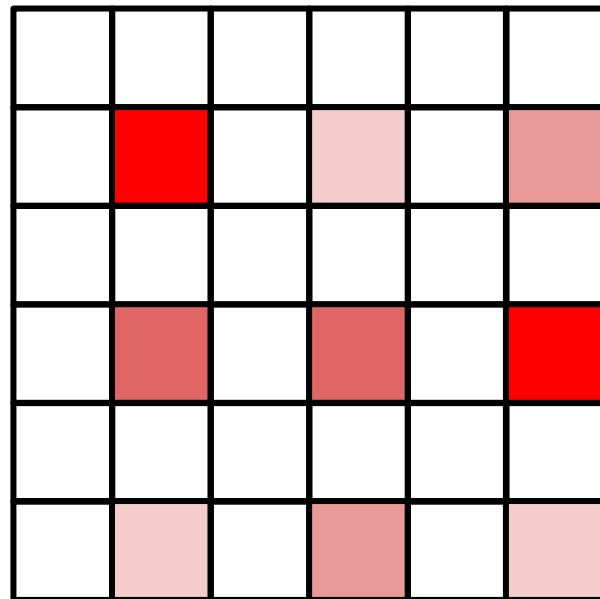
[Slide by Steve Seitz]

# Today's Agenda - Overview of Cameras

- Pinhole Camera model
  - Aperture
  - Camera Obscura
- Cameras with lenses
  - Thin lens equation
  - Depth of field
  - Field of view
- Digital cameras
  - Bayer filters
  - Debayering



# Debayering



	100		10		30
	50		50		100
	10		30		10

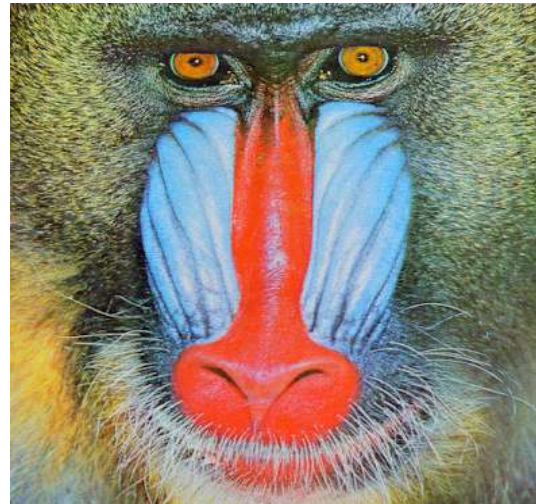
1/4 of pixels see red light (e.g.)

- Q: how do you get red at every pixel?
- A: Need to interpolate -- called *debayering*

[Slide by Steve Seitz]

## RGB images (three channel)

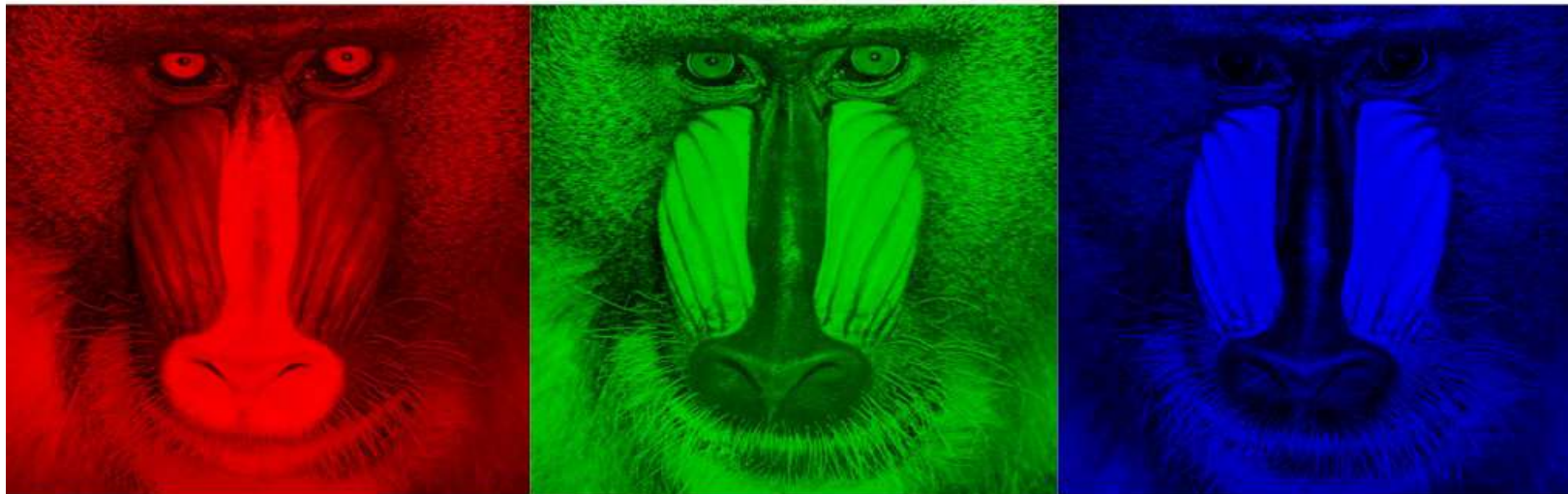
What we see



What we get out of the camera

[Slide by Steve Seitz]

## From now on: what to do with these RGB images



[Slide by Steve Seitz]

**MAI4CAREU**

Master programmes in Artificial  
Intelligence 4 Careers in Europe



**CYENS**  
CENTRE OF EXCELLENCE



# Thank you.







University of Cyprus – MSc Artificial Intelligence

# MAI644 – COMPUTER VISION

## Lecture 4: Interpolation – Resizing

**Melinos Averkiou**

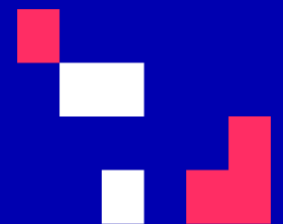
CYENS Centre of Excellence

University of Cyprus - Department of Computer Science

[m.averkiou@cyens.org.cy](mailto:m.averkiou@cyens.org.cy)



**CYENS**  
CENTRE OF EXCELLENCE



# Last time

- Pinhole Camera model
  - Aperture
  - Camera Obscura
- Cameras with lenses
  - Thin lens equation
  - Depth of field
  - Field of view
- Digital cameras
  - Bayer filters
  - Debayering

# Today's Agenda

- Image basics
  - What is an image – addressing pixels
  - Image as a function – image coordinates
- Image interpolation
  - Nearest neighbor
  - Bilinear
  - Bicubic
- Image resizing
  - Enlarge
  - Shrink

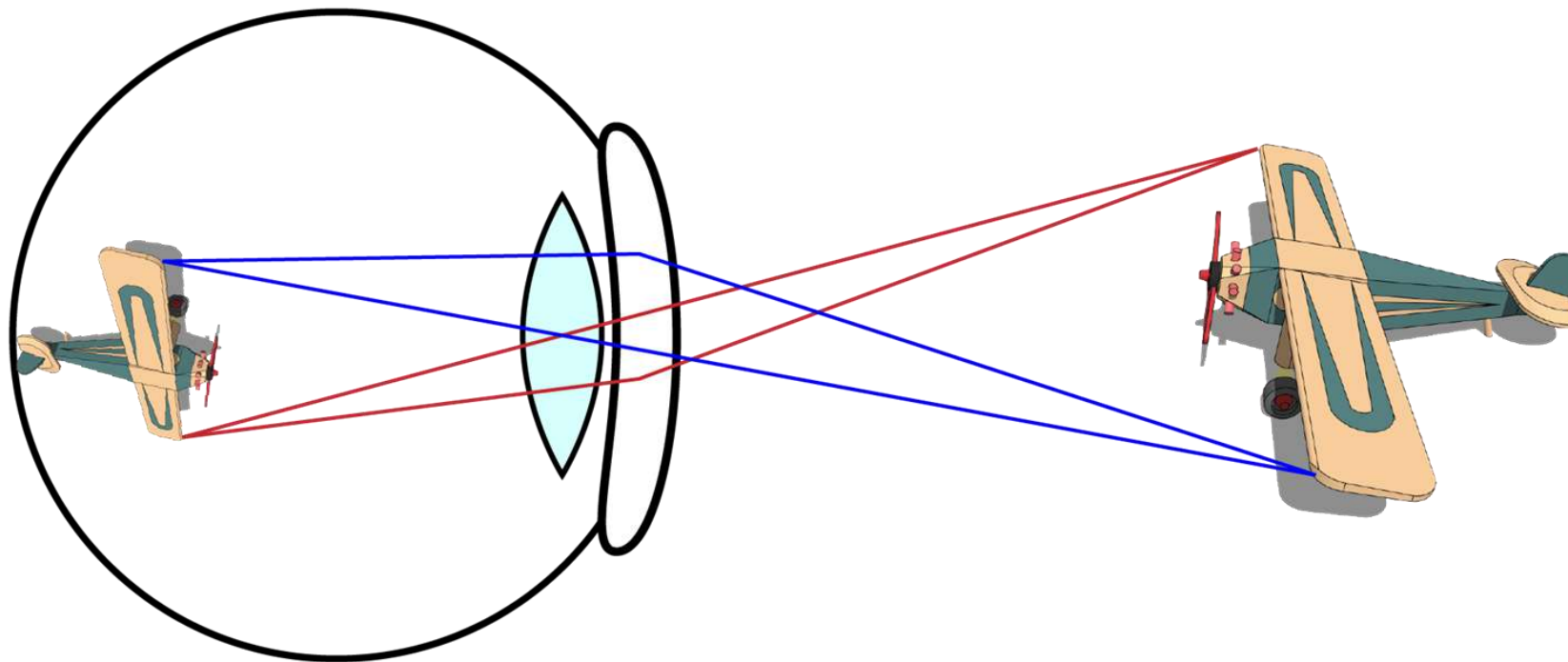
**[material based on Joseph Redmon's course]**

# Today's Agenda

- Image basics
  - What is an image – addressing pixels
  - Image as a function – image coordinates
- Image interpolation
  - Nearest neighbor
  - Bilinear
  - Bicubic
- Image resizing
  - Enlarge
  - Shrink

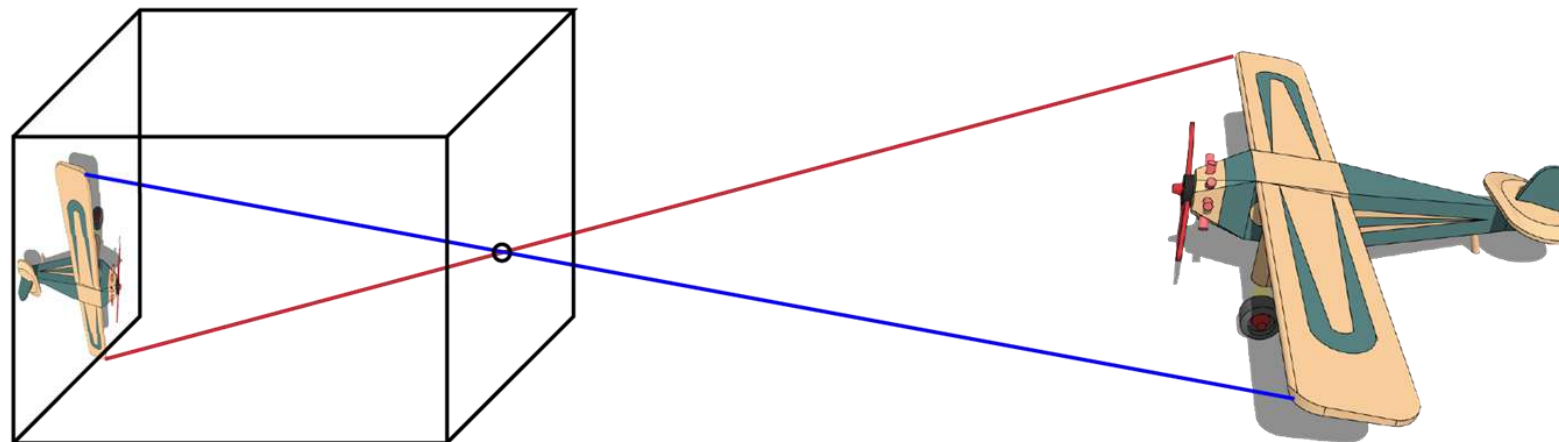


## Eyes: projection onto retina



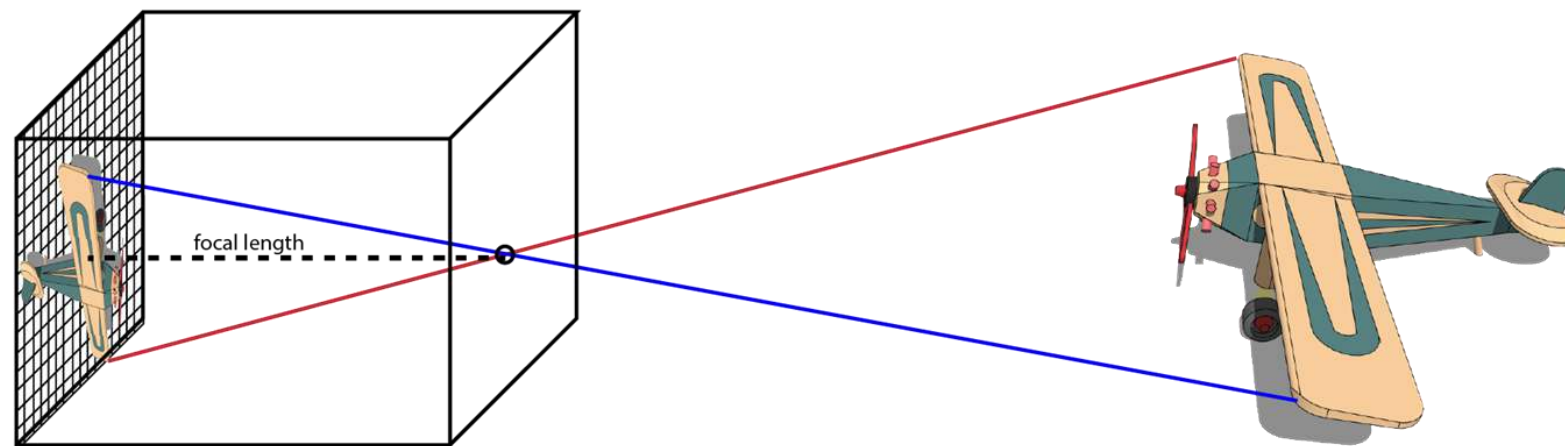


## Model: pinhole camera

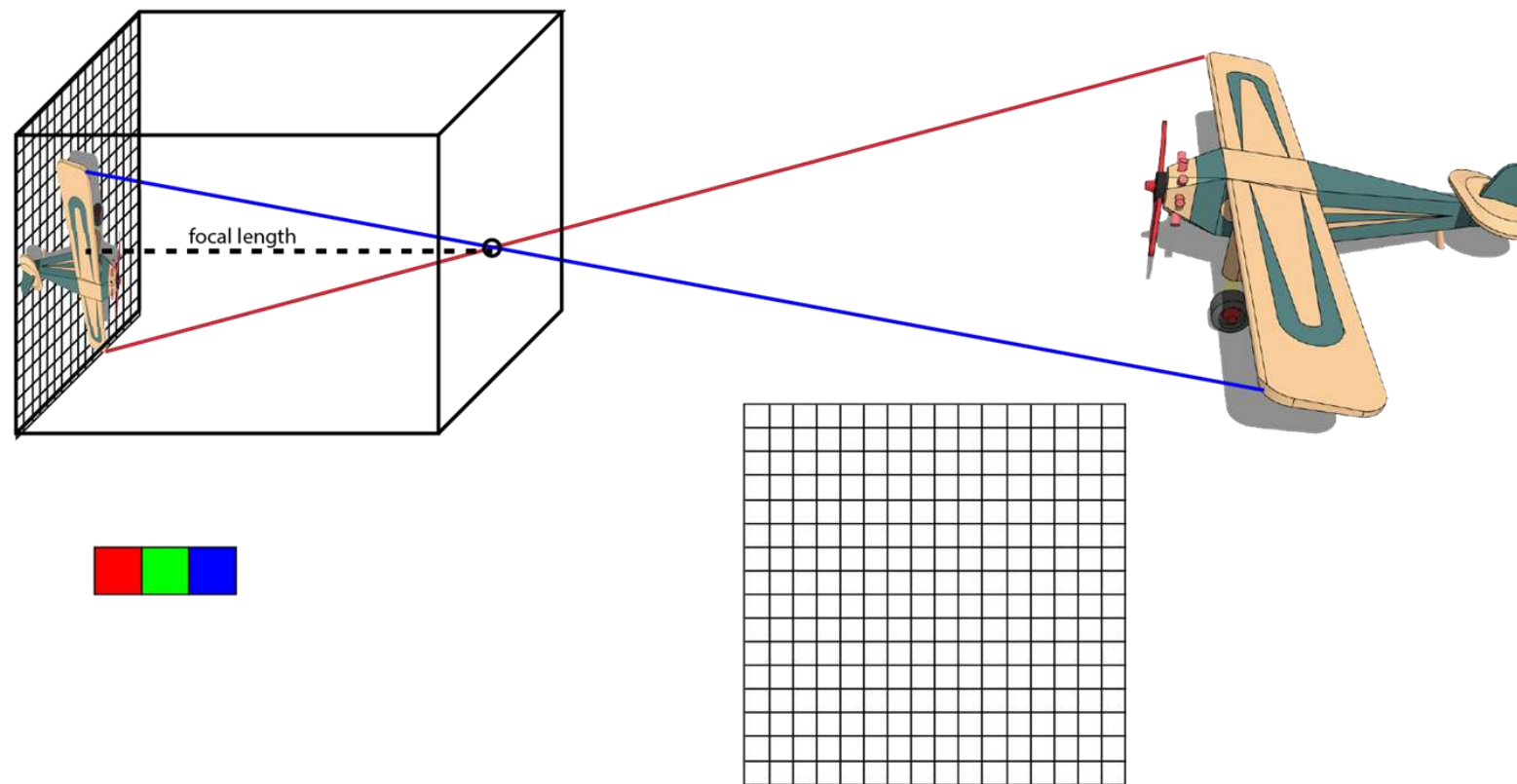




## At each point we record incident light

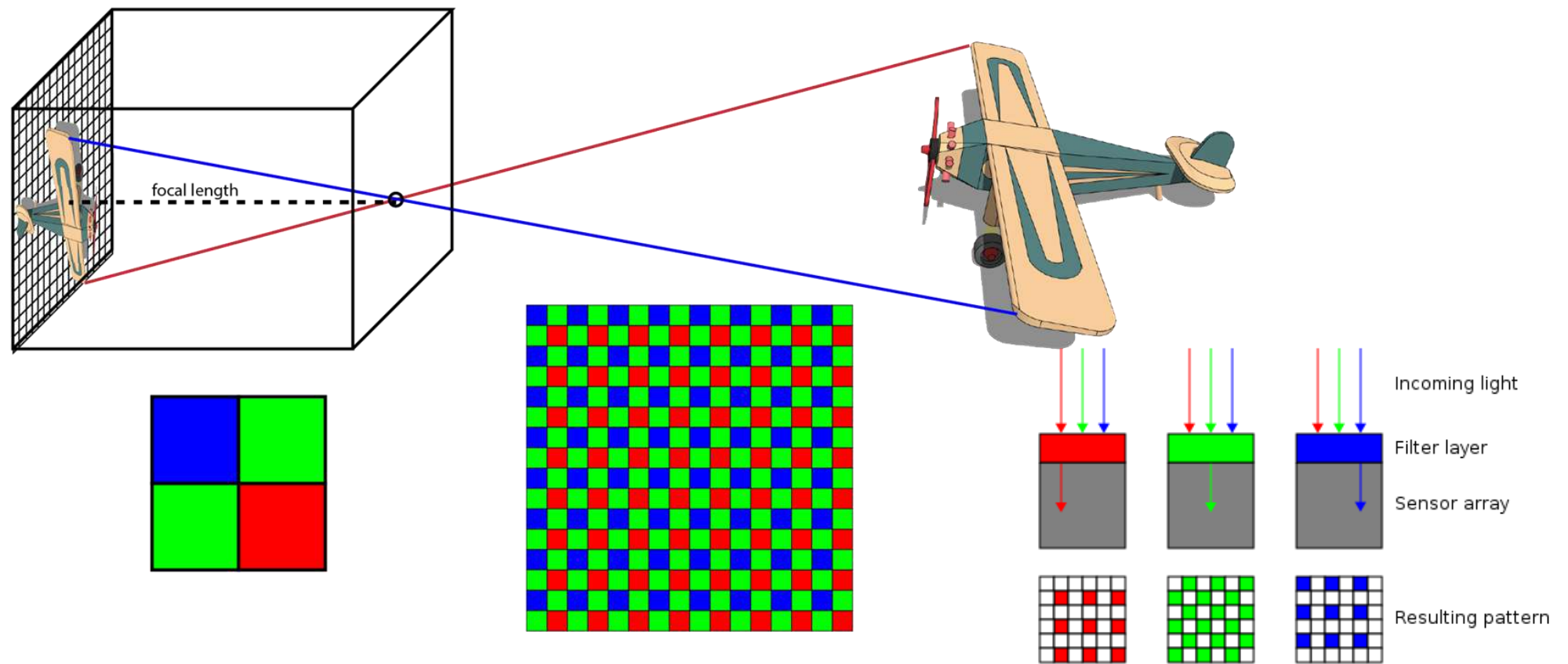


## How do we record color?

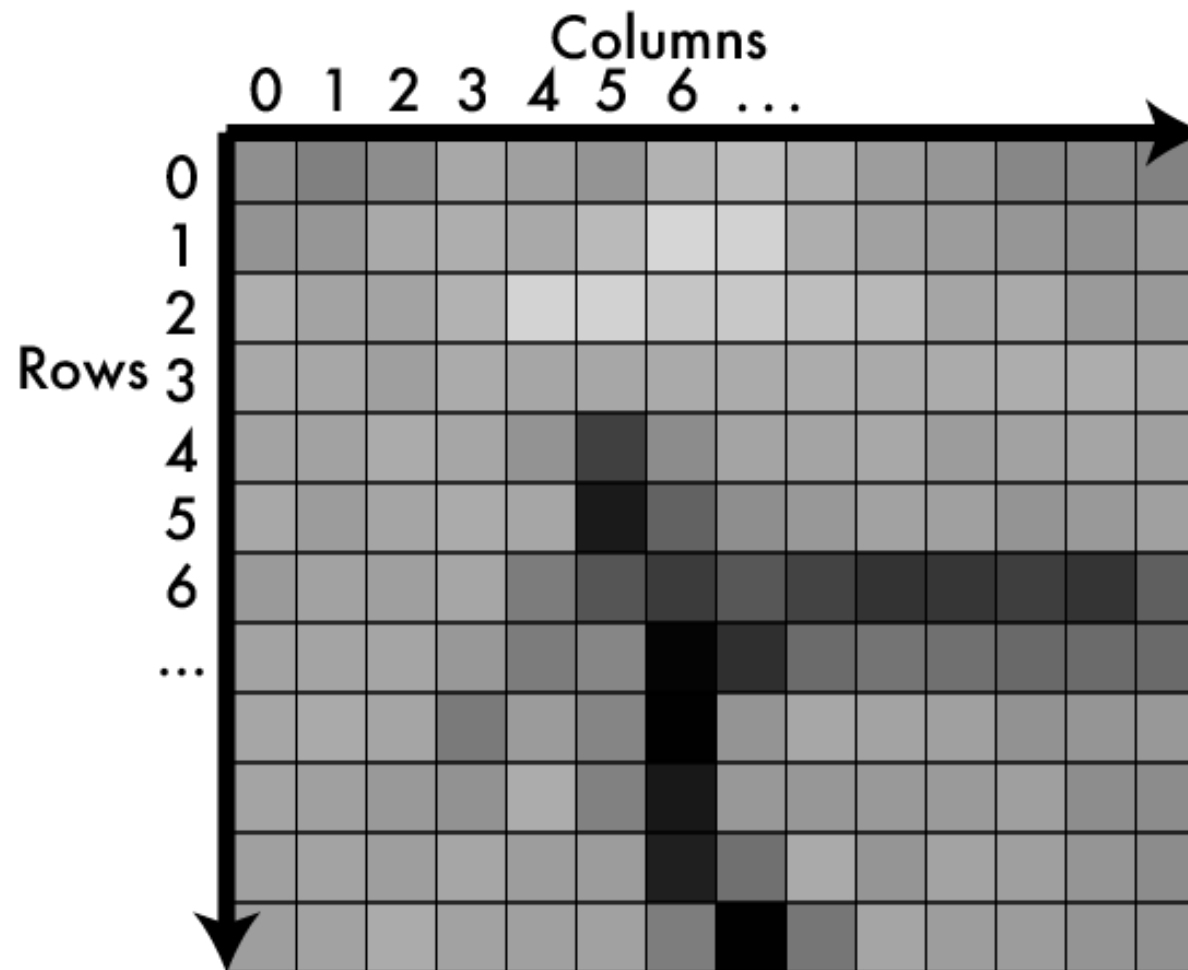




# Bayer pattern for CMOS sensors



# An image is a matrix of light



## Values in matrix = how much light

		Columns													
		0	1	2	3	4	5	6	...						
Rows	0	100	102	107	102	132	146	136	156	148	122	115	104	105	103
	1	100	102	107	102	132	146	136	156	148	122	115	104	105	103
	2	100	102	107	102	132	146	136	156	148	122	115	104	105	103
	3	100	102	107	102	132	146	136	156	148	122	115	104	105	103
	4	100	102	107	102	132	146	136	156	148	122	115	104	105	103
	5	100	102	107	102	132	30	60	156	148	122	115	104	105	103
	6	100	102	107	102	132	40	20	50	32	20	20	24	30	62
	...	100	102	107	102	132	71		156	51	57	57	58	62	58
	100	102	107	102	132	69		156	148	122	115	104	105	103	
	100	102	107	102	132	89	12	156	148	122	115	104	105	103	
	100	102	107	102	132	146	13	45	148	122	115	104	105	103	
	100	102	107	102	132	146	46		42	122	115	104	105	103	

## Values in matrix = how much light

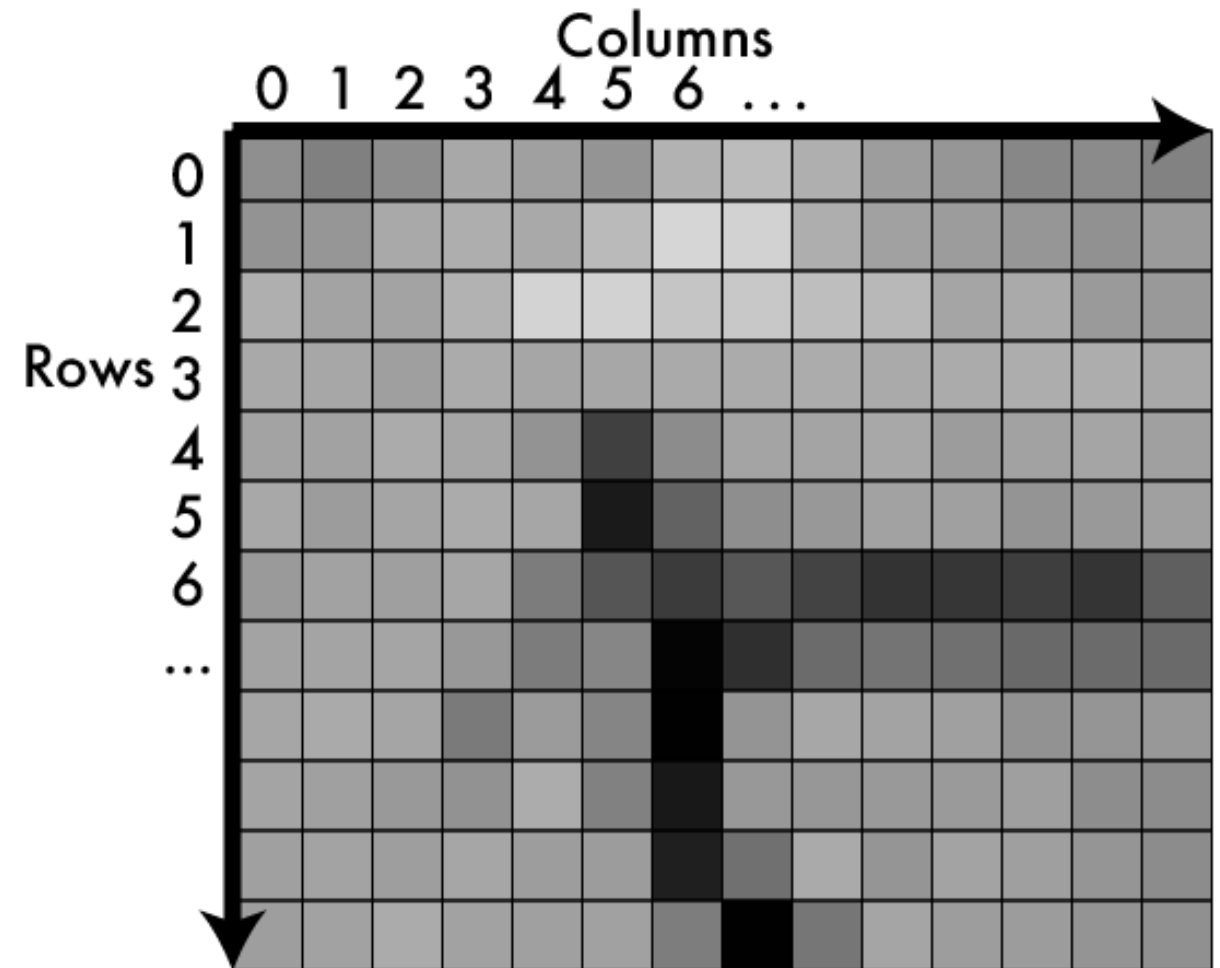
- Higher = more light
- Lower = less light
- Bounded
  - No light = 0
  - Sensor/device limit = max
  - Typical ranges:
    - [0-255], fit into byte
    - [0-1], floating point
- Called pixels

		Columns													
		0	1	2	3	4	5	6	...						
Rows	0	100	102	107	102	132	146	136	156	148	122	115	104	105	103
	1	100	102	107	102	132	146	136	156	148	122	115	104	105	103
	2	100	102	107	102	132	146	136	156	148	122	115	104	105	103
	3	100	102	107	102	132	146	136	156	148	122	115	104	105	103
	4	100	102	107	102	132	146	136	156	148	122	115	104	105	103
	5	100	102	107	102	132	30	60	156	148	122	115	104	105	103
	6	100	102	107	102	132	40	20	50	32	20	20	24	30	62
	...	100	102	107	102	132	71		156	51	57	57	58	62	58
		100	102	107	102	132	69		156	148	122	115	104	105	103
		100	102	107	102	132	89	12	156	148	122	115	104	105	103
	100	102	107	102	132	146	13	45	148	122	115	104	105	103	
	100	102	107	102	132	146	46		42	122	115	104	105	103	

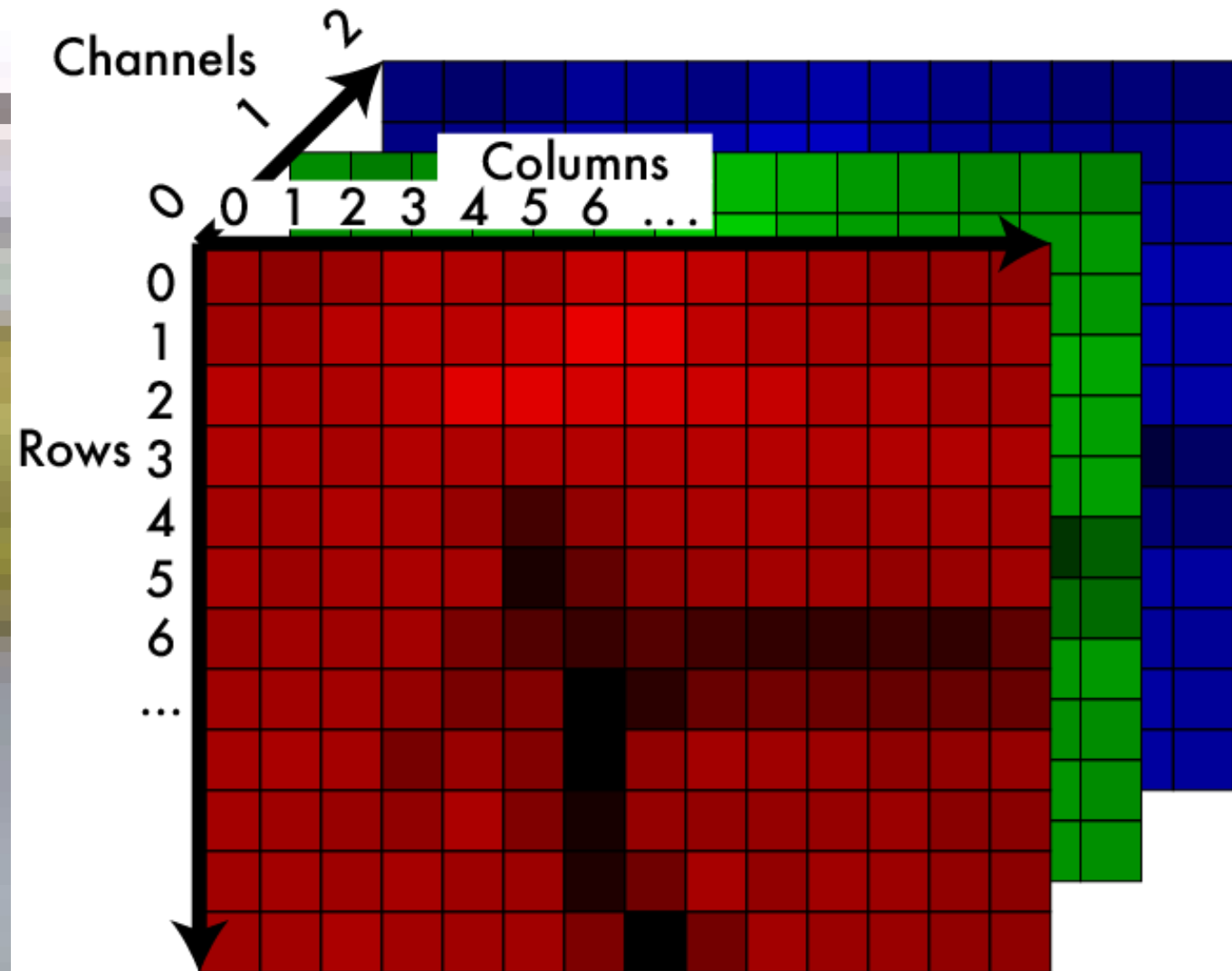


# Addressing pixels

- Ways to index:
  - (x,y)
    - Like cartesian coordinates
    - (3,6) is column 3 row 6
  - (r,c)
    - Like matrix notation
    - (3,6) is row 3 column 6
- We use (x,y)
  - Arbitrary
  - Only thing that matters is consistency



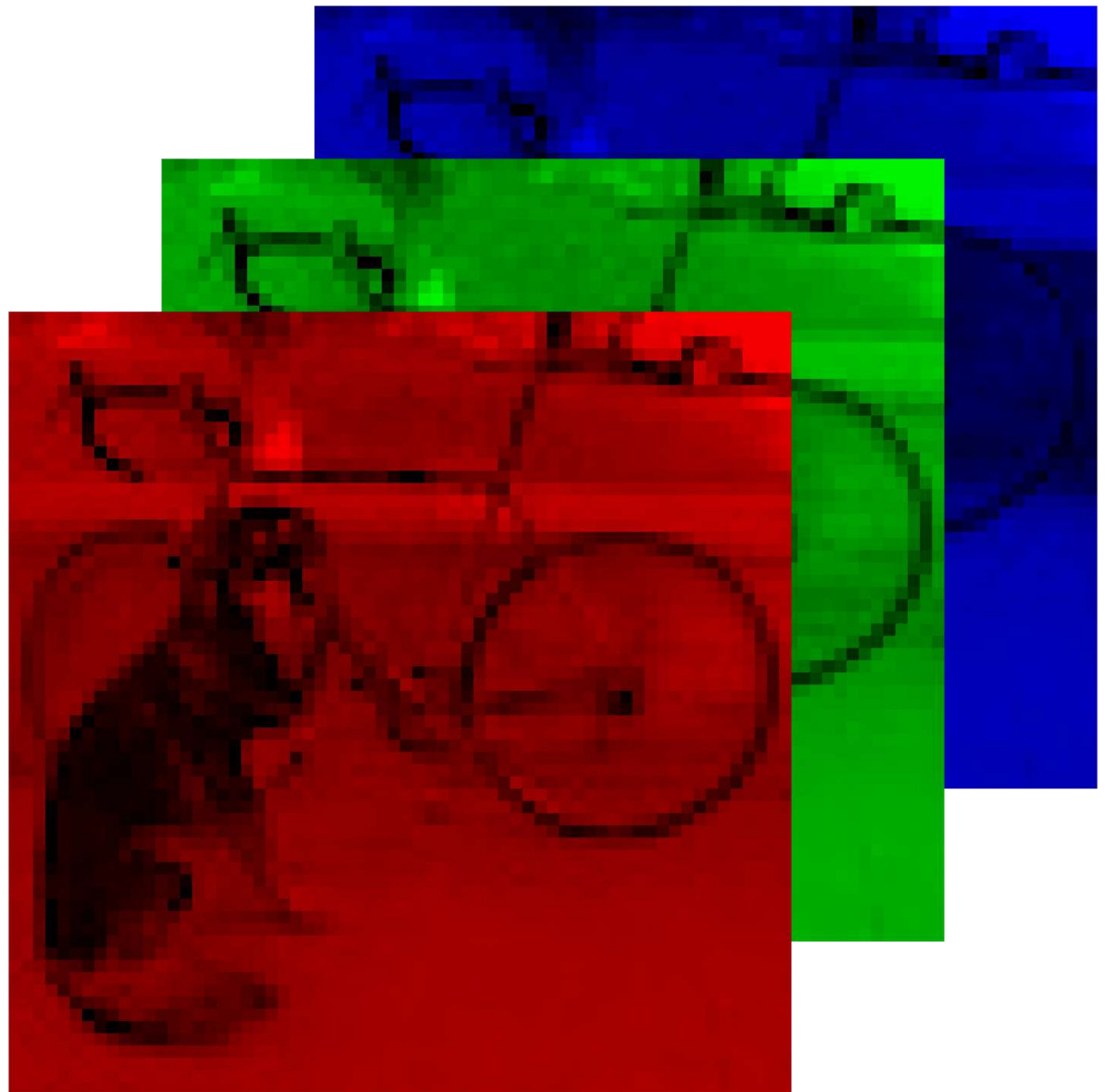
# Color image: 3d tensor in colorspace



# RGB information in separate “channels”

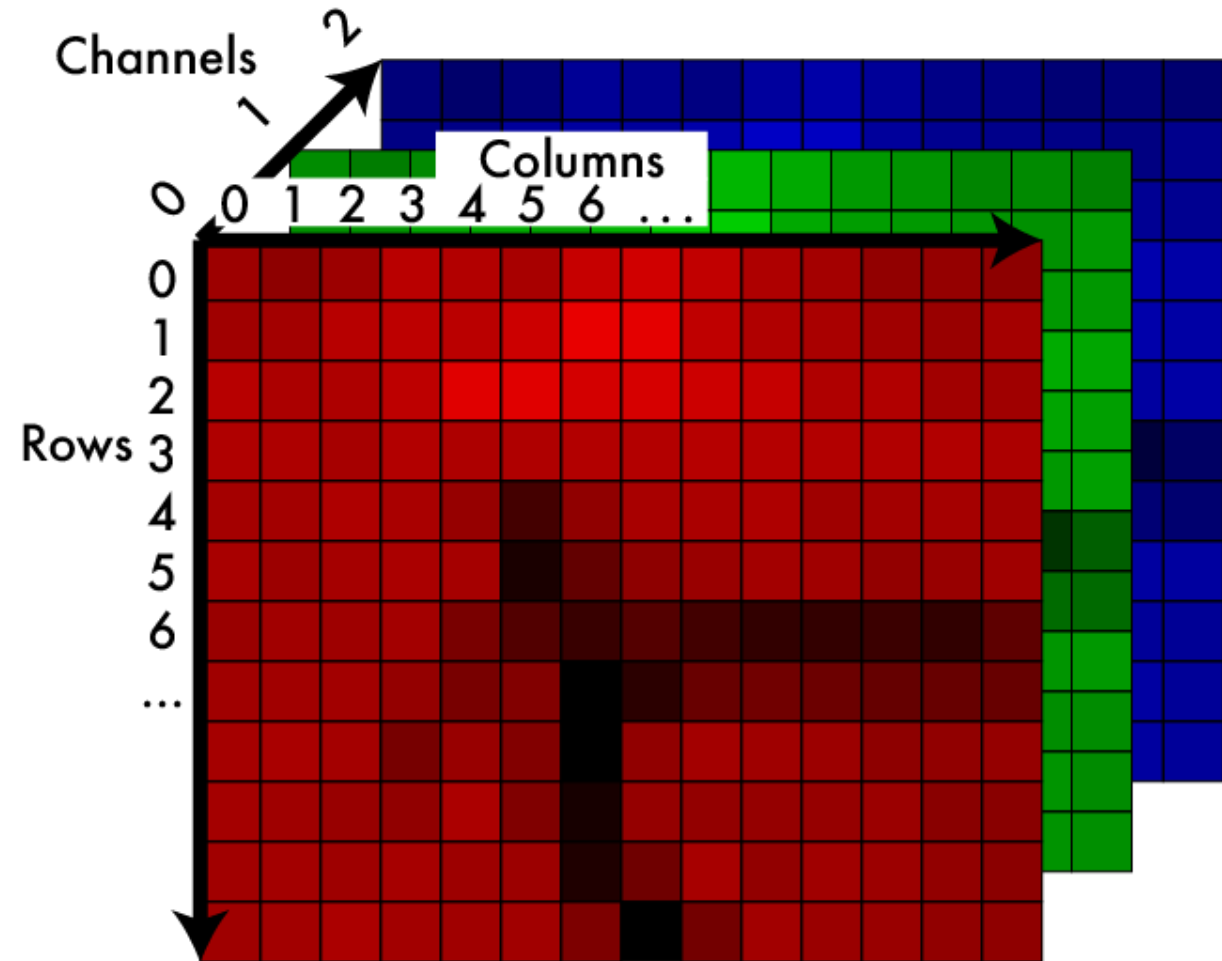
Remember: we can match “real” colors using a mix of primaries.

Each channel encodes one primary. Adding the light produced from each primary mimics the original color.



# Addressing pixels

- We use  $(x,y,c)$ 
  - $(1,2,0)$ :
    - column 1, row 2, channel 0
- Still doesn't matter, just be consistent
- Also for size:
  - 1920 x 1080 x 3 image:
    - 1920 px wide
    - 1080 px tall
    - 3 channels





# Today's Agenda

- Image basics
  - What is an image – addressing pixels
  - Image as a function – image coordinates
- Image interpolation
  - Nearest neighbor
  - Bilinear
  - Bicubic
- Image resizing
  - Enlarge
  - Shrink

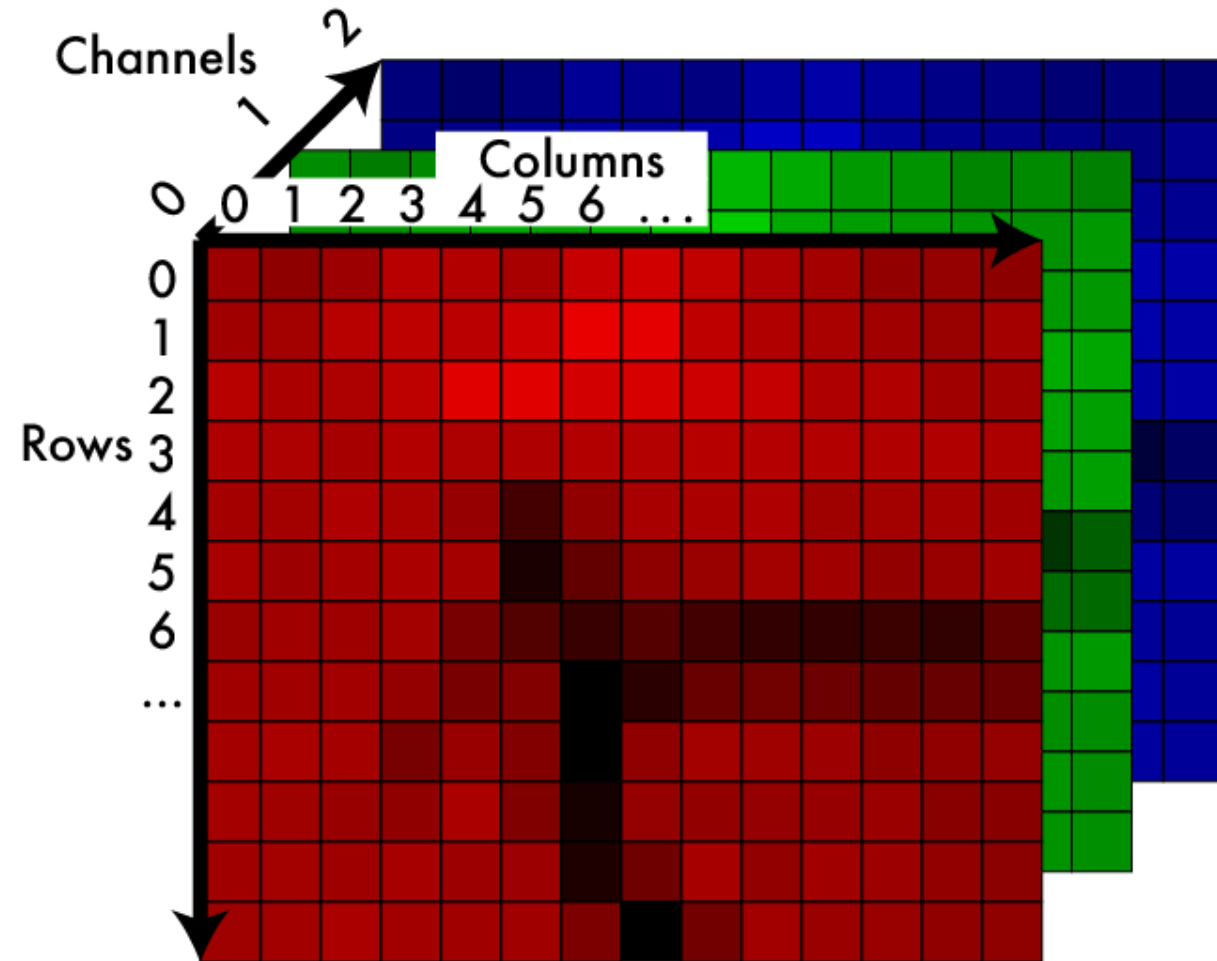
# An image is like a function

An image is a mapping from indices to pixel value:

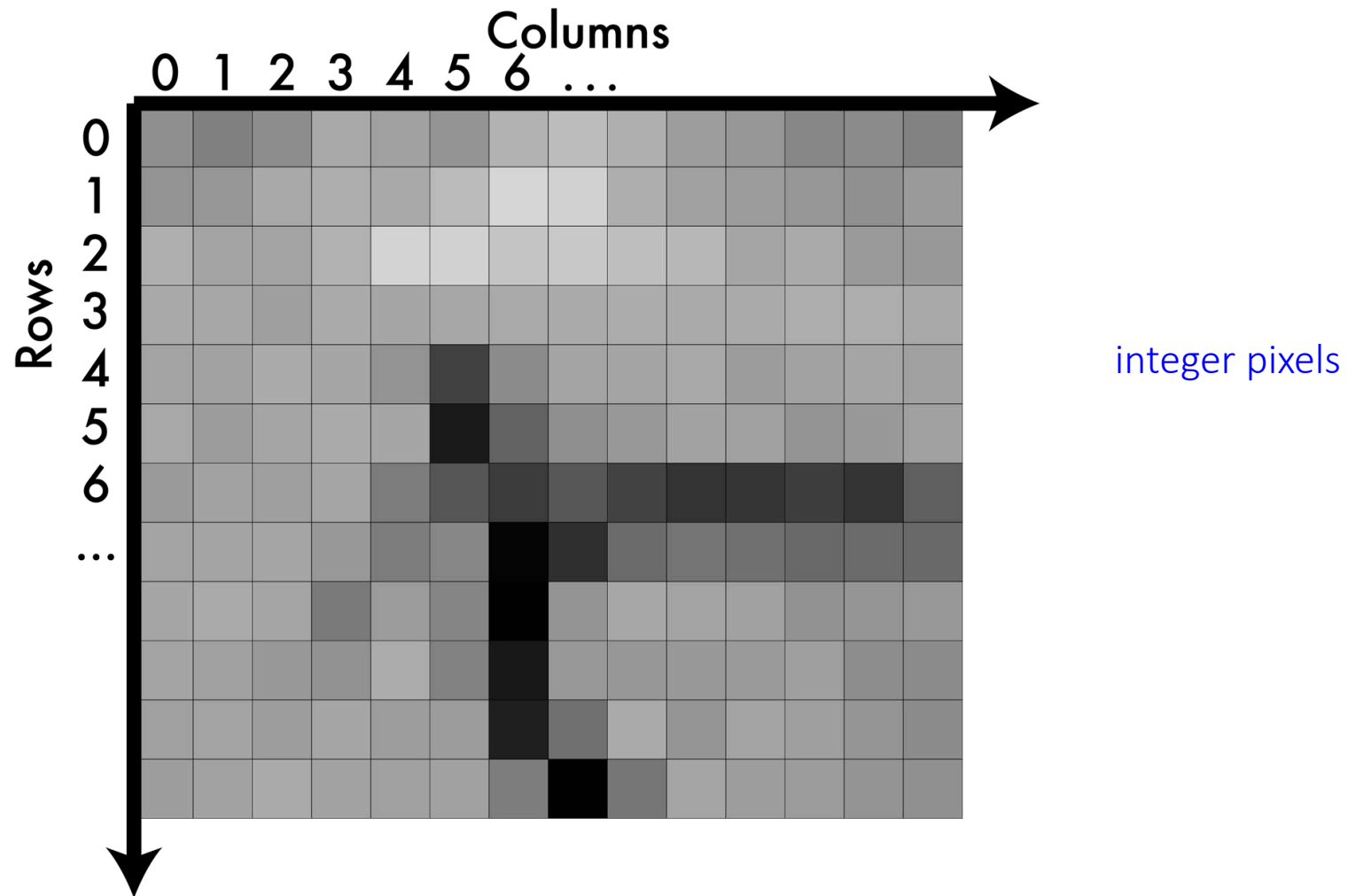
- $Im: I \times I \times I \rightarrow R$

We may want to pass in non-integers:

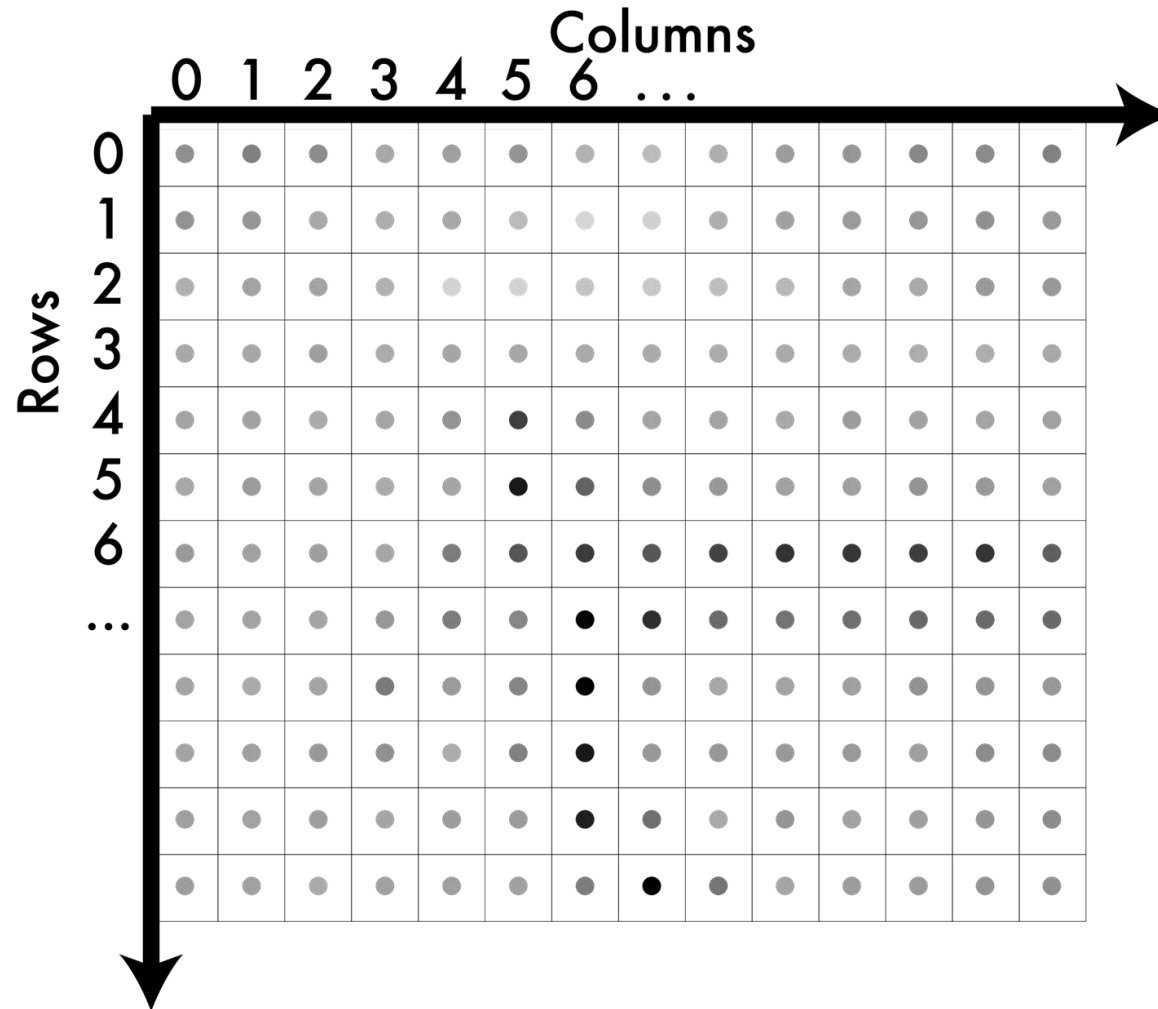
- $Im': R \times R \times I \rightarrow R$



## A note on coordinates in images



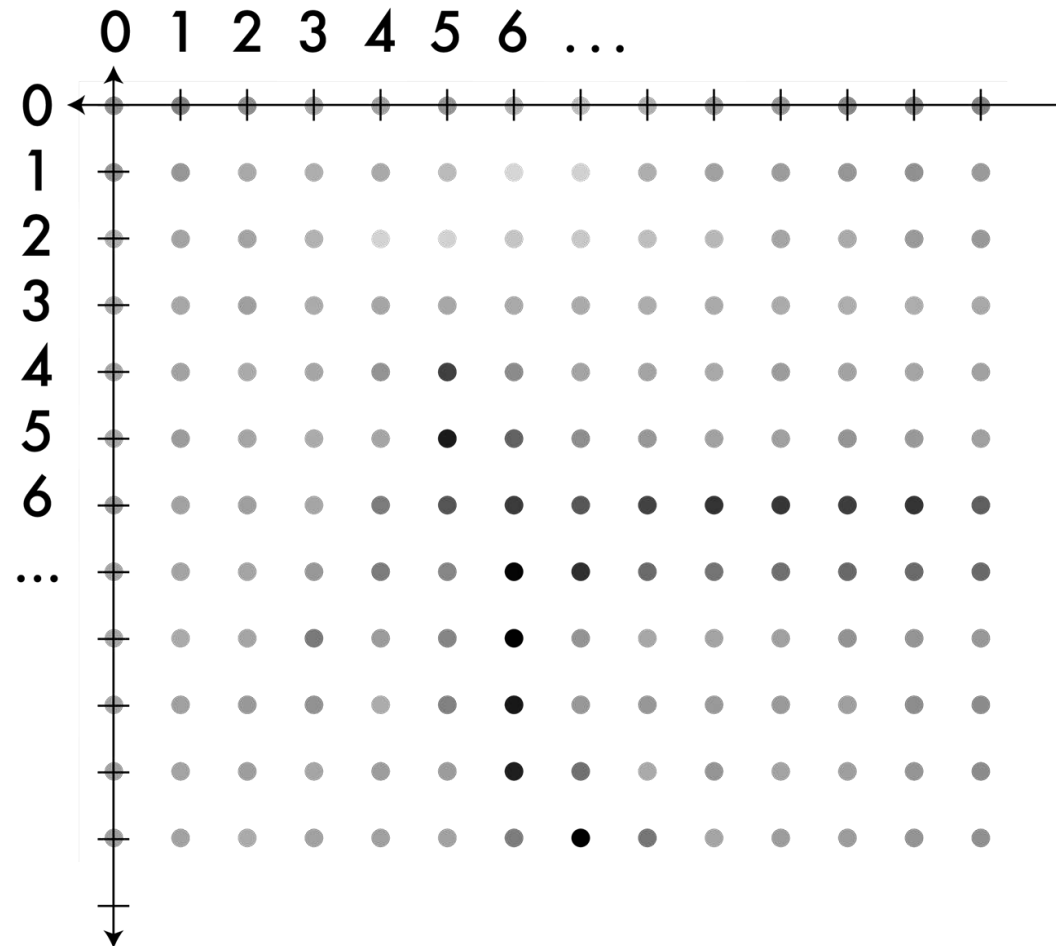
# A note on coordinates in images



We can think of their values as being at the centers.

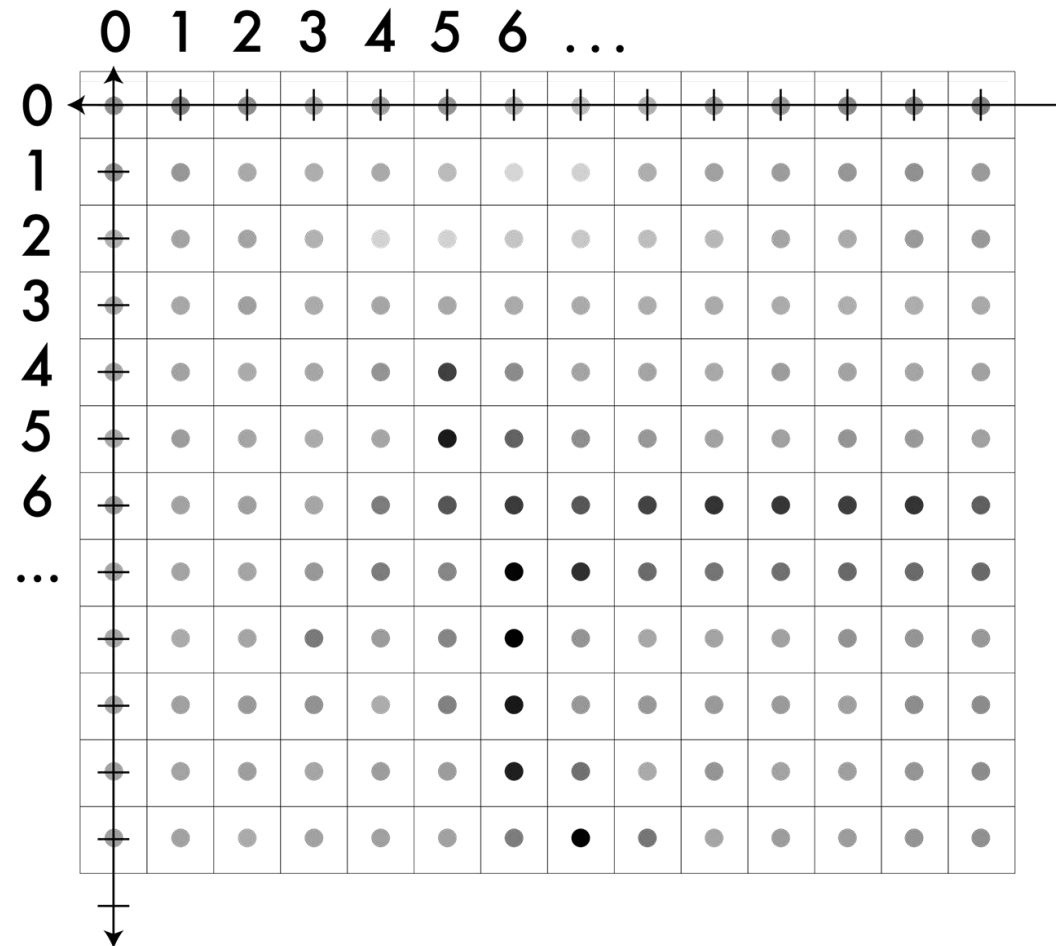


# A note on coordinates in images



Now we can move to a real coordinate system.

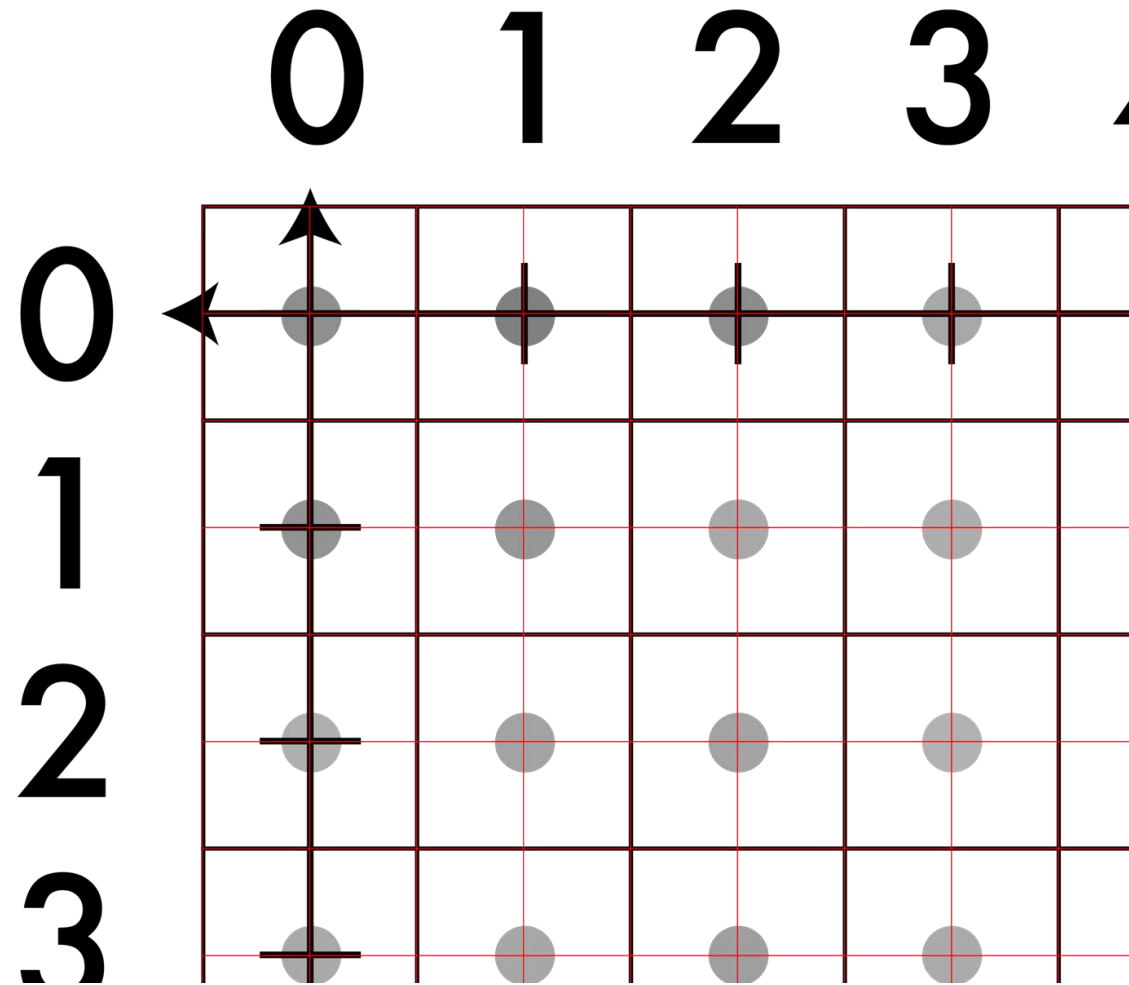
# A note on coordinates in images



On the image

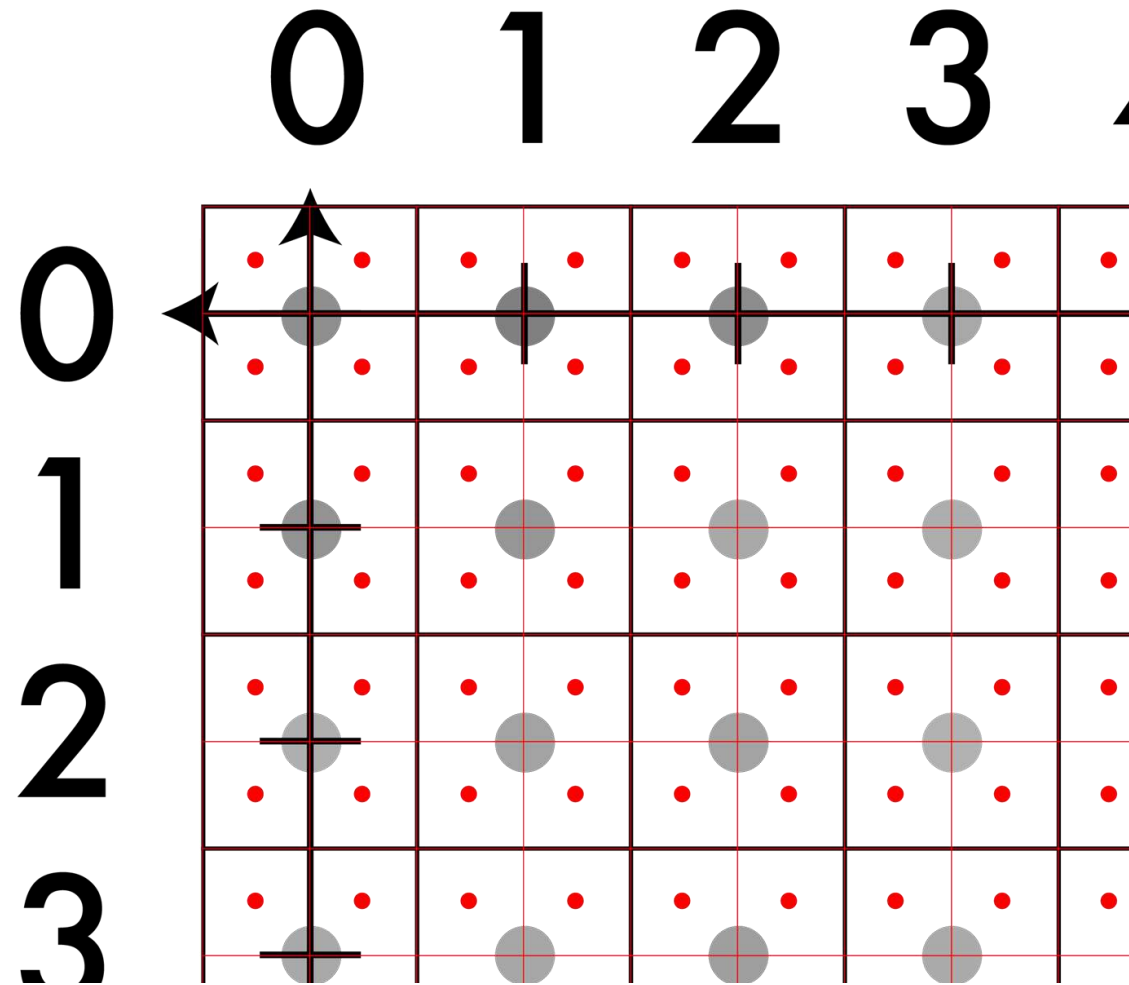
# A note on coordinates in images

So, the value of the pixel  $(x,y)$  is now centered at  $(x,y)$ .



# A note on coordinates in images

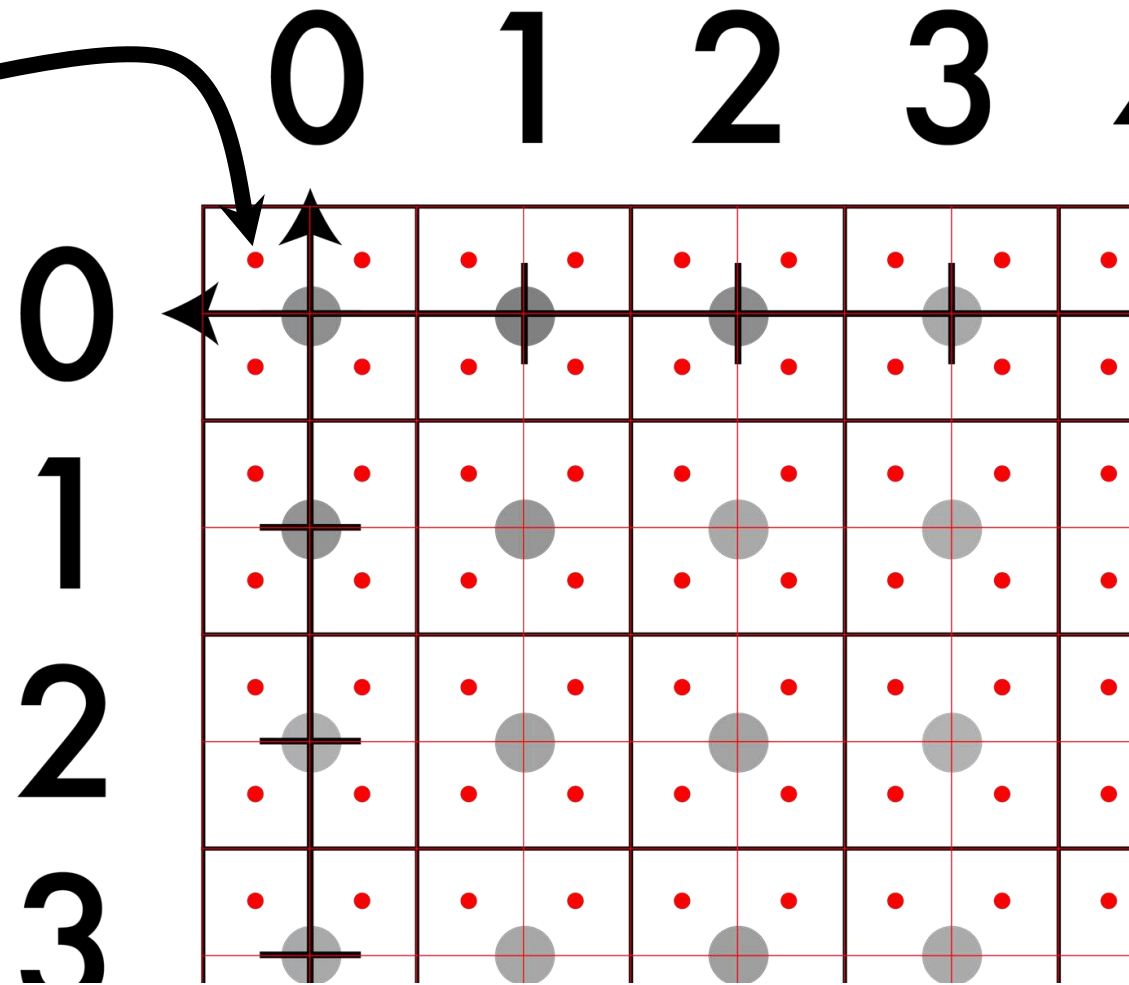
But there are other real-valued points.





# A note on coordinates in images

This point is:  
(-.25, -.25)



# Today's Agenda

- Image basics
  - What is an image – addressing pixels
  - Image as a function – image coordinates
- Image interpolation
  - Nearest neighbor
  - Bilinear
  - Bicubic
- Image resizing
  - Enlarge
  - Shrink

# Interpolation

How do we find out the VALUE of a non-integer point, when the image only comes with integer points, i.e. (25,45,3).

Two simple ideas:

1. Nearest-Neighbor Interpolation
2. Bilinear Interpolation

# Nearest neighbor: what it sounds like

$$f(x,y,z) = \text{Im}(\text{round}(x), \text{round}(y), z)$$

- Looks blocky
- Note:  $z$  is still int



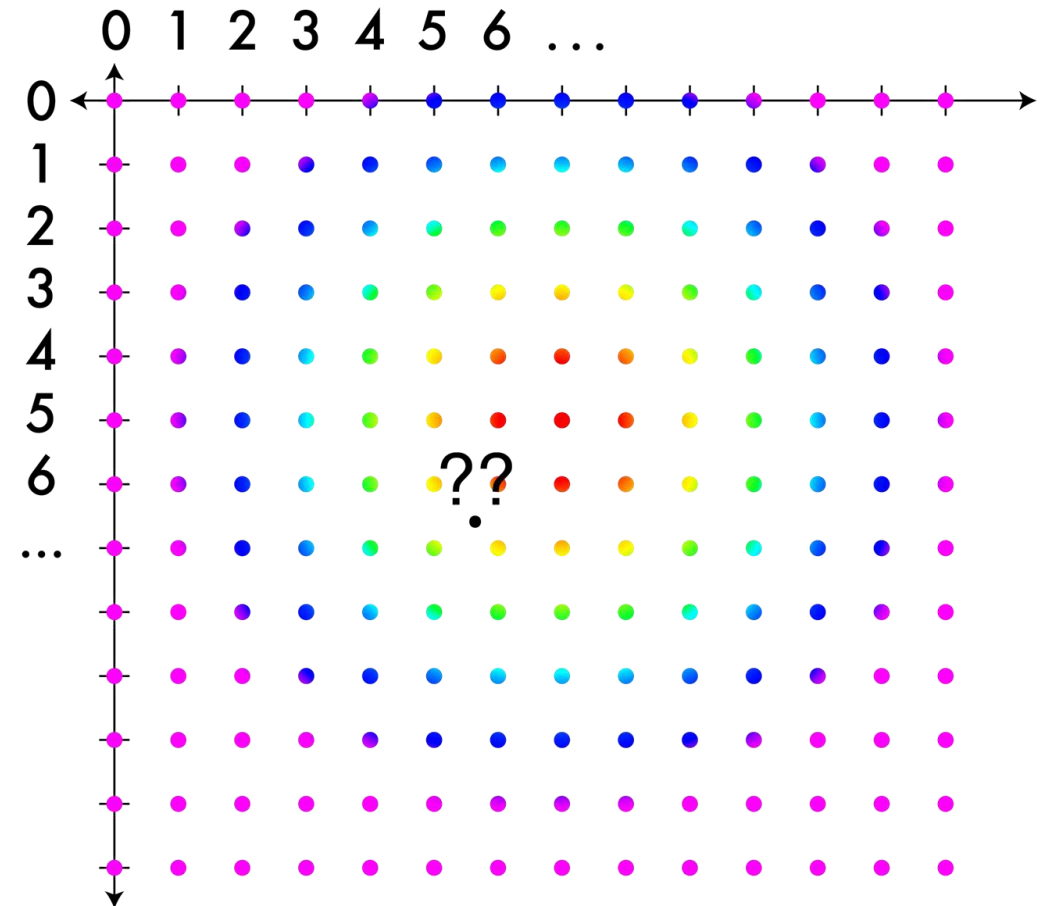
# Today's Agenda

- Image basics
  - What is an image – addressing pixels
  - Image as a function – image coordinates
- Image interpolation
  - Nearest neighbor
  - Bilinear
  - Bicubic
- Image resizing
  - Enlarge
  - Shrink



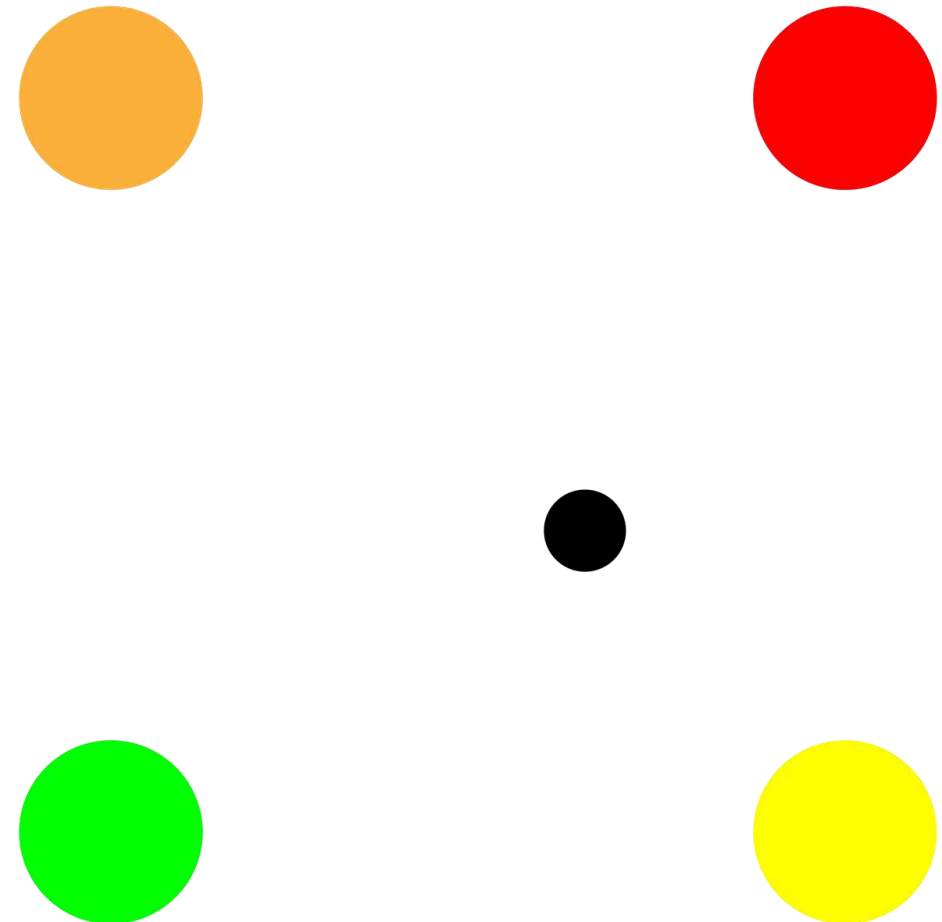
Bilinear interpolation: for grids, pretty good

This time find the closest pixels in a box



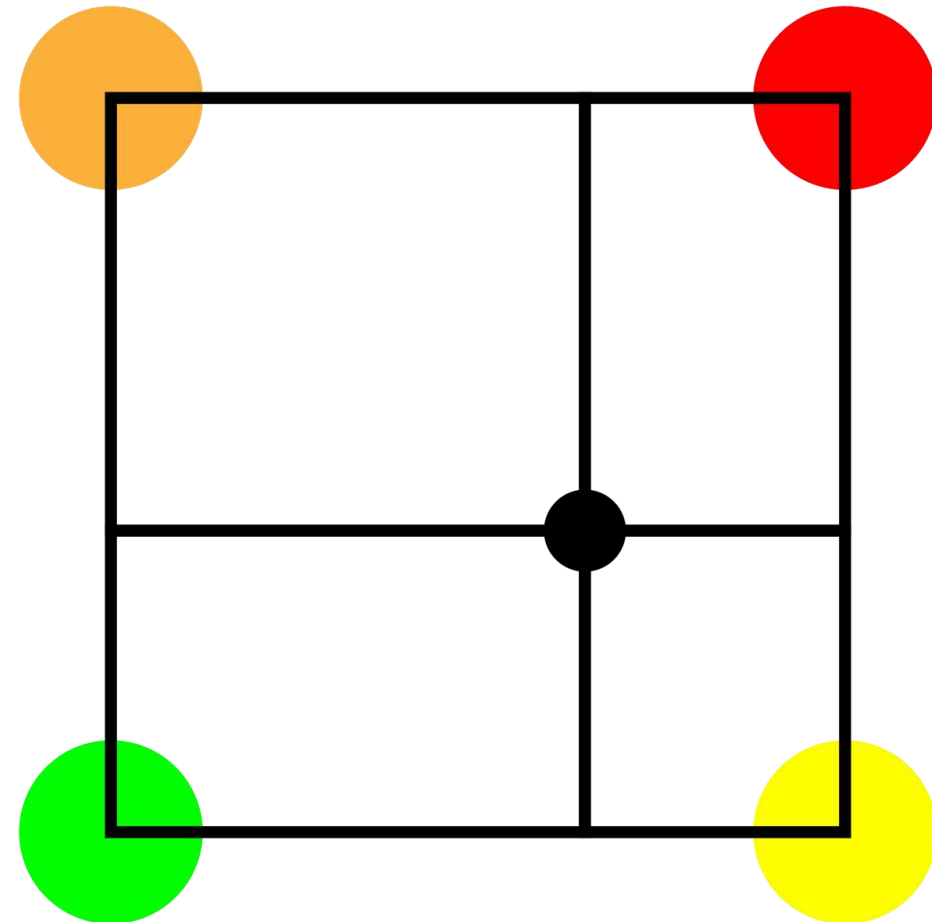
Bilinear interpolation: for grids, pretty good

This time find the closest pixels in  
a box



Bilinear interpolation: for grids, pretty good

This time find the closest pixels in  
a box



Bilinear interpolation: for grids, pretty good

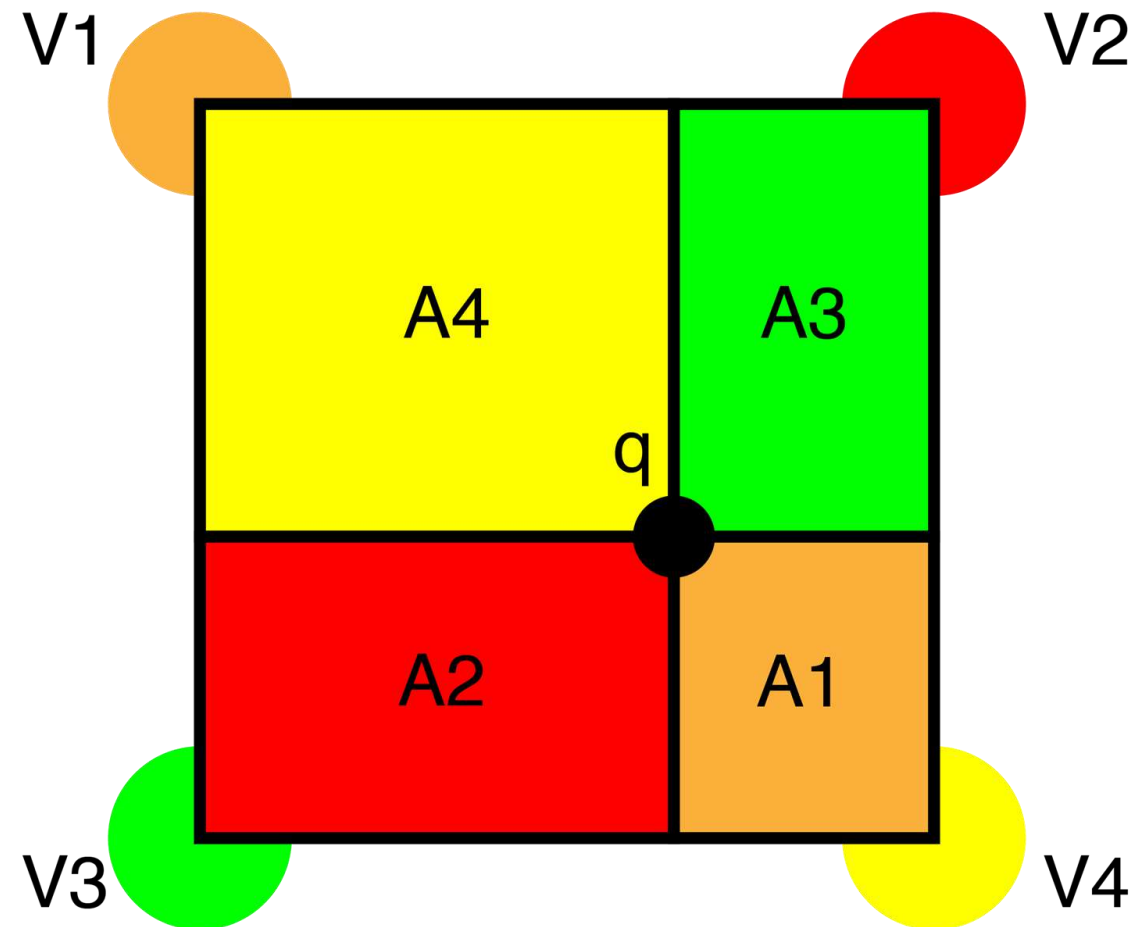
This time find the closest pixels in  
a box

Weighted sum based on area of  
opposite rectangle

$$q = V1 * A1 + V2 * A2 + V3 * A3 + V4 * A4$$

Need to normalize!

Or do we?



## Bilinear interpolation: for grids, pretty good

$$q = V1 * A1 + V2 * A2 + V3 * A3 + V4 * A4$$

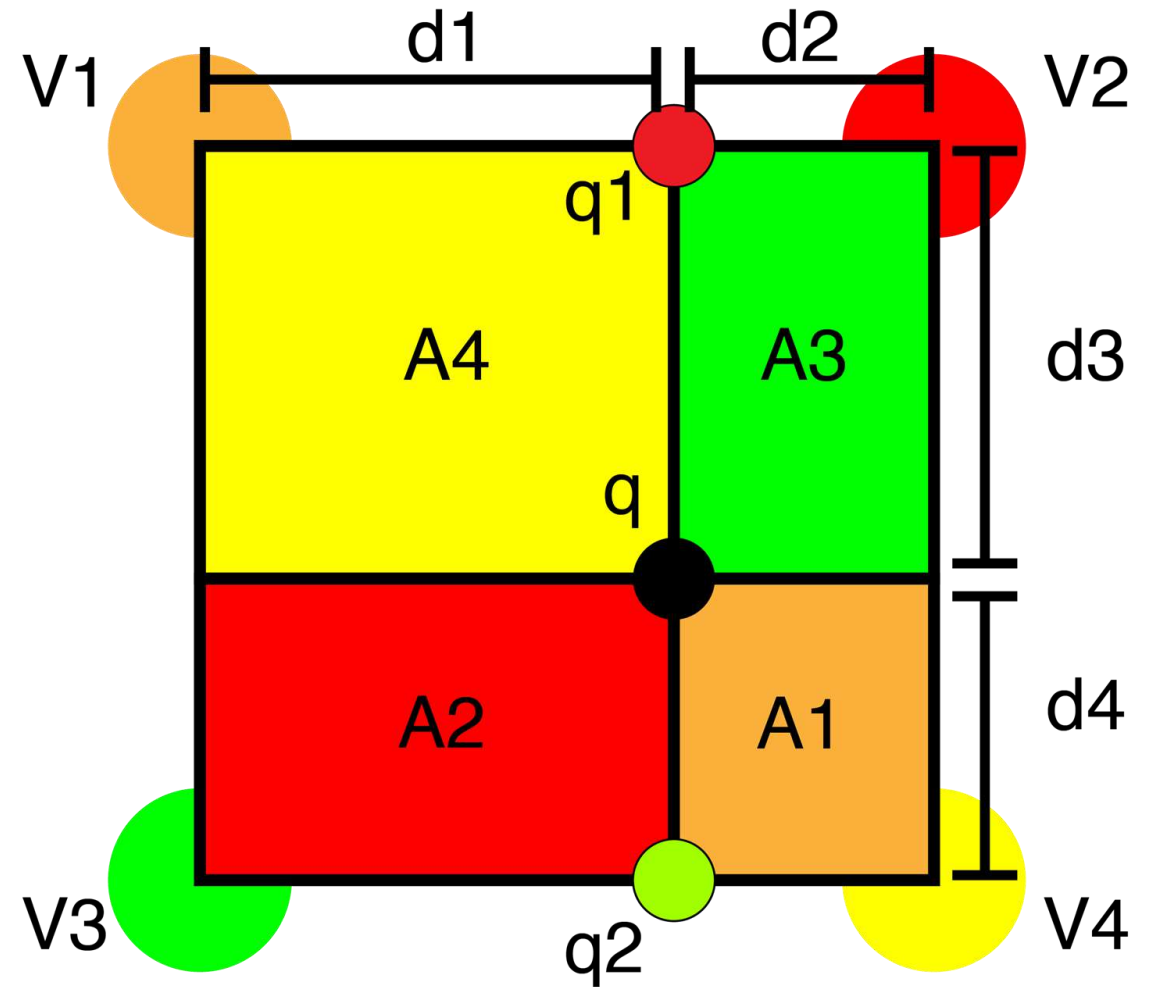
$$A1 = d2 * d4$$

$$A2 = d1 * d4$$

$$A3 = d2 * d3$$

$$A4 = d1 * d3$$

$$\Rightarrow q = V1 * d2 * d4 + V2 * d1 * d4 + V3 * d2 * d3 + V4 * d1 * d3$$





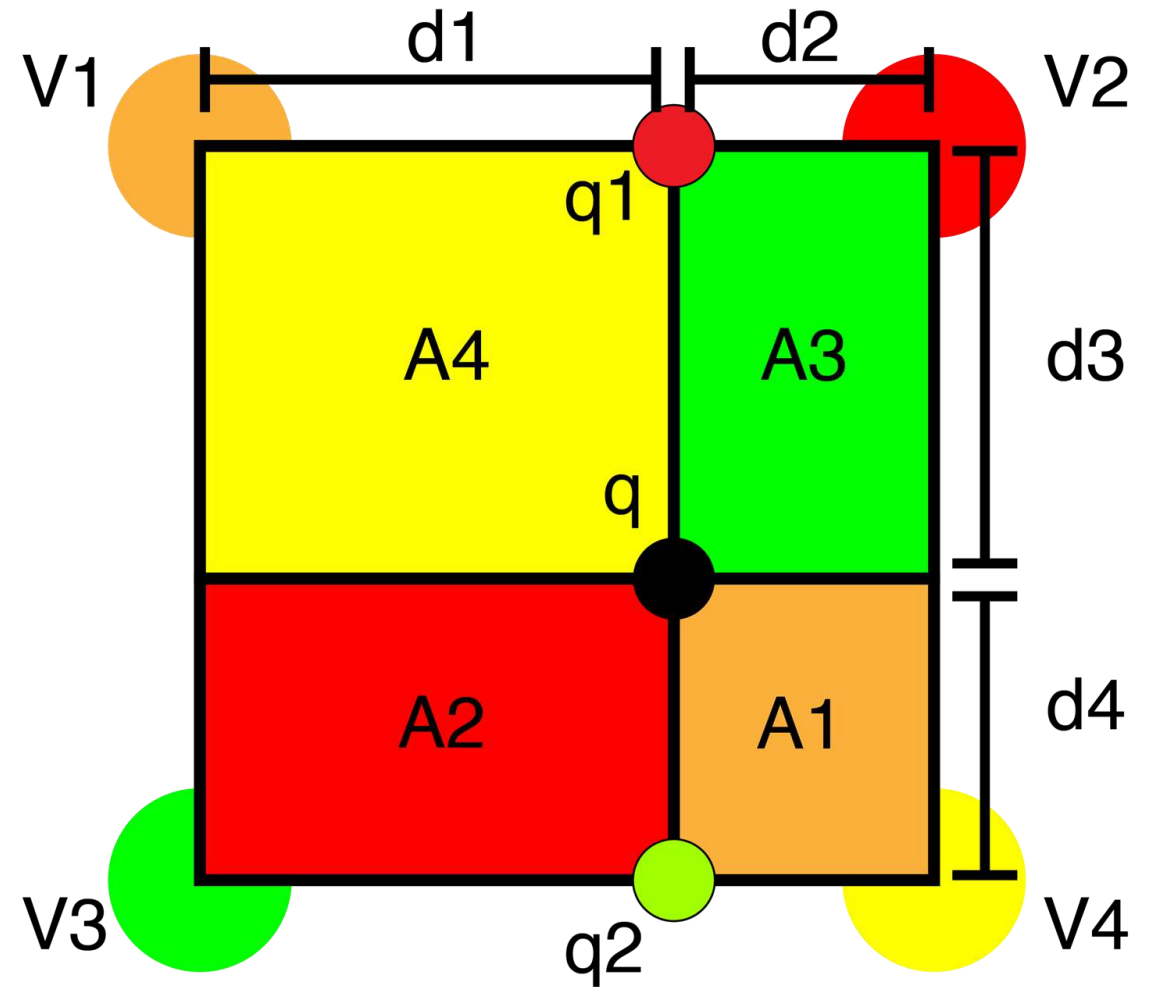
Bilinear interpolation: for grids, pretty good

Alternatively, linear interpolation of linear interpolates

$$q1 = V1 * d2 + V2 * d1$$

$$q2 = V3 * d2 + V4 * d1$$

$$q = q1 * d4 + q2 * d3$$



Bilinear interpolation: for grids, pretty good

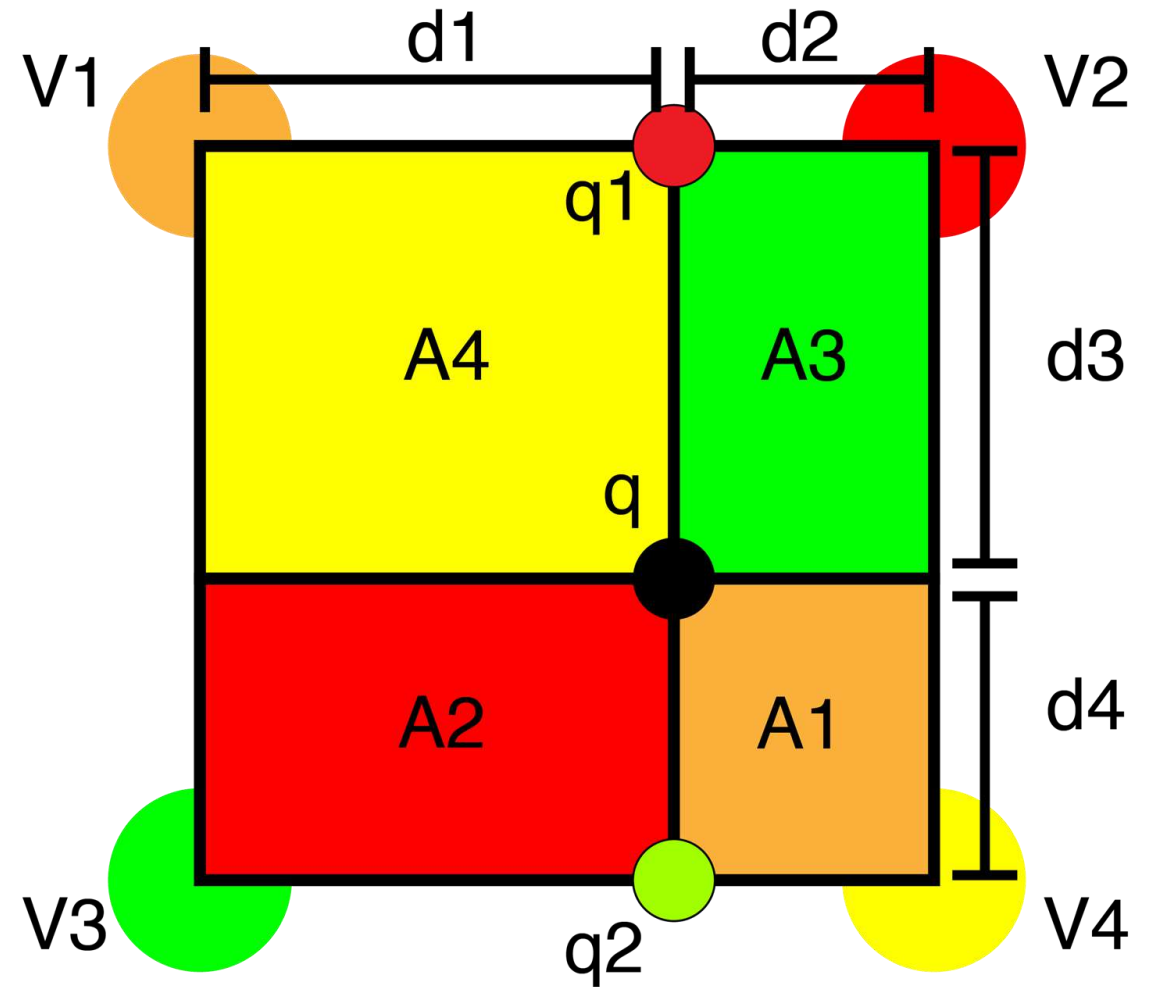
$$q_1 = V_1 * d_2 + V_2 * d_1$$

$$q_2 = V_3 * d_2 + V_4 * d_1$$

$$q = q_1 * d_4 + q_2 * d_3$$

Equivalent:

$$q = q_1 * d_4 + q_2 * d_3$$



Bilinear interpolation: for grids, pretty good

$$q_1 = V_1 * d_2 + V_2 * d_1$$

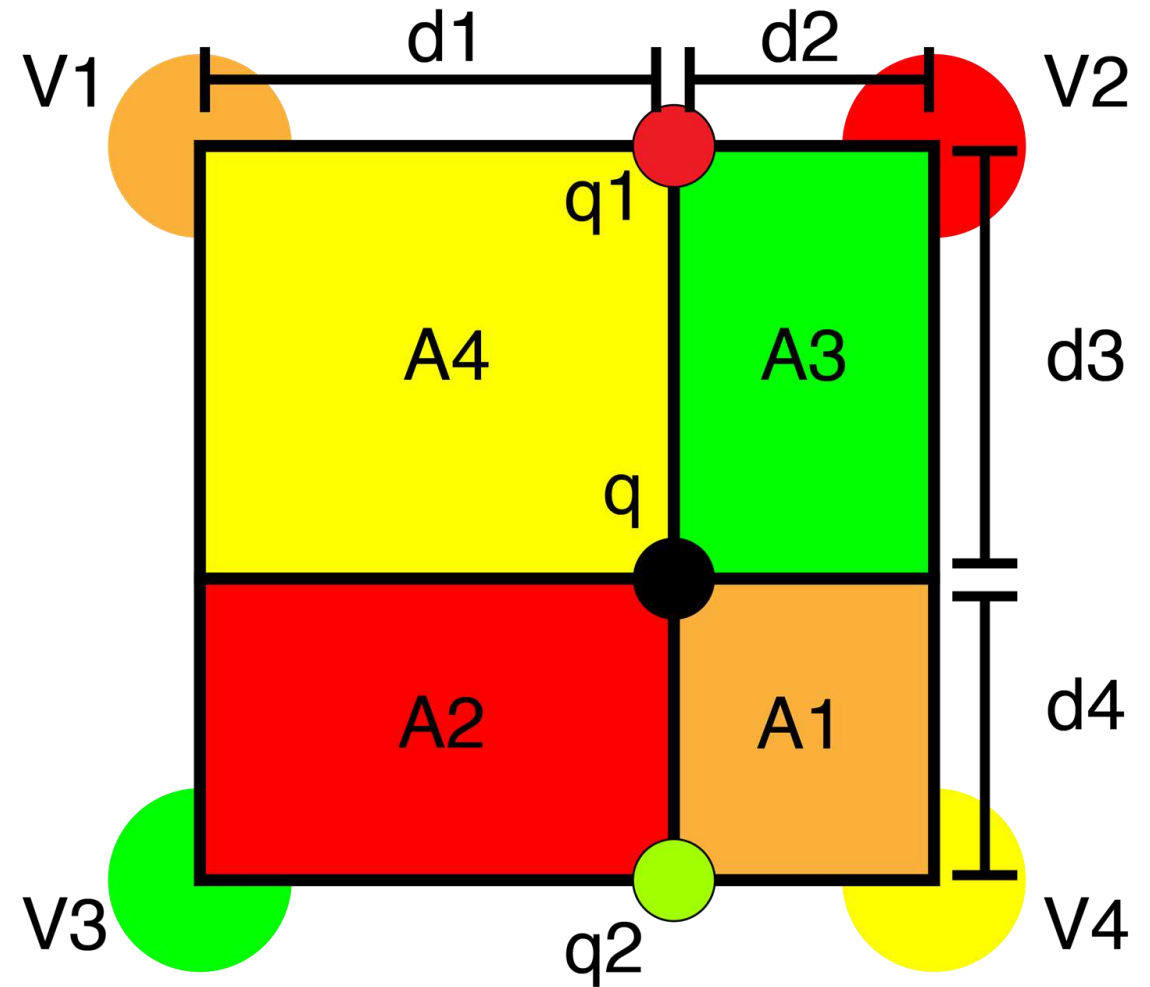
$$q_2 = V_3 * d_2 + V_4 * d_1$$

$$q = q_1 * d_4 + q_2 * d_3$$

Equivalent:

$$q = q_1 * d_4 + q_2 * d_3$$

$$q = (V_1 * d_2 + V_2 * d_1) * d_4 + (V_3 * d_2 + V_4 * d_1) * d_3 \text{ (subst)}$$



## Bilinear interpolation: for grids, pretty good

$$q_1 = V_1 * d_2 + V_2 * d_1$$

$$q_2 = V_3 * d_2 + V_4 * d_1$$

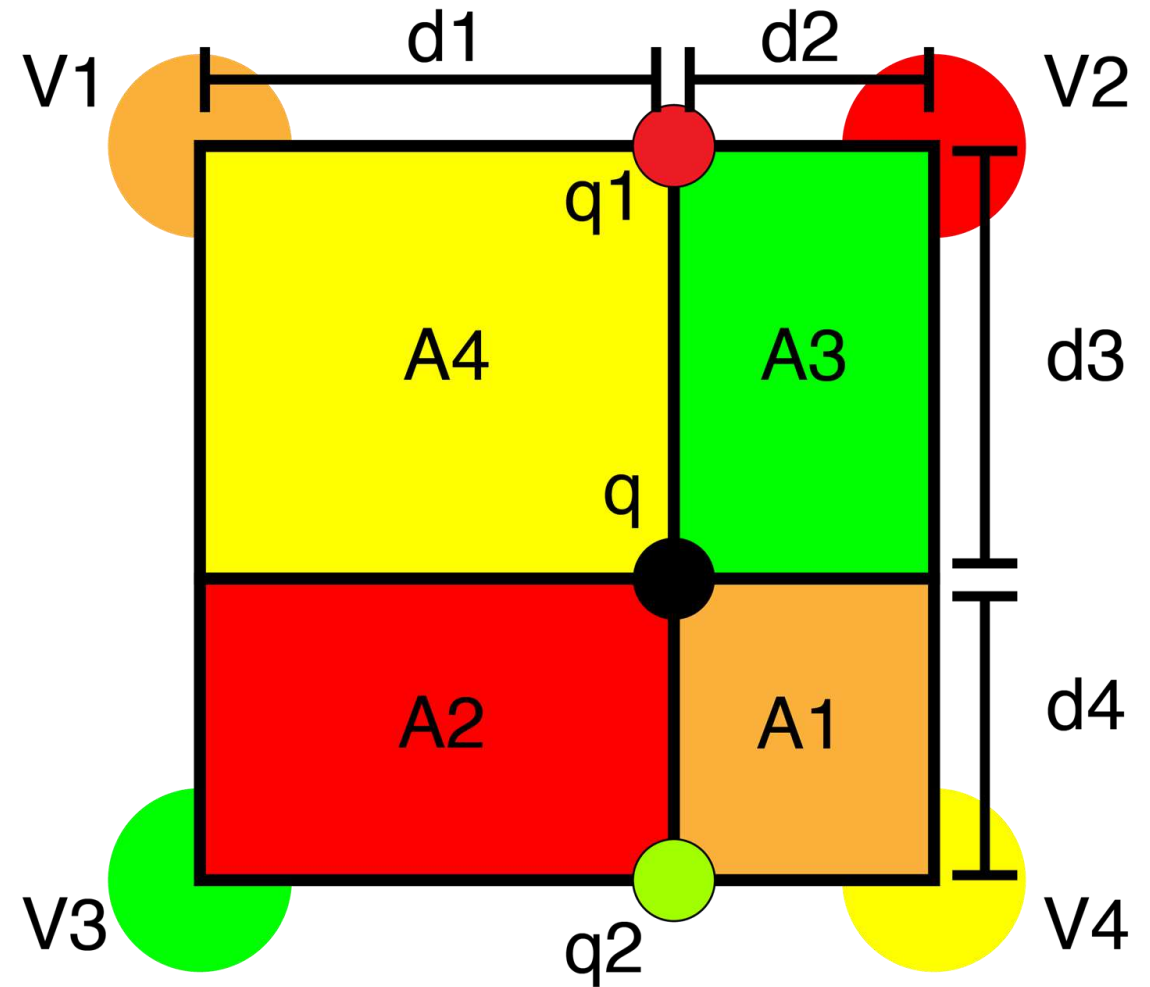
$$q = q_1 * d_4 + q_2 * d_3$$

Equivalent:

$$q = q_1 * d_4 + q_2 * d_3$$

$$q = (V_1 * d_2 + V_2 * d_1) * d_4 + (V_3 * d_2 + V_4 * d_1) * d_3 \text{ (subst)}$$

$$q = V_1 * d_2 * d_4 + V_2 * d_1 * d_4 + V_3 * d_2 * d_3 + V_4 * d_1 * d_3 \text{ (distribution)}$$



## Bilinear interpolation: for grids, pretty good

$$q_1 = V_1 * d_2 + V_2 * d_1$$

$$q_2 = V_3 * d_2 + V_4 * d_1$$

$$q = q_1 * d_4 + q_2 * d_3$$

Equivalent:

$$q = q_1 * d_4 + q_2 * d_3$$

$$q = (V_1 * d_2 + V_2 * d_1) * d_4 + (V_3 * d_2 + V_4 * d_1) * d_3 \text{ (subst)}$$

$$q = V_1 * d_2 * d_4 + V_2 * d_1 * d_4 + V_3 * d_2 * d_3 + V_4 * d_1 * d_3 \text{ (distribution)}$$

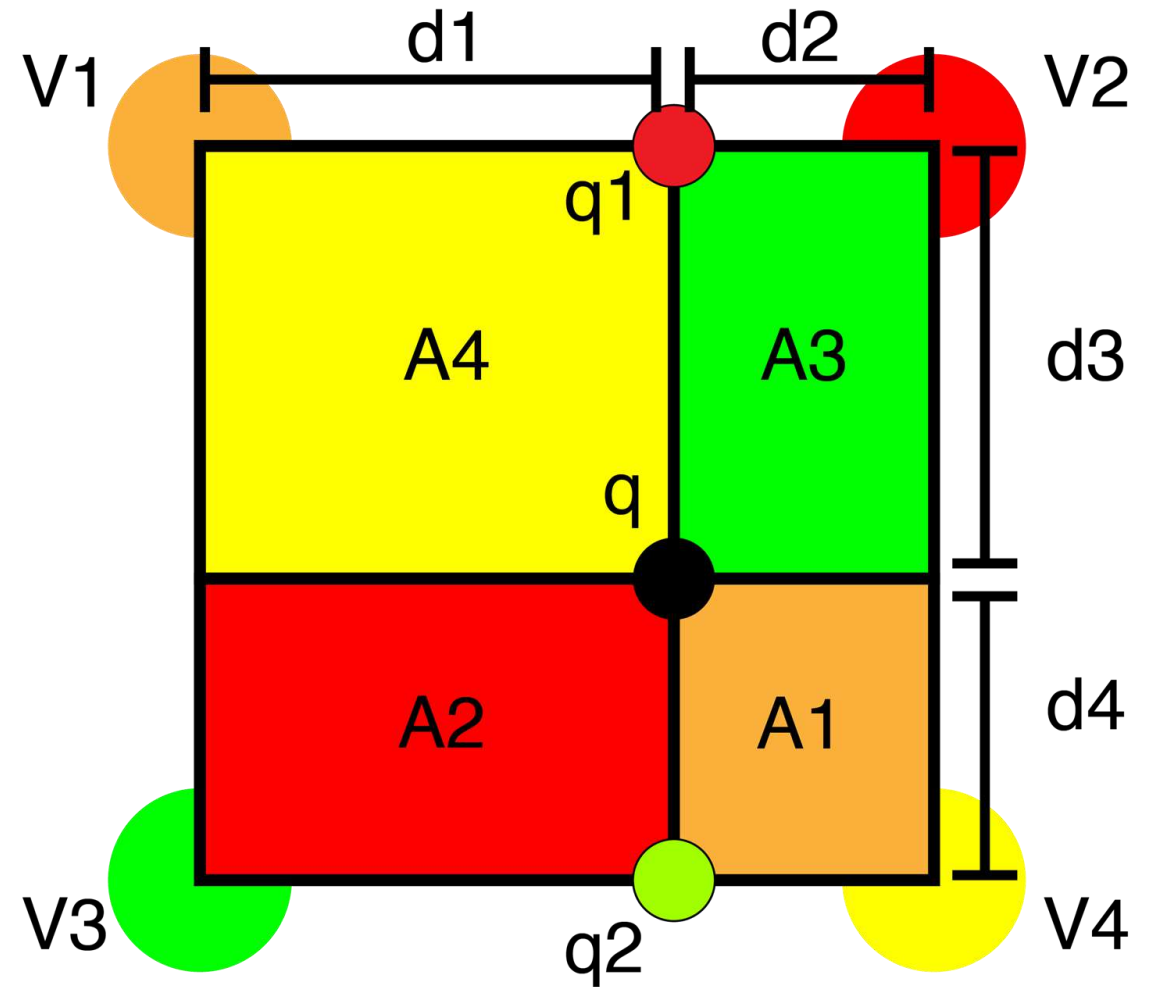
Recall:

$$A_1 = d_2 * d_4$$

$$A_2 = d_1 * d_4$$

$$A_3 = d_2 * d_3$$

$$A_4 = d_1 * d_3$$





Bilinear interpolation: for grids, pretty good

$$q_1 = V_1 * d_2 + V_2 * d_1$$

$$q_2 = V_3 * d_2 + V_4 * d_1$$

$$q = q_1 * d_4 + q_2 * d_3$$

Equivalent:

$$q = q_1 * d_4 + q_2 * d_3$$

$$q = (V_1 * d_2 + V_2 * d_1) * d_4 + (V_3 * d_2 + V_4 * d_1) * d_3 \text{ (subst)}$$

$$q = V_1 * d_2 * d_4 + V_2 * d_1 * d_4 + V_3 * d_2 * d_3 + V_4 * d_1 * d_3 \text{ (distribution)}$$

Recall:

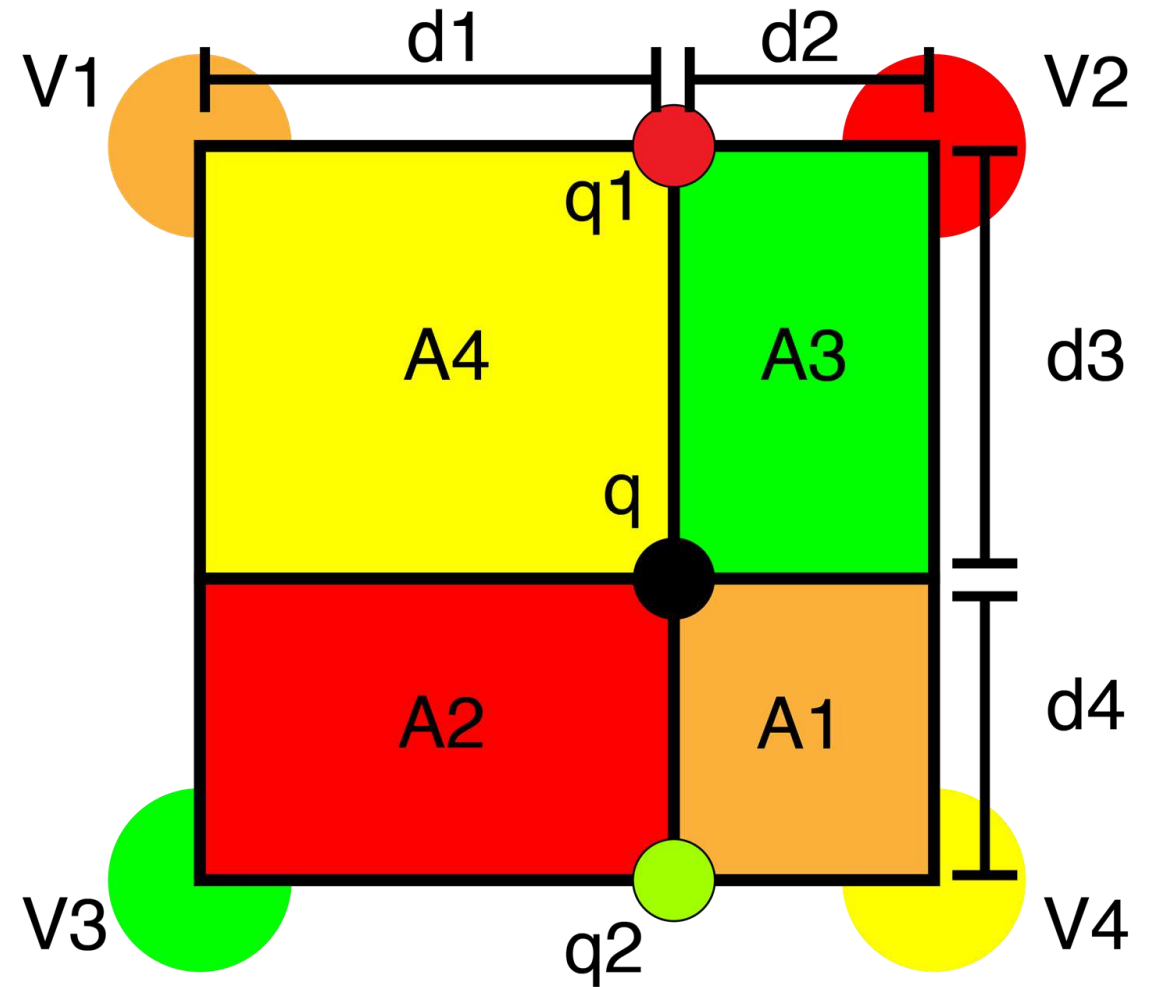
$$A_1 = d_2 * d_4$$

$$A_2 = d_1 * d_4$$

$$A_3 = d_2 * d_3$$

$$A_4 = d_1 * d_3$$

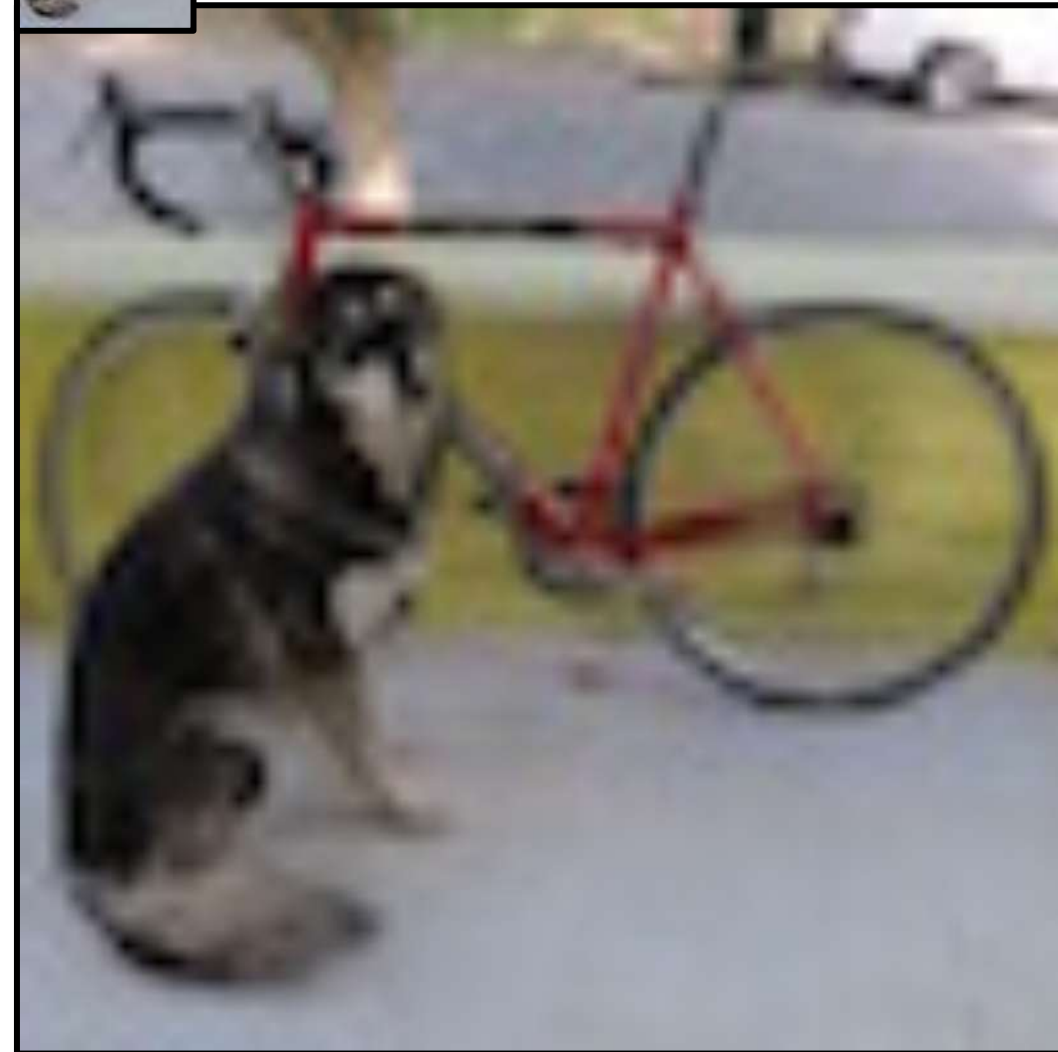
$$q = V_1 * A_1 + V_2 * A_2 + V_3 * A_3 + V_4 * A_4$$





## Bilinear interpolation: for grids, pretty good

- Smoother than NN
- More complex
  - 4 lookups
  - Some math
- Often the right tradeoff of speed vs final result

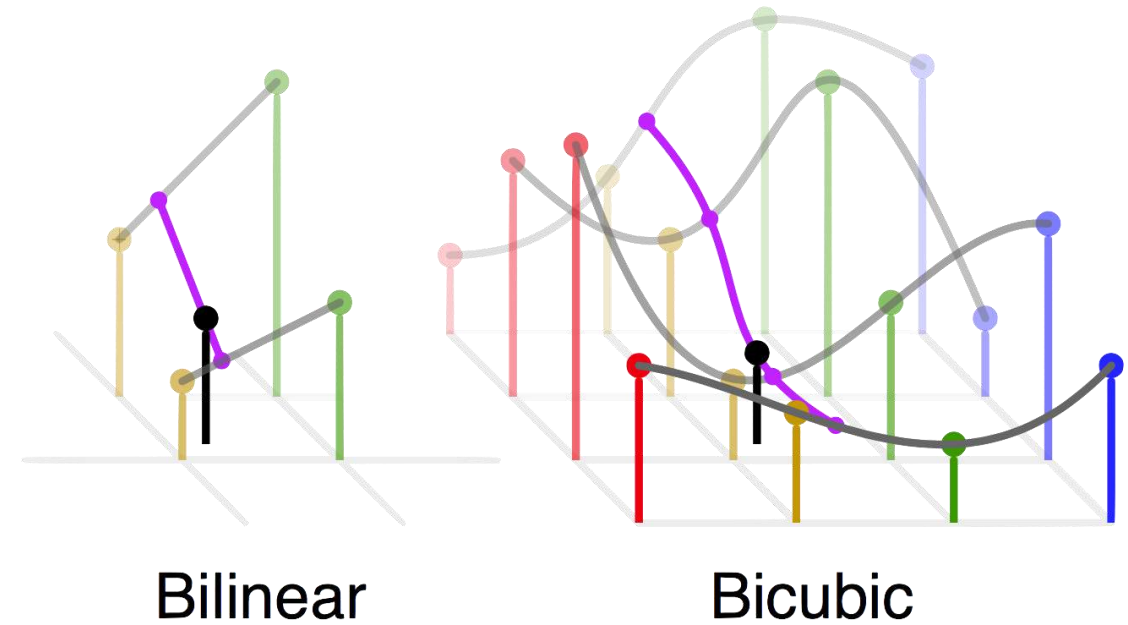


# Today's Agenda

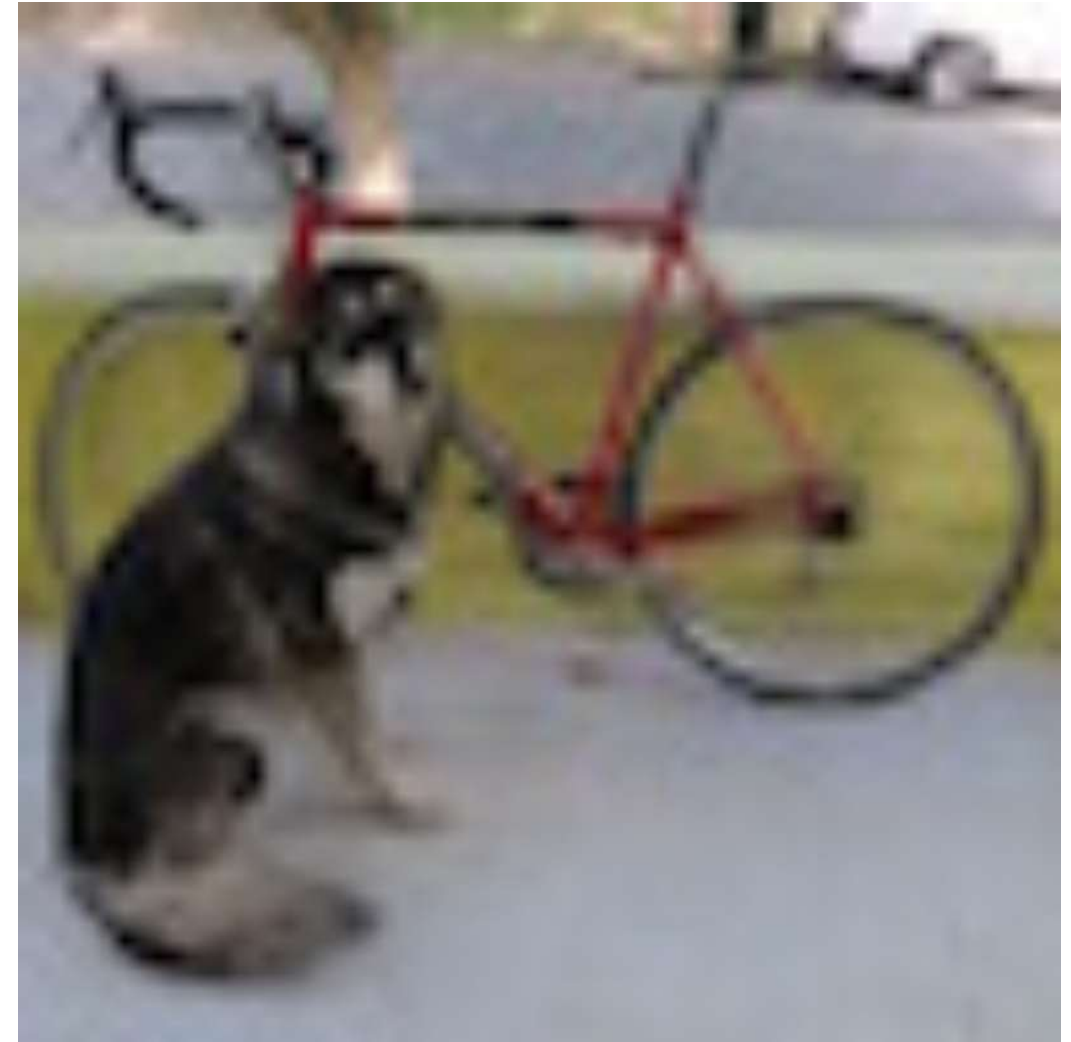
- Image basics
  - What is an image – addressing pixels
  - Image as a function – image coordinates
- Image interpolation
  - Nearest neighbor
  - Bilinear
  - Bicubic
- Image resizing
  - Enlarge
  - Shrink

## Bicubic sampling: more complex, maybe better?

- A cubic interpolation of 4 cubic interpolations
- Smoother than bilinear, no “star”
- 16 nearest neighbors
- Fit 3rd order poly:
  - $f(x) = a + bx + cx^2 + dx^3$
- Interpolate along axis
- Fit another poly to interpolated values

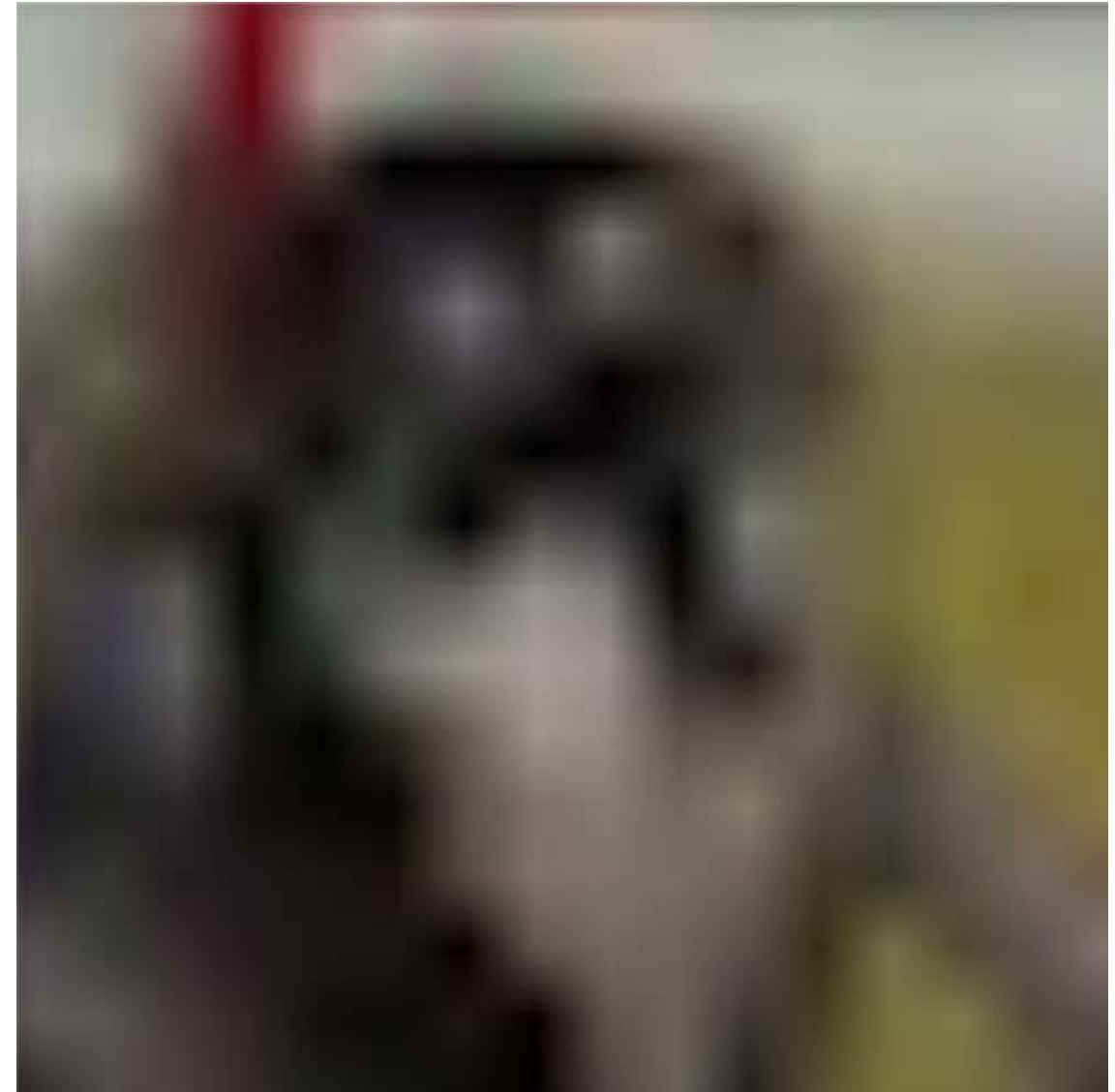
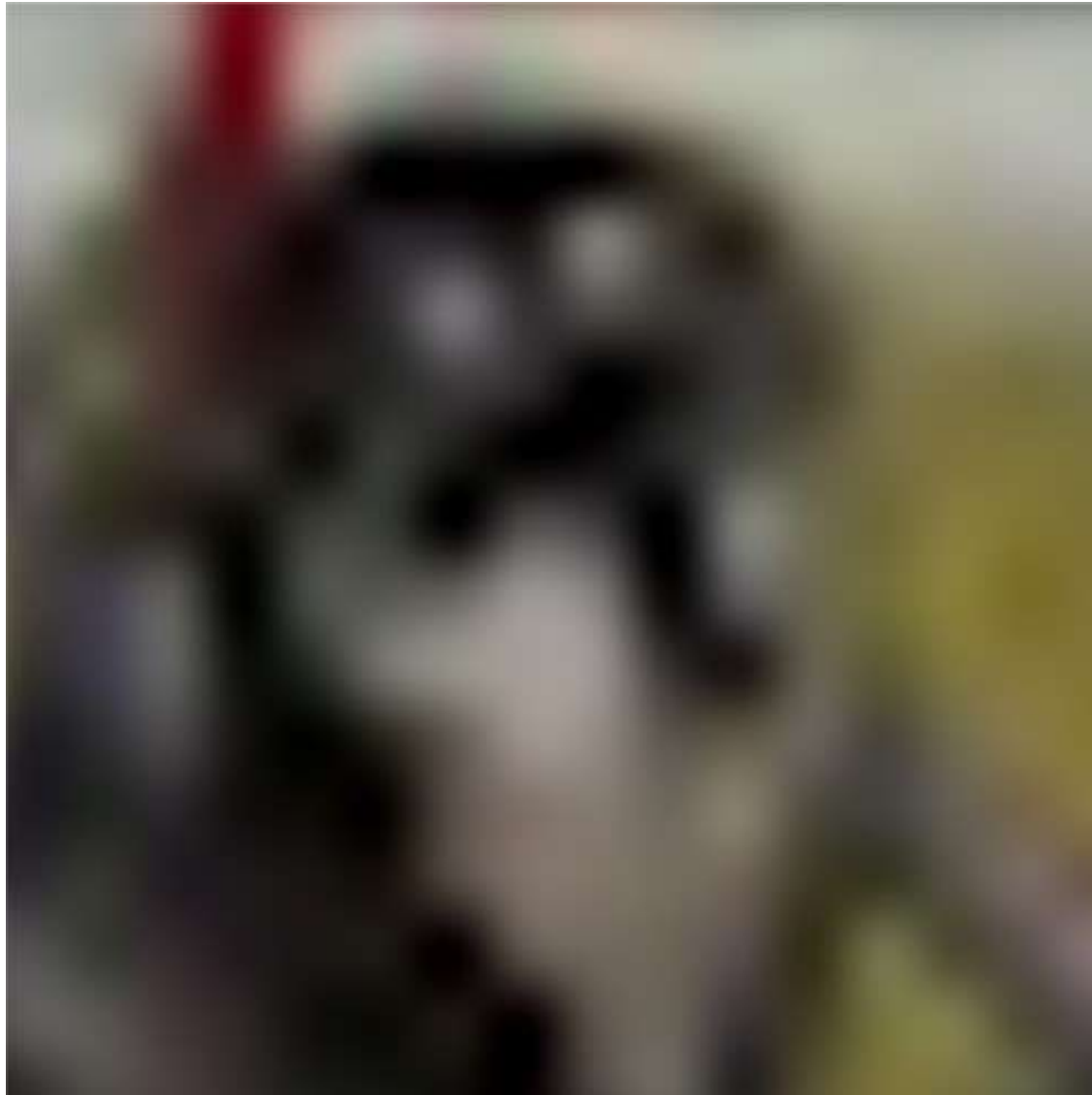


## Bicubic vs bilinear



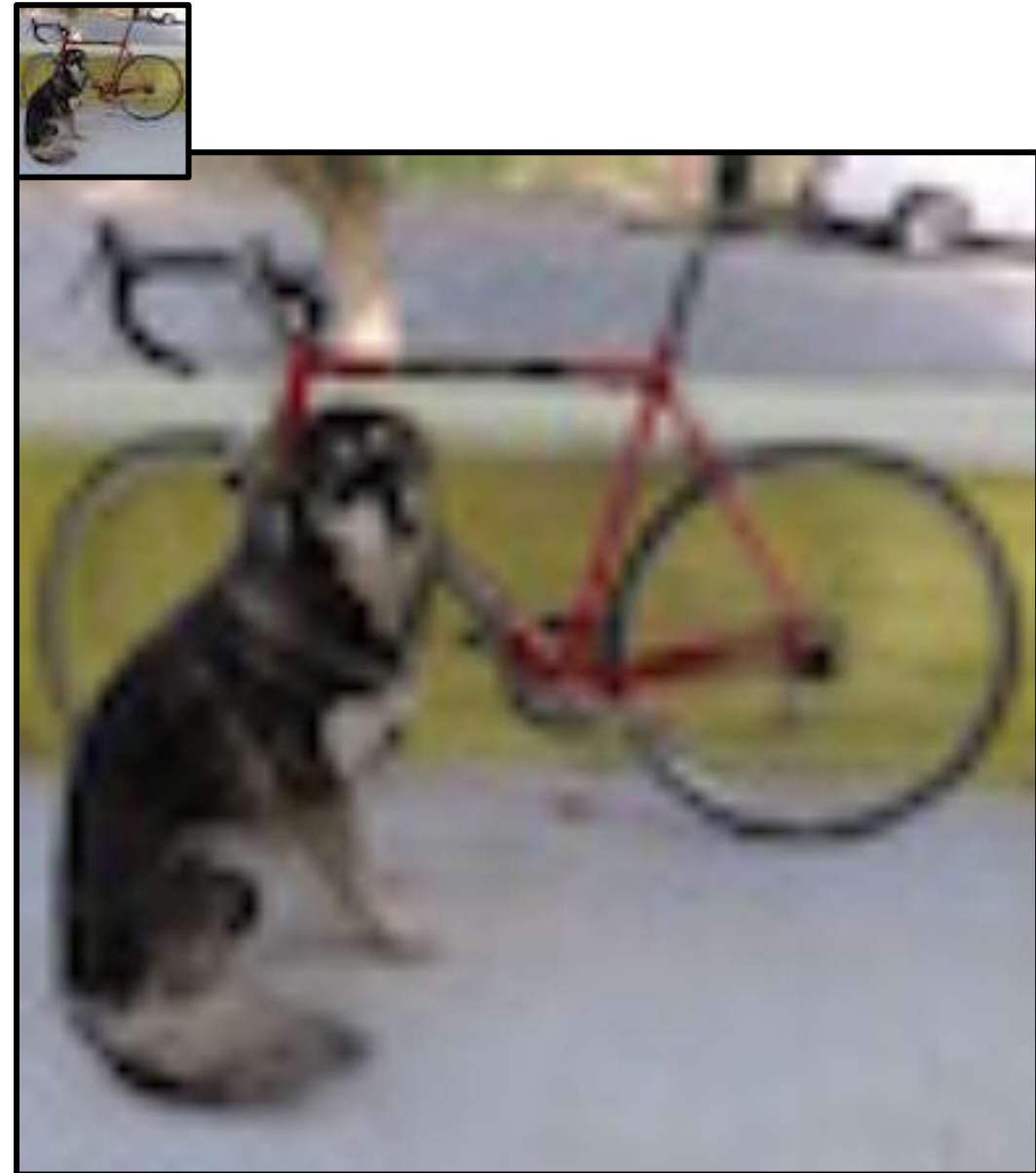


## Bicubic vs bilinear



# Resize algorithm:

- For each pixel in new image:
  - Map to old im coordinates
  - Interpolate value
  - Set new value in image



# So what is this interpolation useful for?

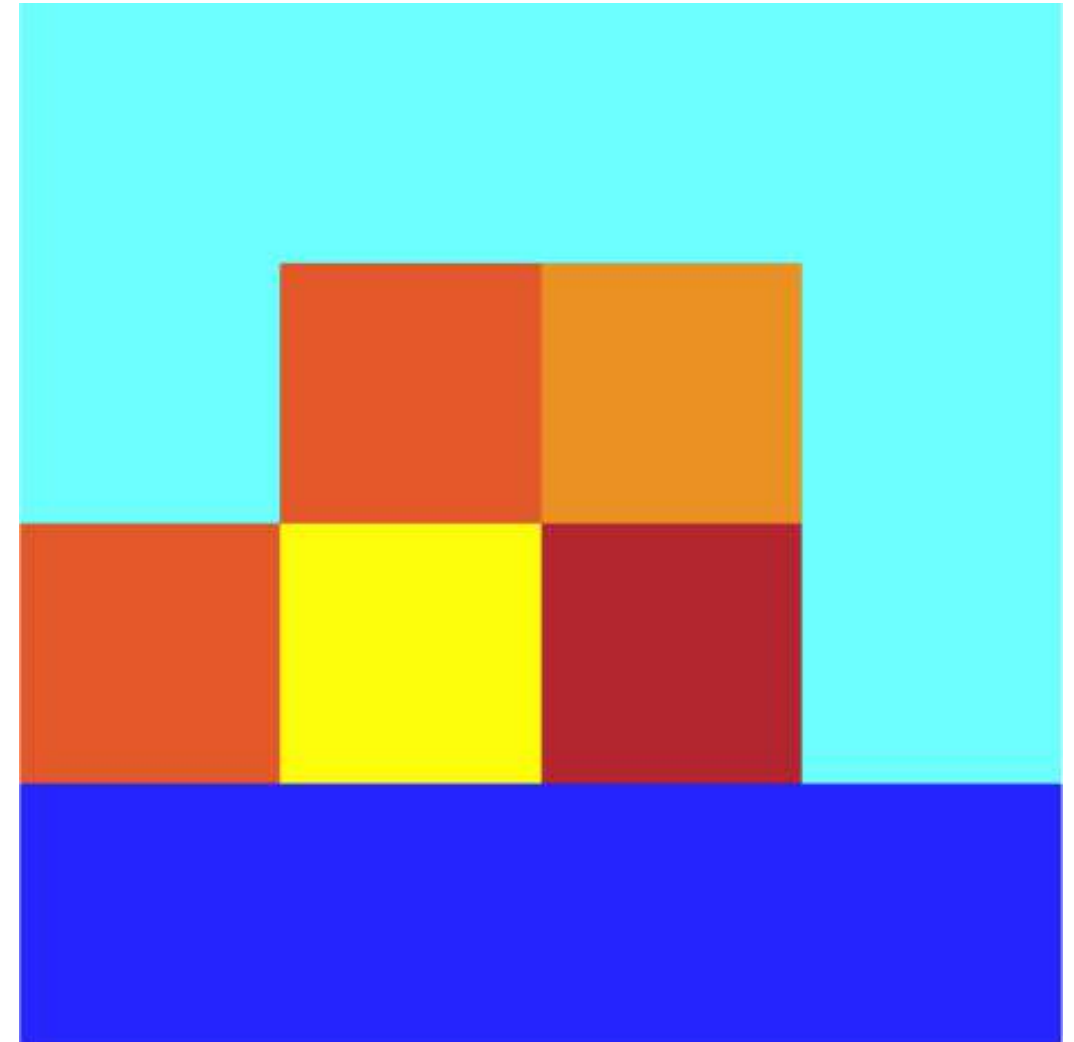
# Today's Agenda

- Image basics
  - What is an image – addressing pixels
  - Image as a function – image coordinates
- Image interpolation
  - Nearest neighbor
  - Bilinear
  - Bicubic
- Image resizing
  - Enlarge
  - Shrink

# Image resizing!

Say we want to increase the size of an image...

This is a beautiful image of a sunset... it's just very small...



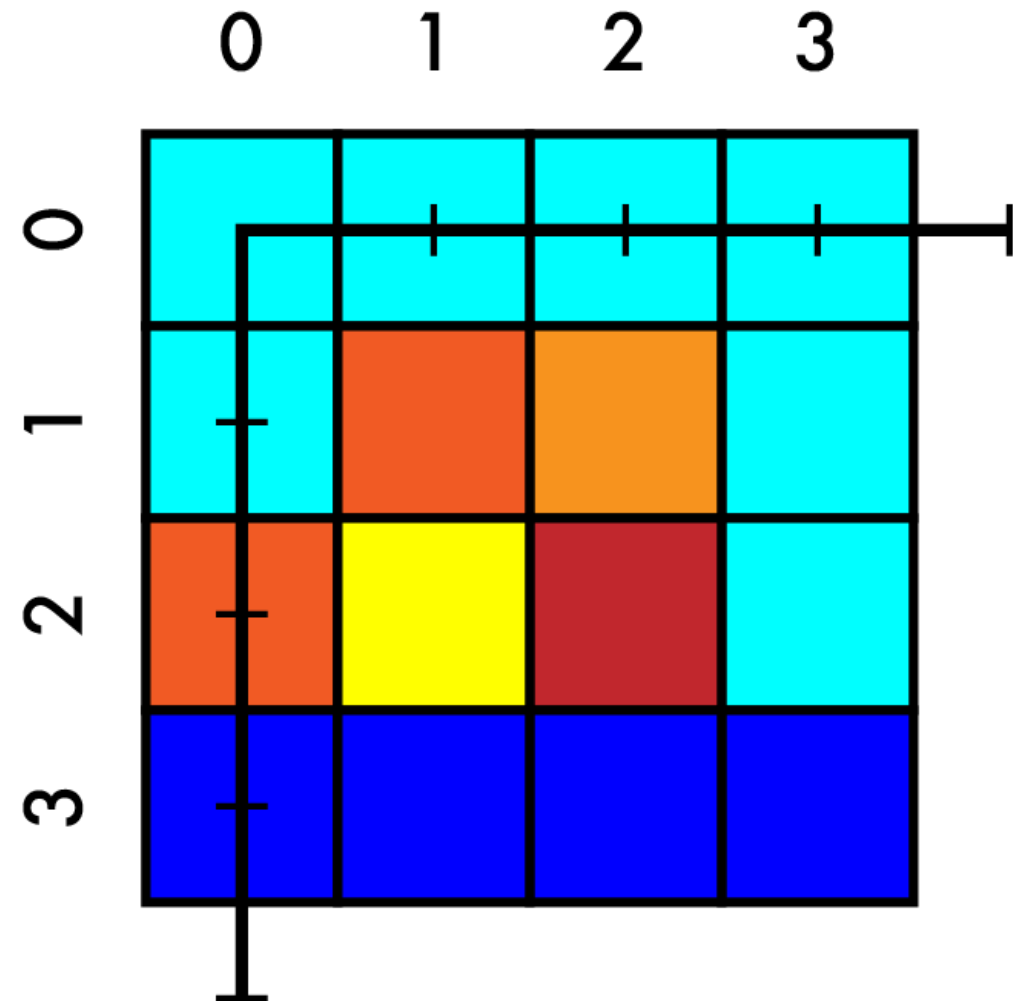


# Image resizing!

Say we want to increase the size of an image...

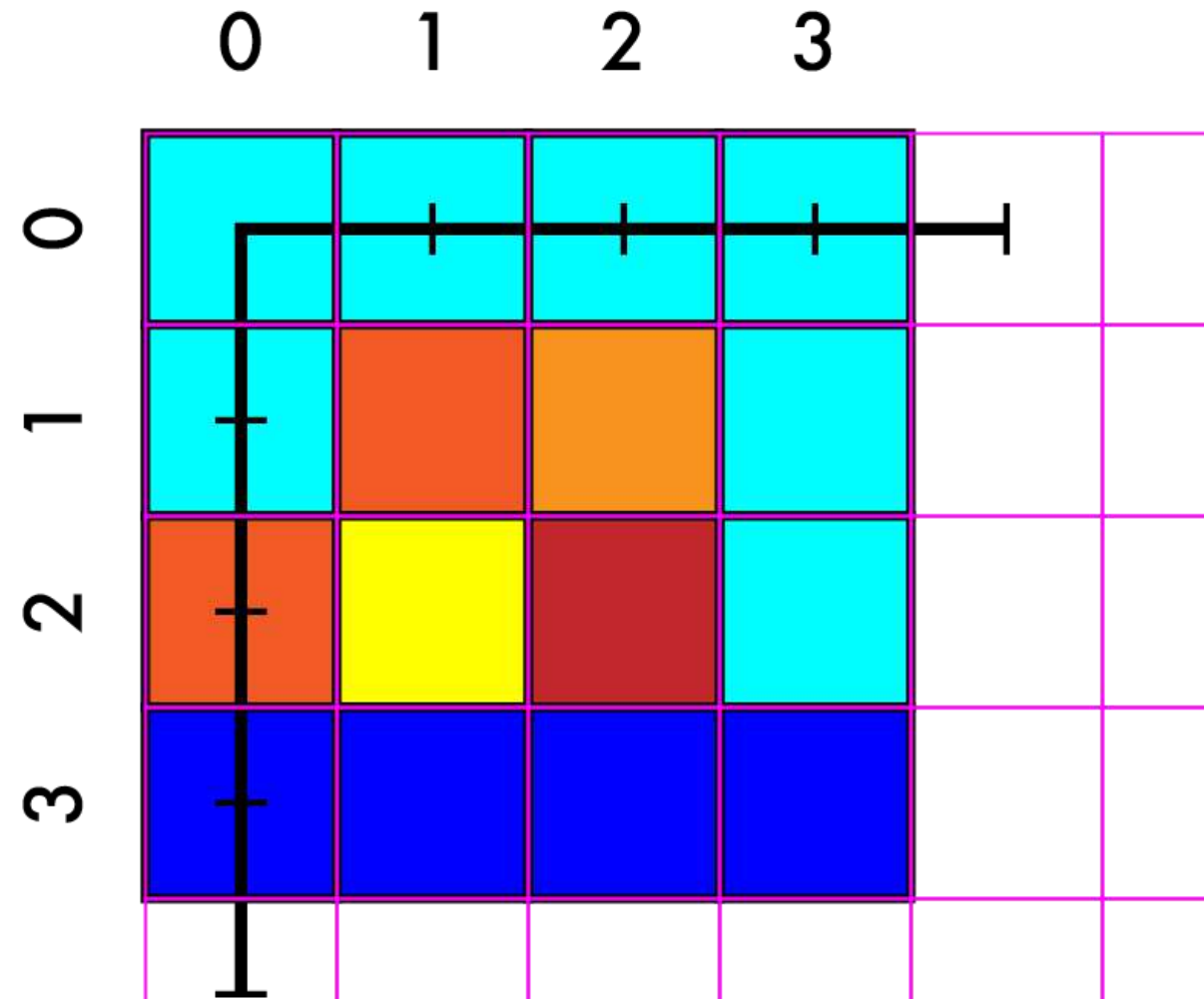
This is a beautiful image of a sunset... it's just very small...

Say we want to increase size 4x4 - > 7x7



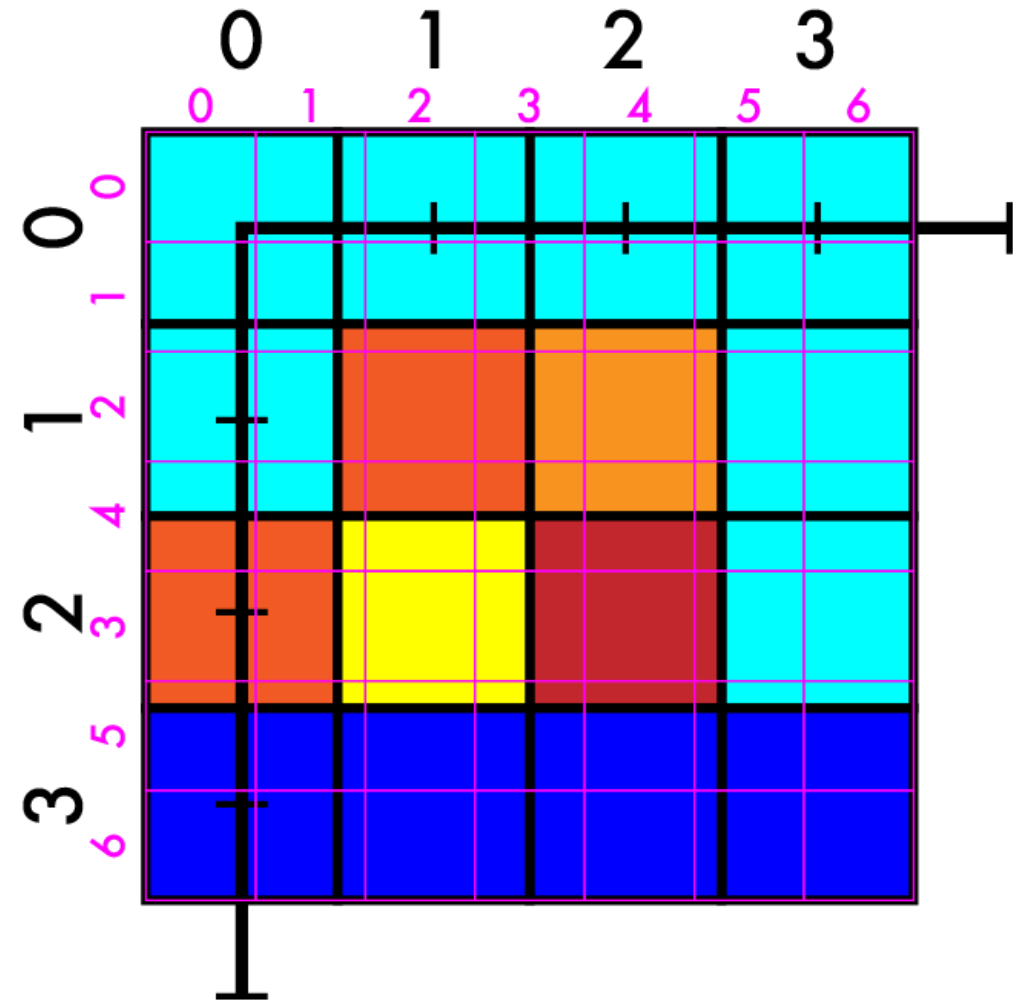
## Resize 4x4 -> 7x7

- Create our new image



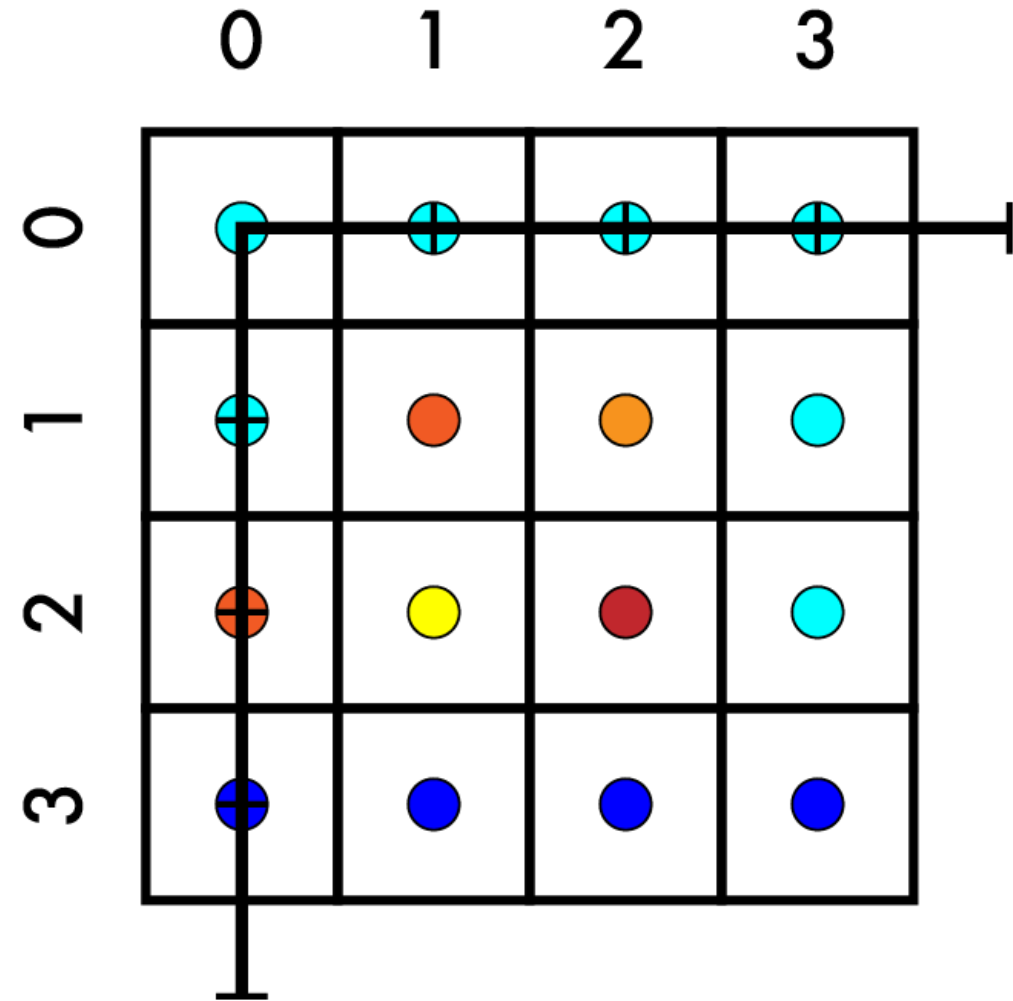
## Resize 4x4 -> 7x7

- Create our new image
- Match up coordinates



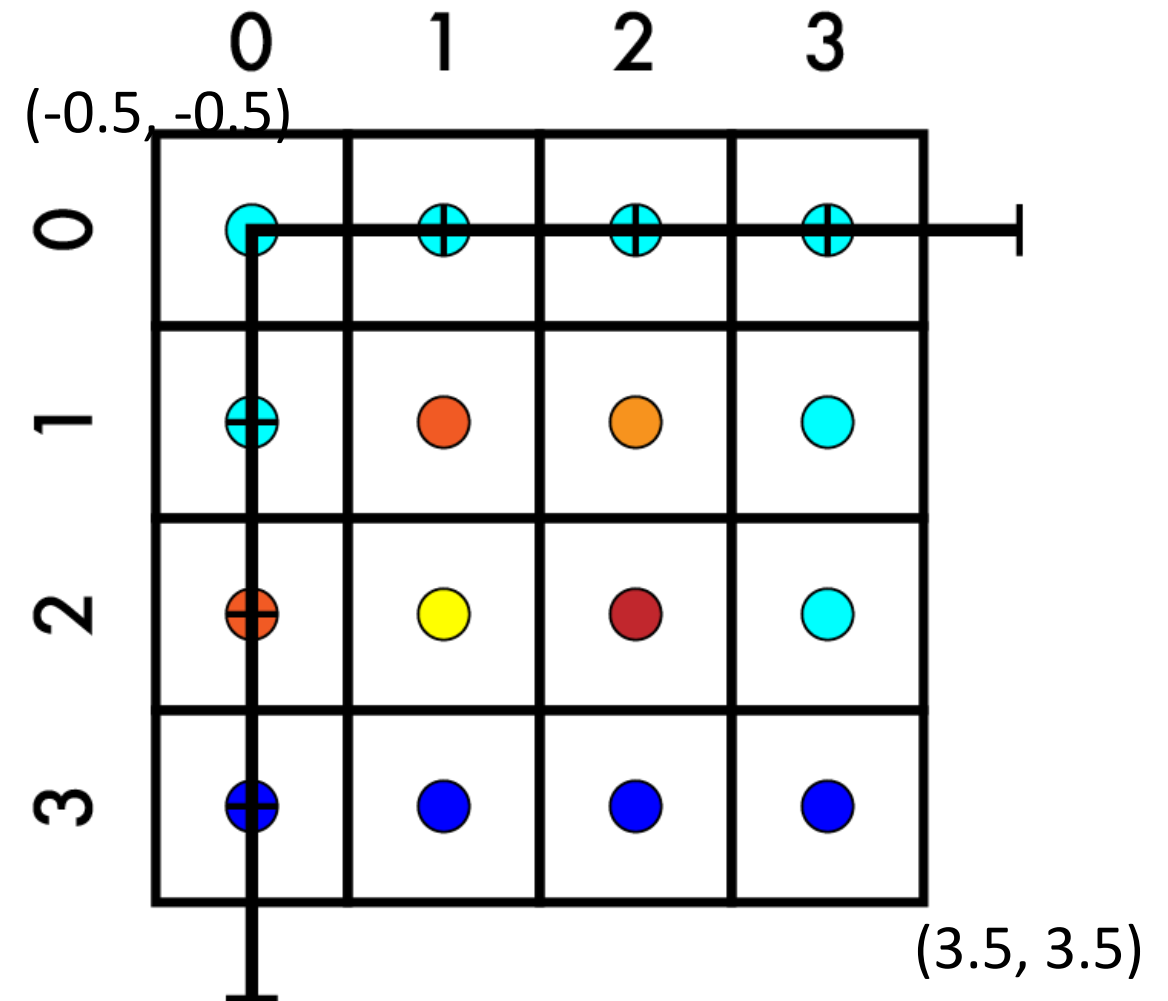
## Resize 4x4 -> 7x7

- Create our new image
- Match up coordinates



## Resize 4x4 -> 7x7

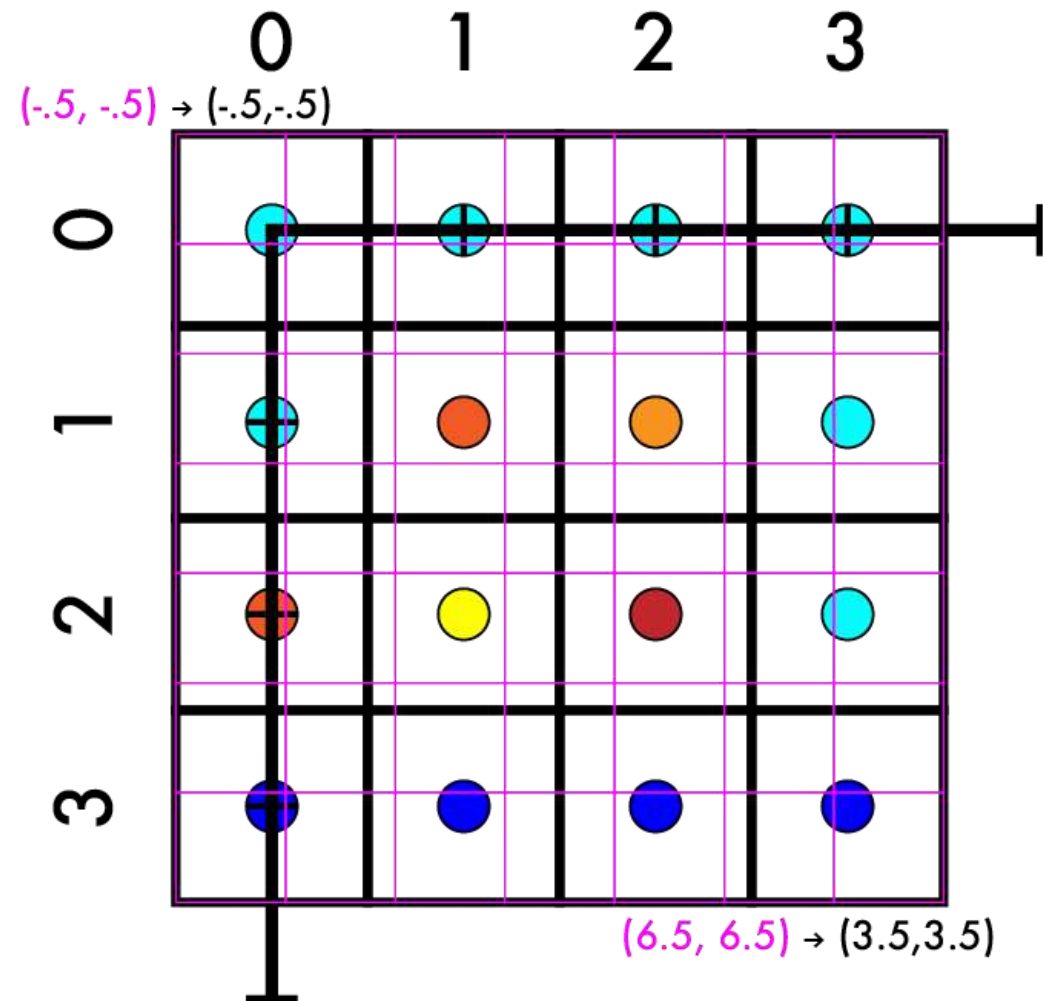
- Create our new image
- Match up coordinates





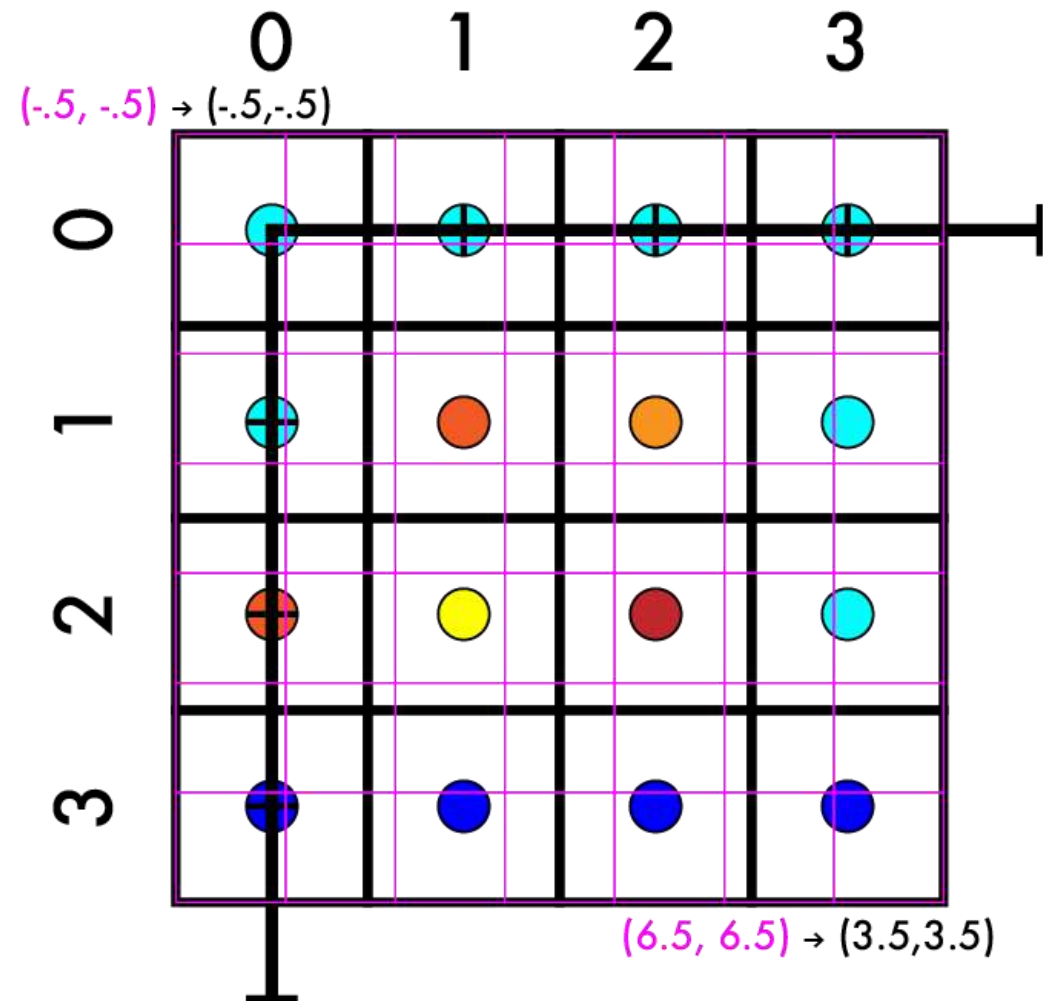
## Resize 4x4 -> 7x7

- Create our new image
- Match up coordinates
  - System of equations
  - $aX + b = Y$
  - $a * -.5 + b = -.5$
  - $a * 6.5 + b = 3.5$



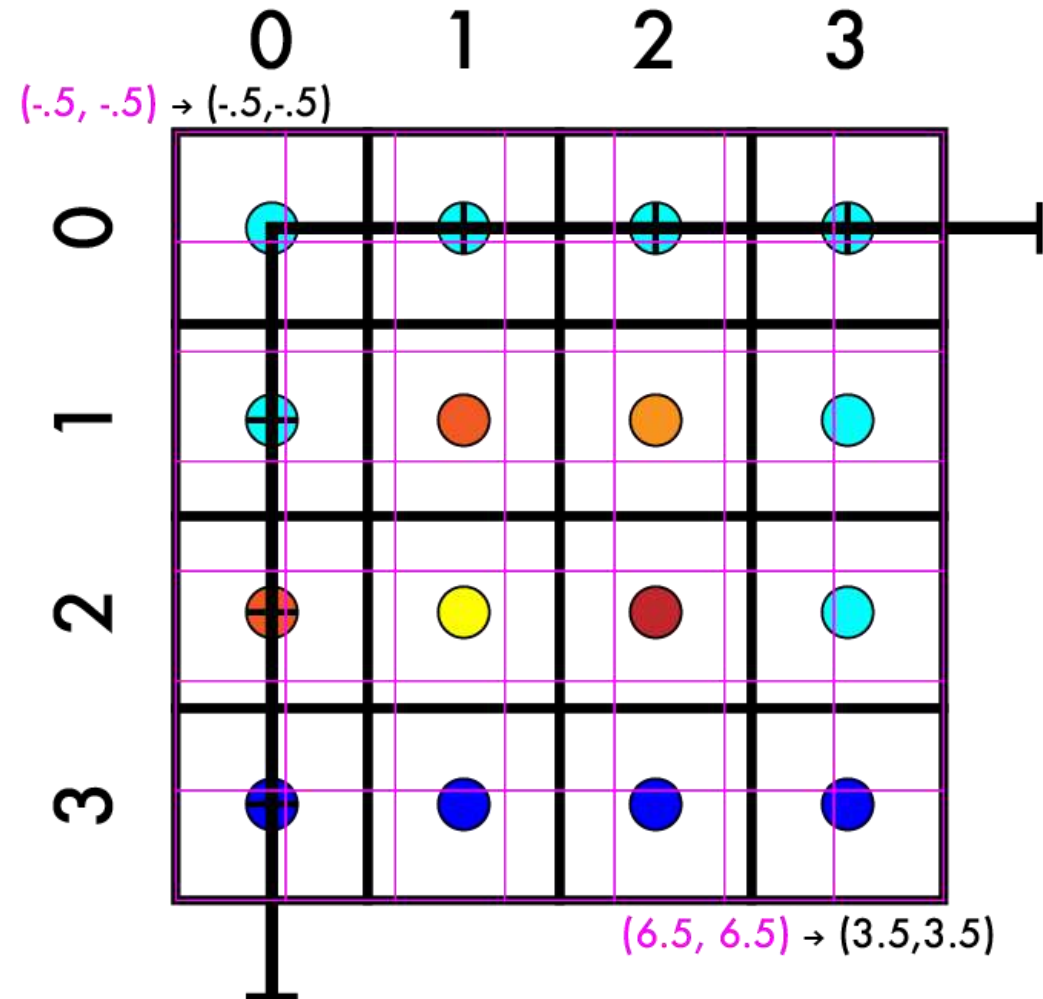
## Resize 4x4 -> 7x7

- Create our new image
- Match up coordinates
  - System of equations
  - $aX + b = Y$
  - $a * -.5 + b = -.5$
  - $a * 6.5 + b = 3.5$ 
    - $a * 7 = 4$



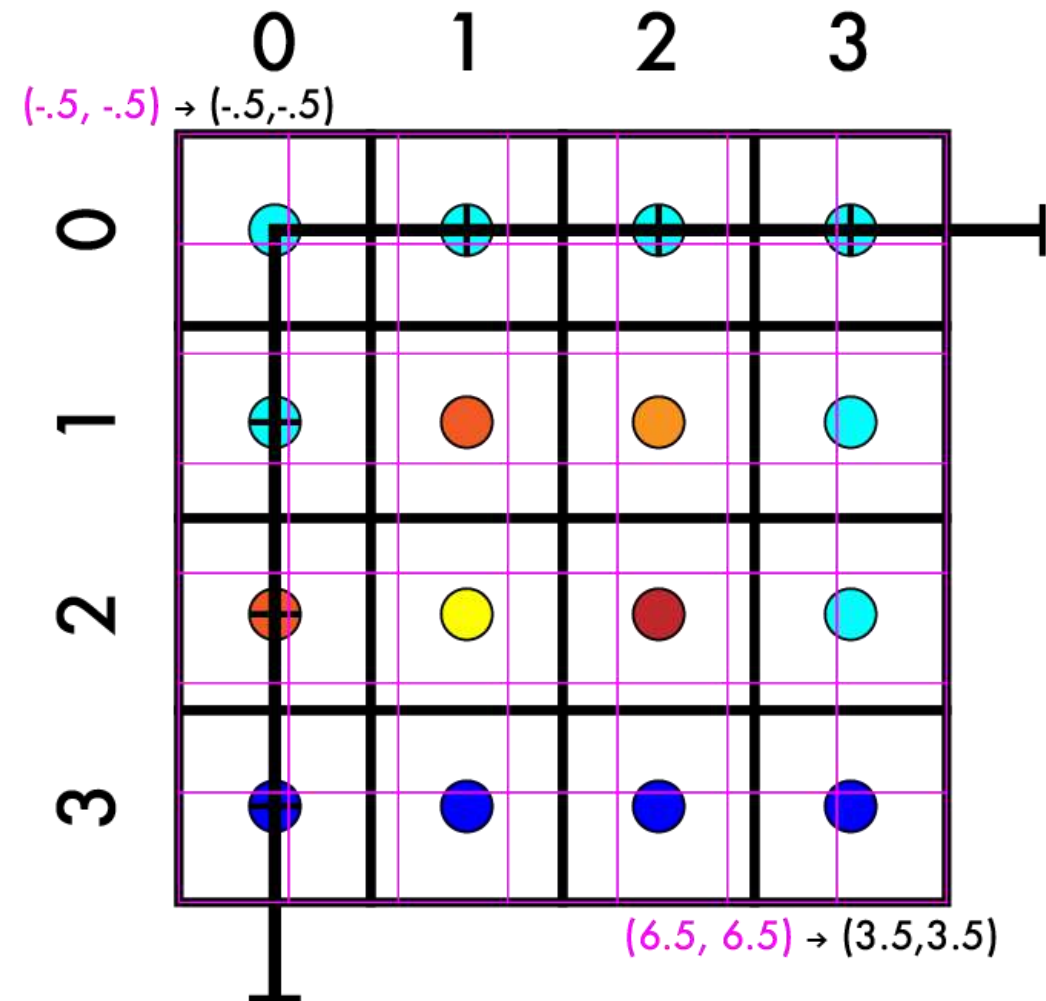
## Resize 4x4 -> 7x7

- Create our new image
- Match up coordinates
  - System of equations
  - $aX + b = Y$
  - $a * -.5 + b = -.5$
  - $a * 6.5 + b = 3.5$ 
    - $a * 7 = 4$
    - $a = 4/7$



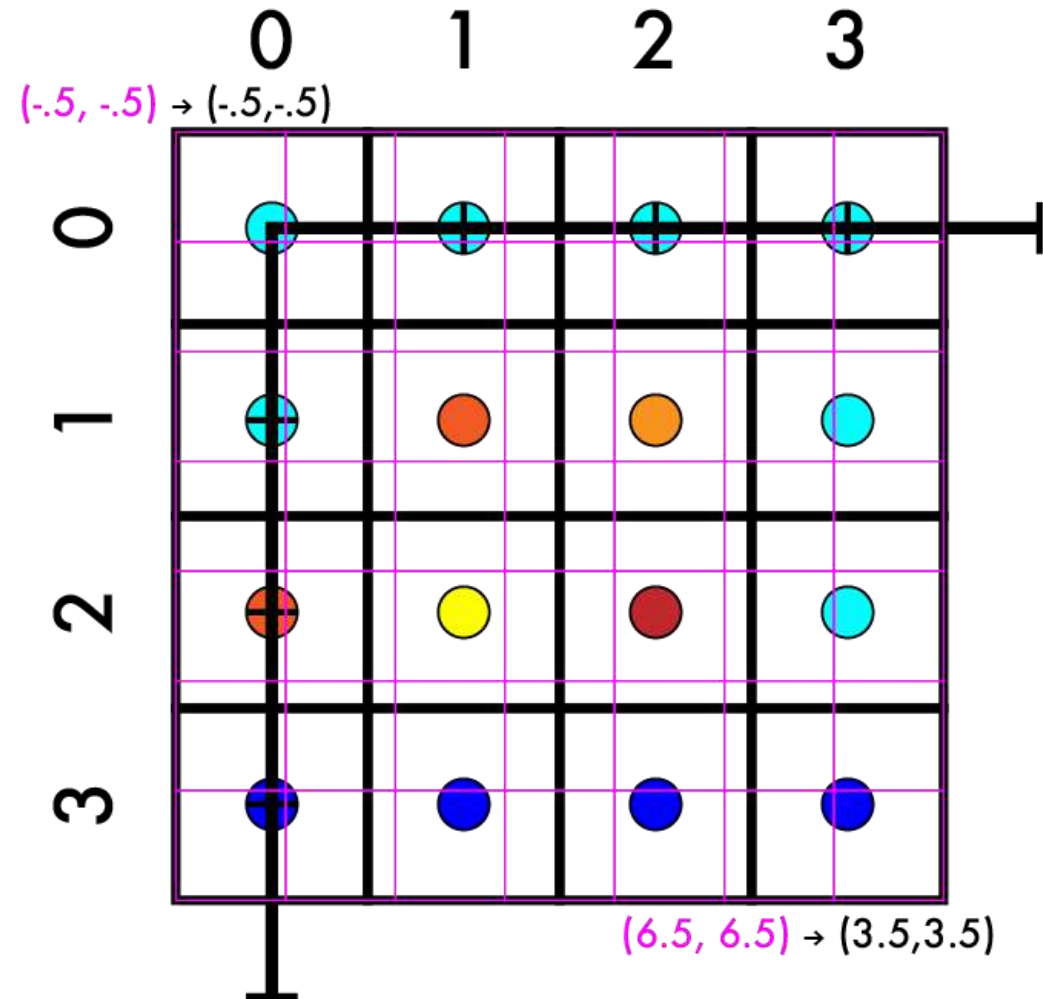
## Resize 4x4 -> 7x7

- Create our new image
- Match up coordinates
  - System of equations
  - $aX + b = Y$
  - $a \cdot -.5 + b = -.5$
  - $a \cdot 6.5 + b = 3.5$
  - $a = 4/7$



## Resize 4x4 -> 7x7

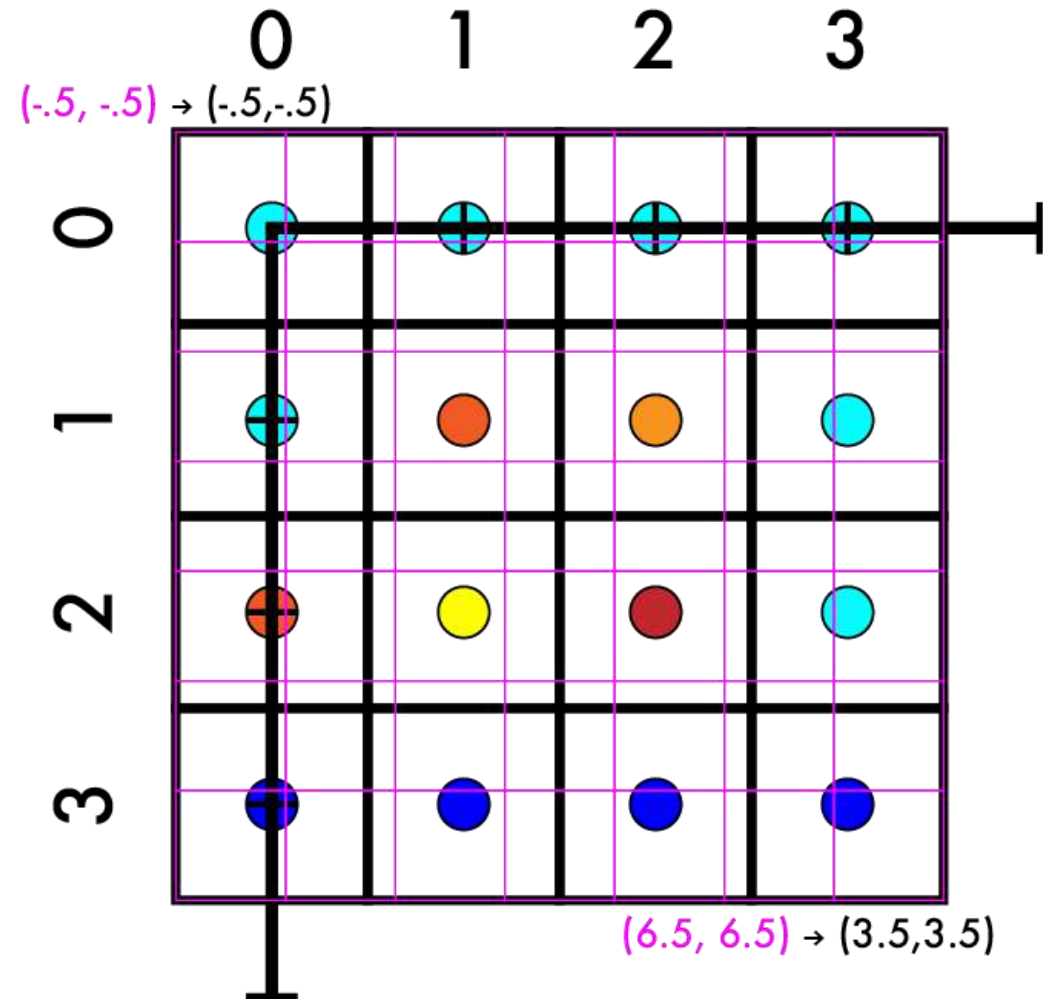
- Create our new image
- Match up coordinates
  - System of equations
  - $aX + b = Y$
  - $a \cdot -.5 + b = -.5$
  - $a \cdot 6.5 + b = 3.5$
  - $a = 4/7$ 
    - $a \cdot -.5 + b = -.5$





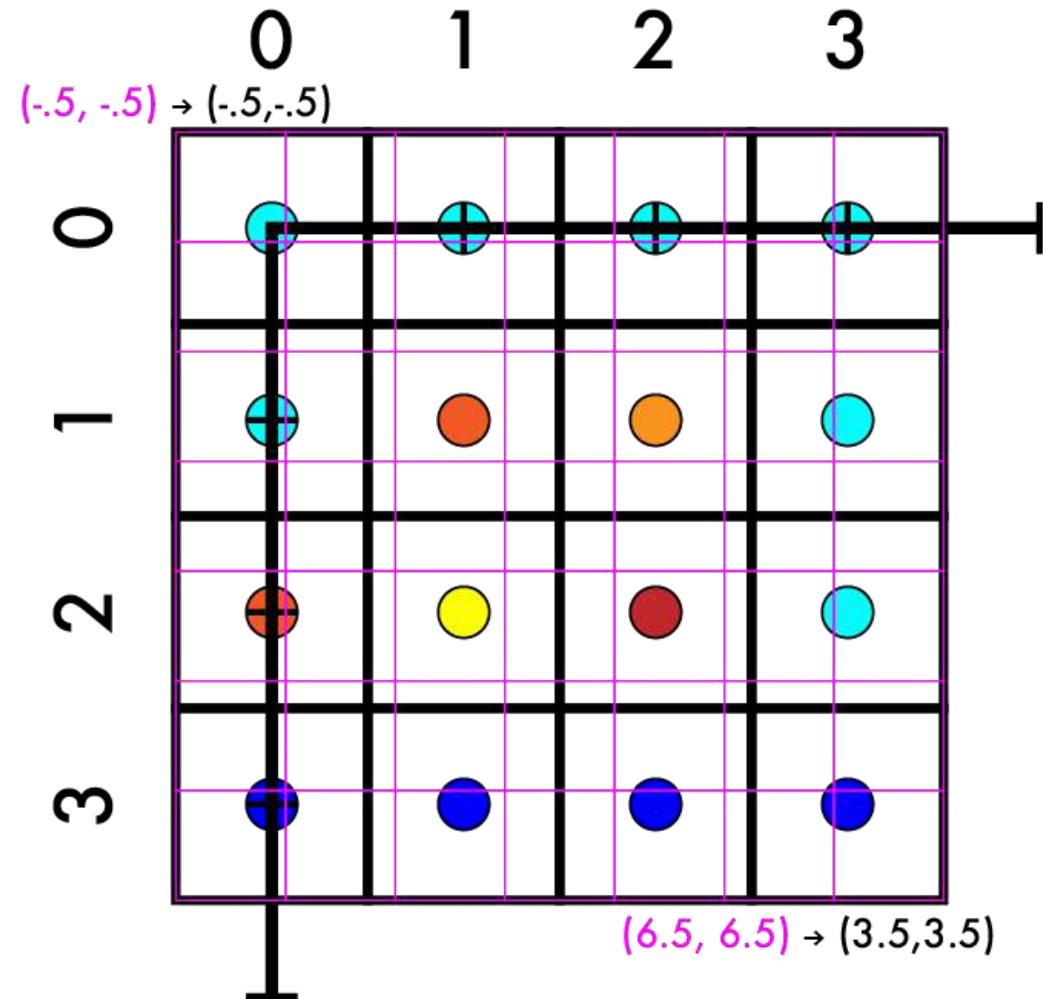
## Resize 4x4 -> 7x7

- Create our new image
- Match up coordinates
  - System of equations
  - $aX + b = Y$
  - $a \cdot -.5 + b = -.5$
  - $a \cdot 6.5 + b = 3.5$
  - $a = 4/7$ 
    - $a \cdot -.5 + b = -.5$
    - $4/7 \cdot -1/2 + b = -1/2$



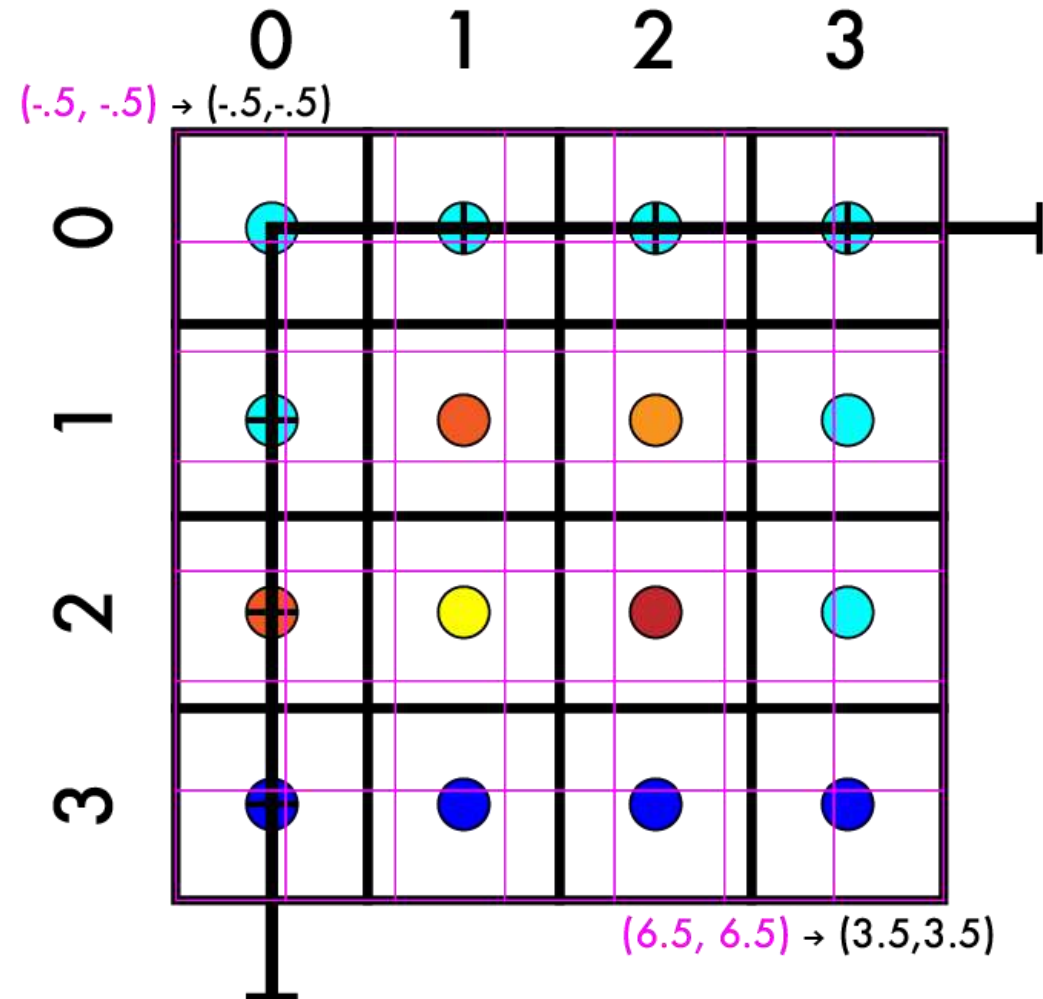
## Resize 4x4 -> 7x7

- Create our new image
- Match up coordinates
  - System of equations
  - $aX + b = Y$
  - $a \cdot -.5 + b = -.5$
  - $a \cdot 6.5 + b = 3.5$
  - $a = 4/7$ 
    - $a \cdot -.5 + b = -.5$
    - $4/7 \cdot -1/2 + b = -1/2$
    - $-4/14 + b = -7/14$



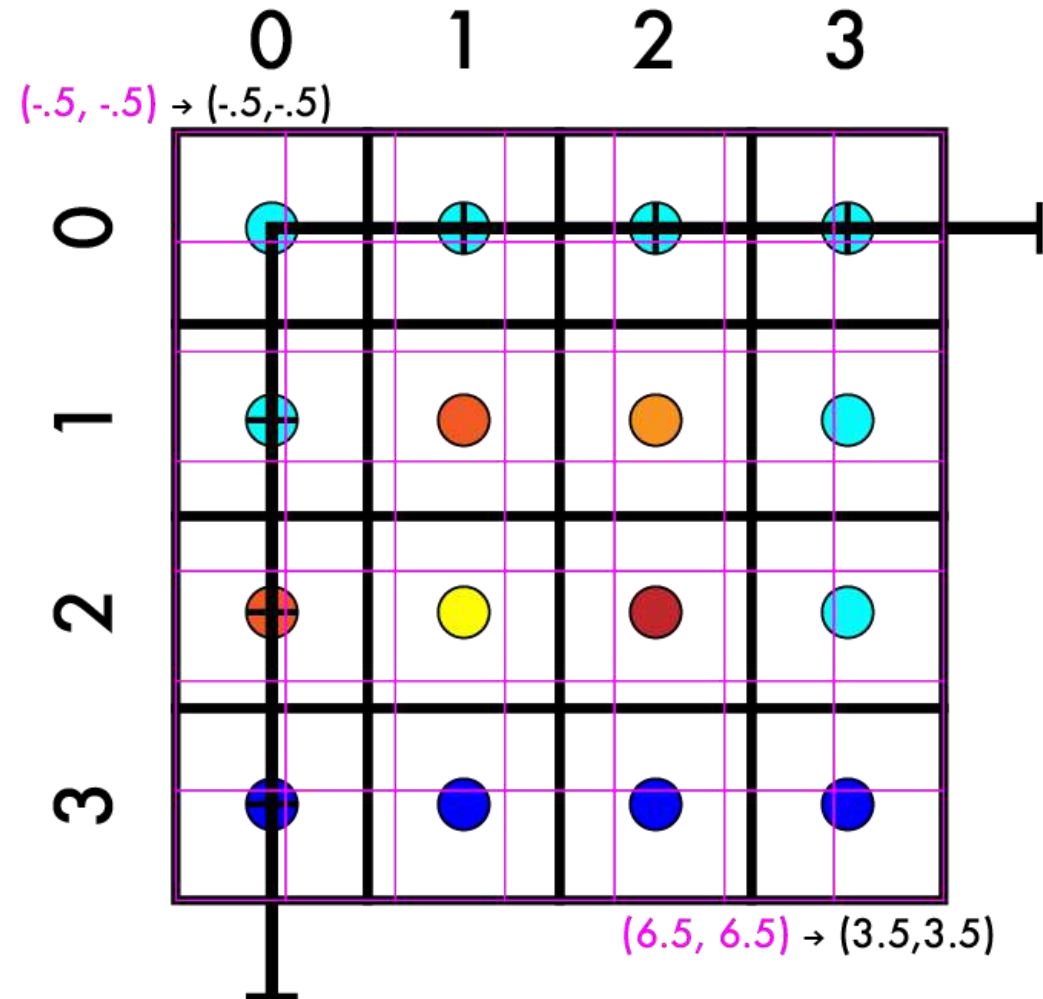
## Resize 4x4 -> 7x7

- Create our new image
- Match up coordinates
  - System of equations
  - $aX + b = Y$
  - $a \cdot -.5 + b = -.5$
  - $a \cdot 6.5 + b = 3.5$
  - $a = 4/7$ 
    - $a \cdot -.5 + b = -.5$
    - $4/7 \cdot -1/2 + b = -1/2$
    - $-4/14 + b = -7/14$
    - $b = -3/14$



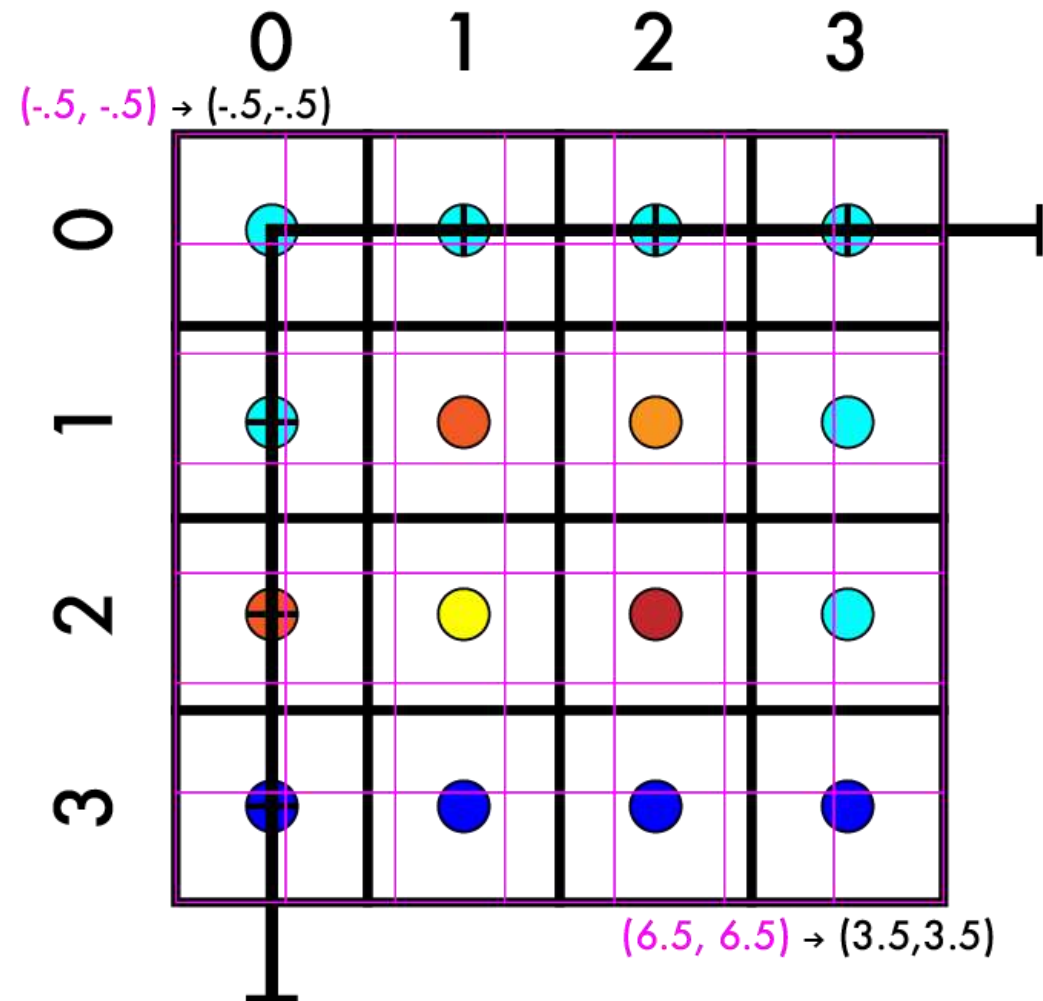
## Resize 4x4 -> 7x7

- Create our new image
- Match up coordinates
  - System of equations
  - $aX + b = Y$
  - $a \cdot -.5 + b = -.5$
  - $a \cdot 6.5 + b = 3.5$
  - $a = 4/7$
  - $b = -3/14$



## Resize 4x4 -> 7x7

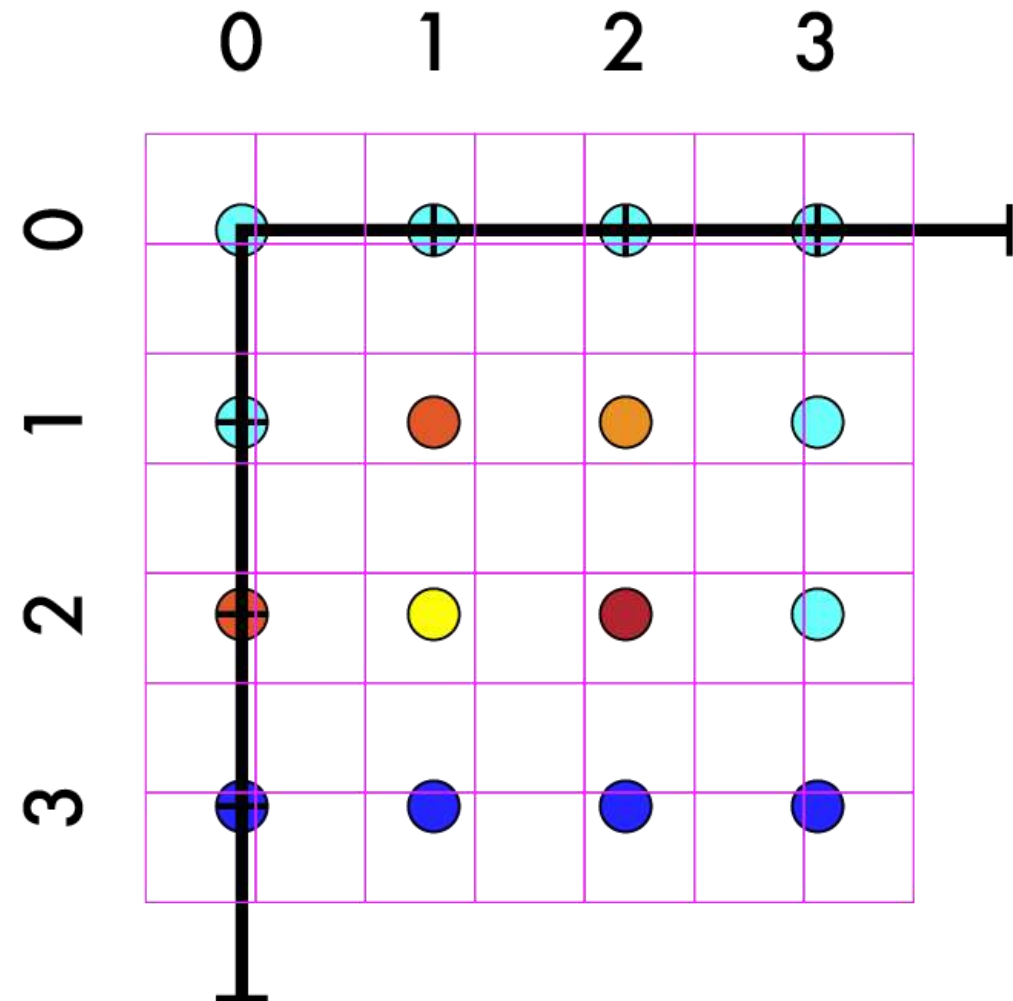
- Create our new image
- Match up coordinates
  - $\frac{4}{7} X - \frac{3}{14} = Y$





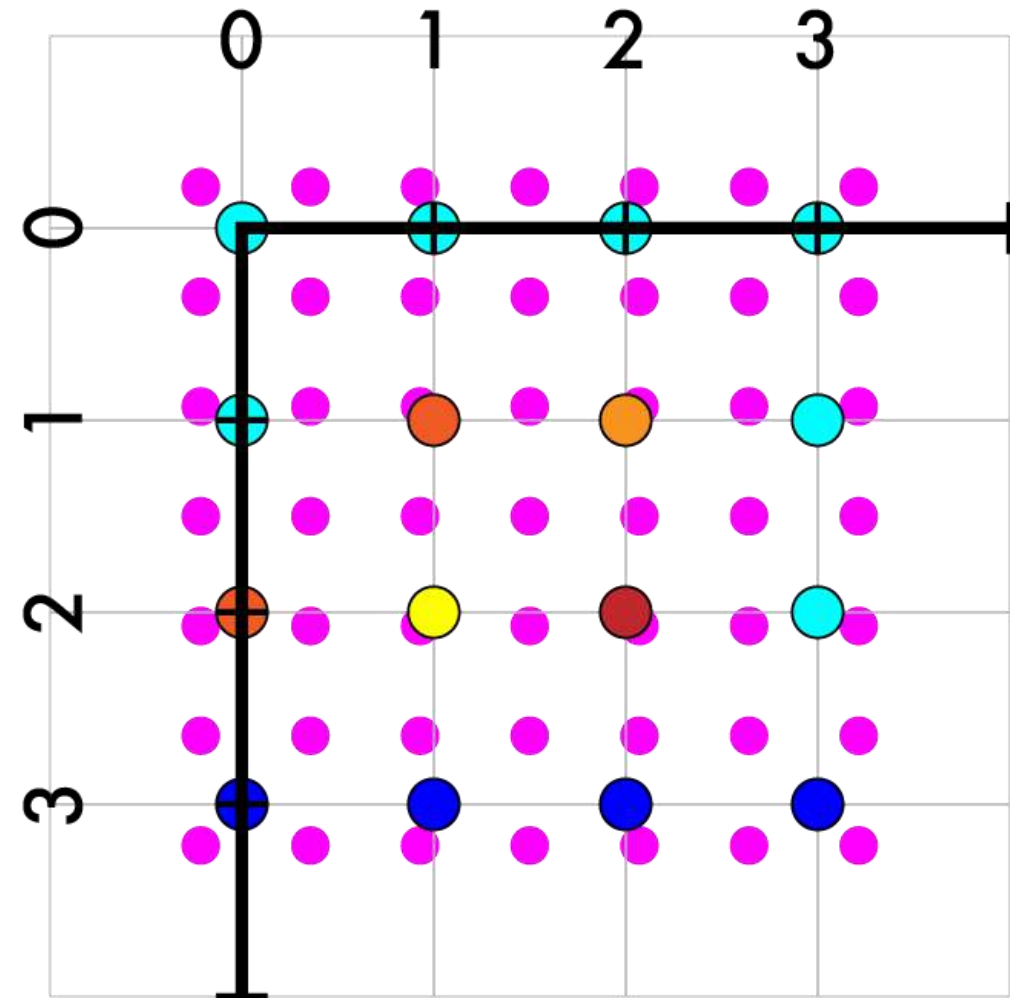
## Resize 4x4 -> 7x7

- Create our new image
- Match up coordinates
  - $\frac{4}{7} X - \frac{3}{14} = Y$



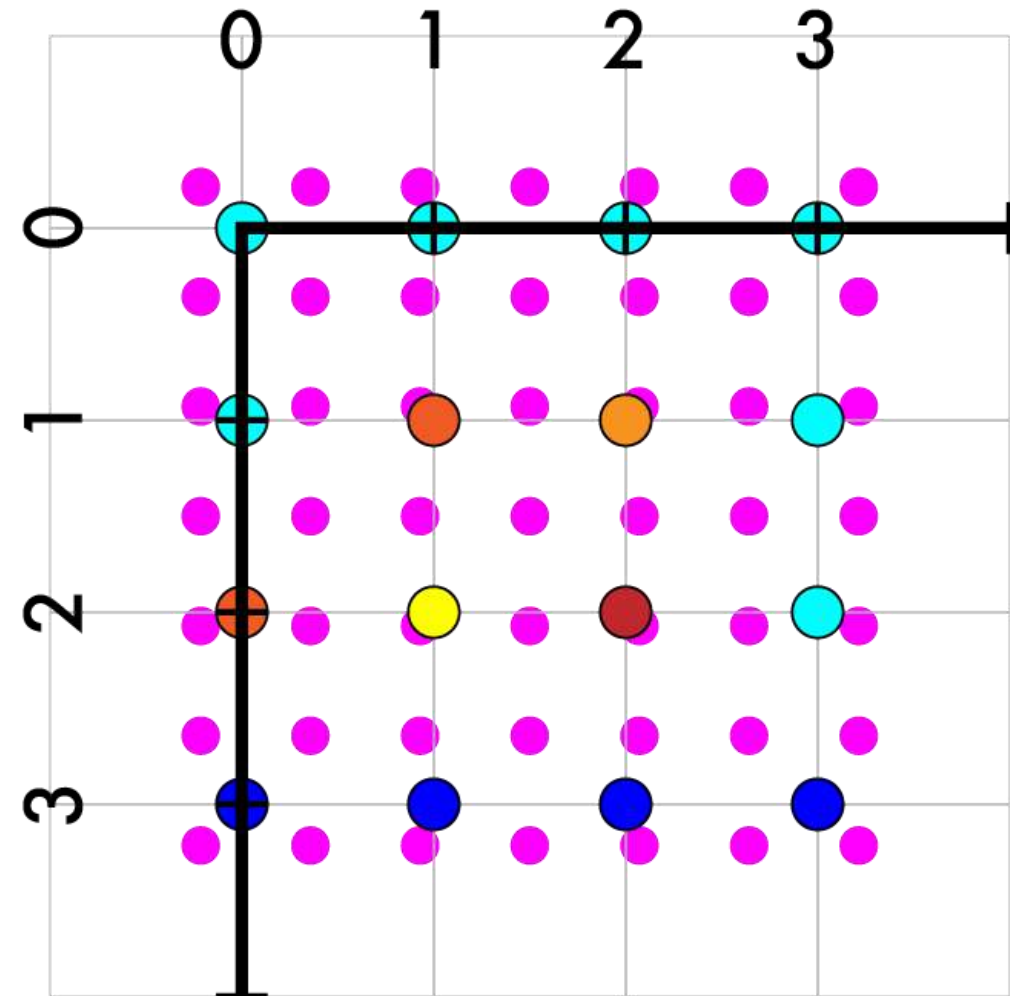
## Resize 4x4 -> 7x7

- Create our new image
- Match up coordinates
  - $\frac{4}{7} X - \frac{3}{14} = Y$
- Iterate over new pts



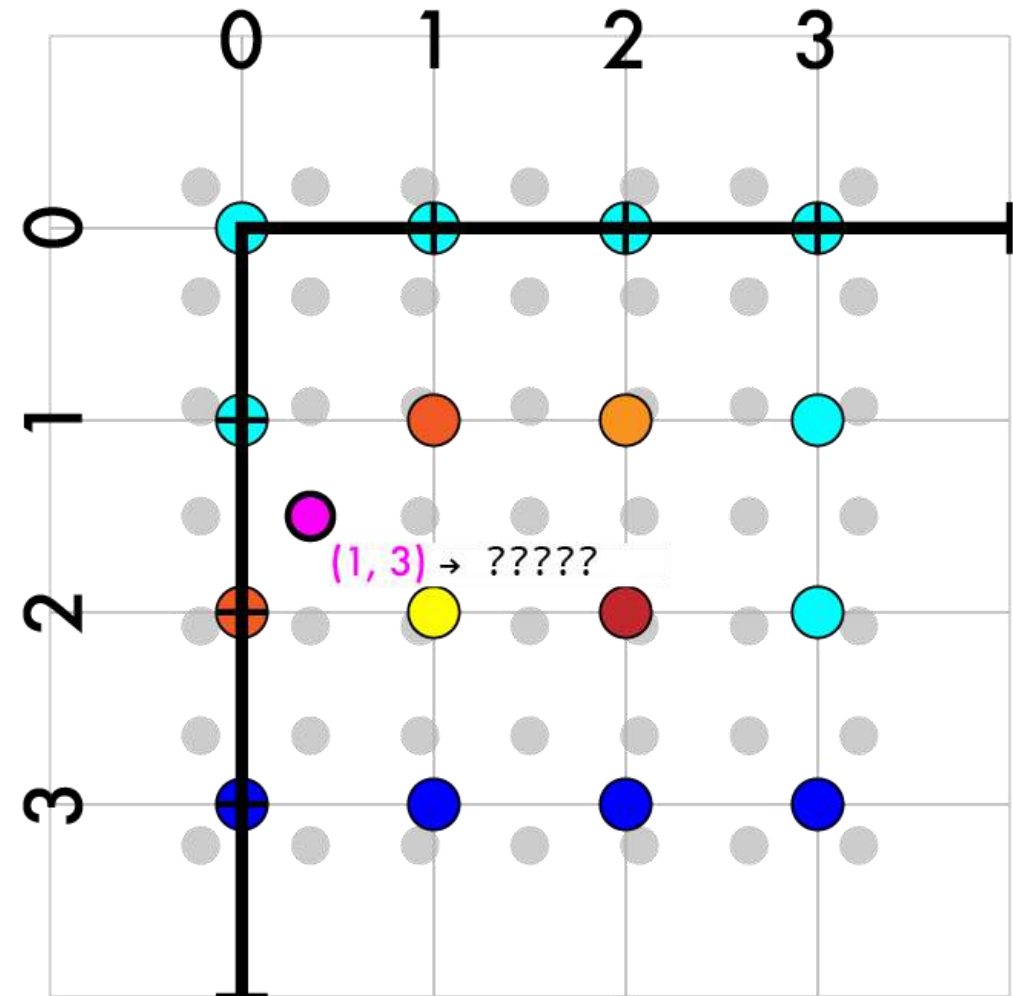
## Resize 4x4 -> 7x7

- Create our new image
- Match up coordinates
  - $4/7 X - 3/14 = Y$
- Iterate over new pts
  - Map to old coords



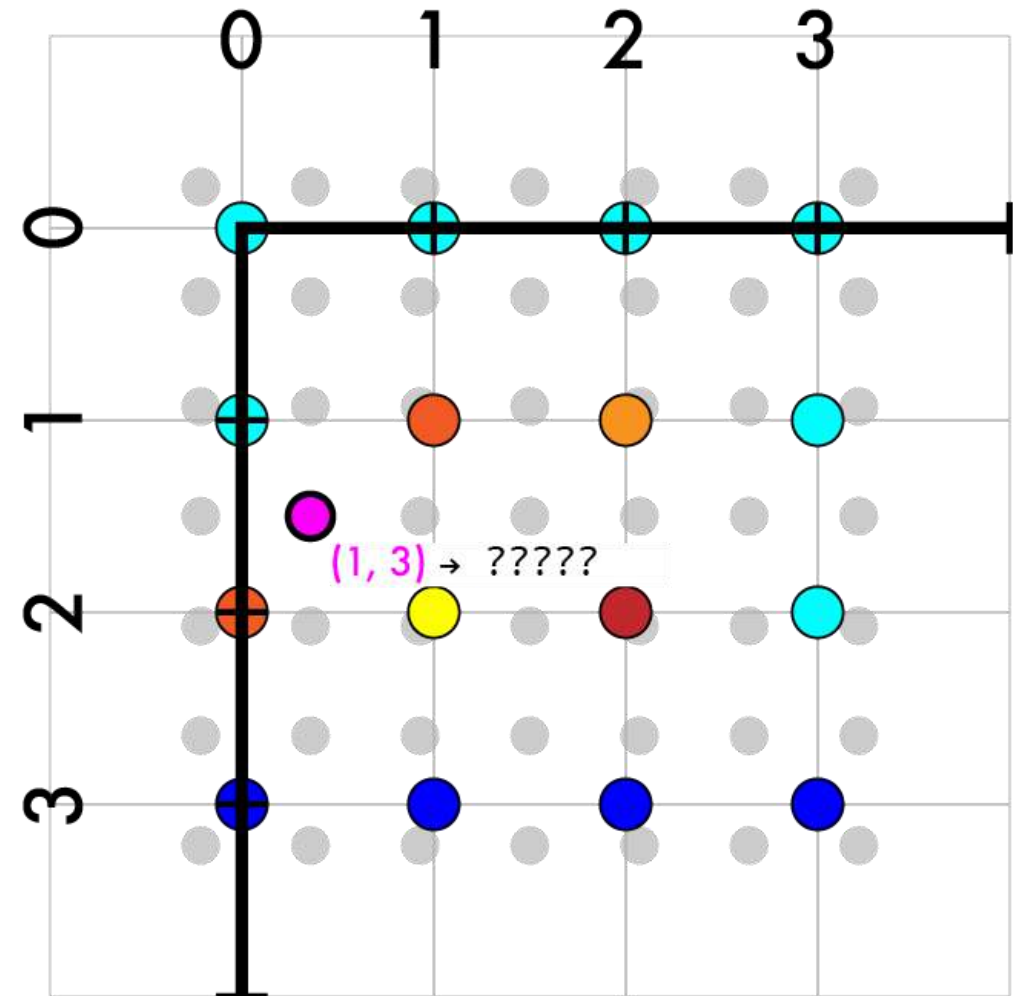
## Resize 4x4 -> 7x7

- Create our new image
- Match up coordinates
  - $\frac{4}{7} X - \frac{3}{14} = Y$
- Iterate over new pts
  - Map to old coords
  - (1, 3)



## Resize 4x4 -> 7x7

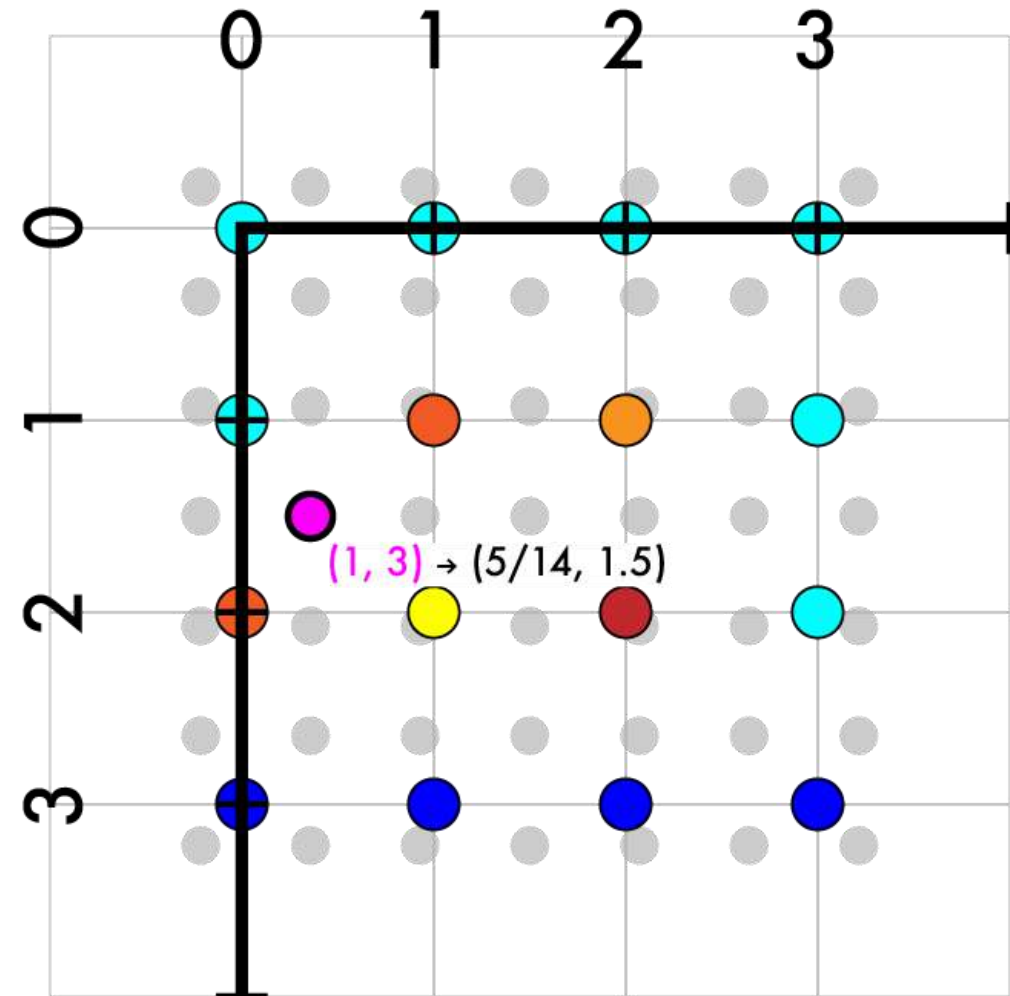
- Create our new image
- Match up coordinates
  - $4/7 X - 3/14 = Y$
- Iterate over new pts
  - Map to old coords
  - (1, 3)
  - $4/7 * 1 - 3/14$
  - $4/7 * 3 - 3/14$





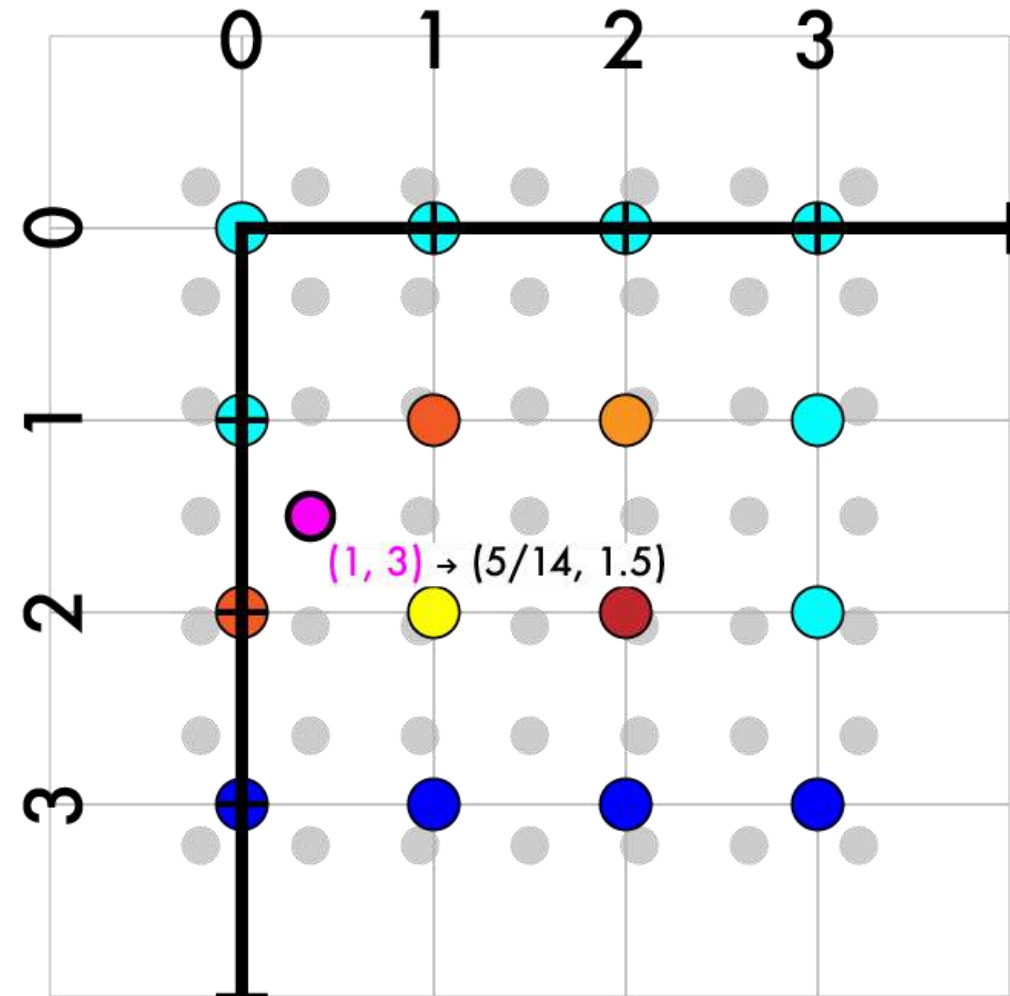
## Resize 4x4 -> 7x7

- Create our new image
- Match up coordinates
  - $4/7 X - 3/14 = Y$
- Iterate over new pts
  - Map to old coords
  - (1, 3)
  - $4/7 * 1 - 3/14$
  - $4/7 * 3 - 3/14$
  - (5/14, 21/14)



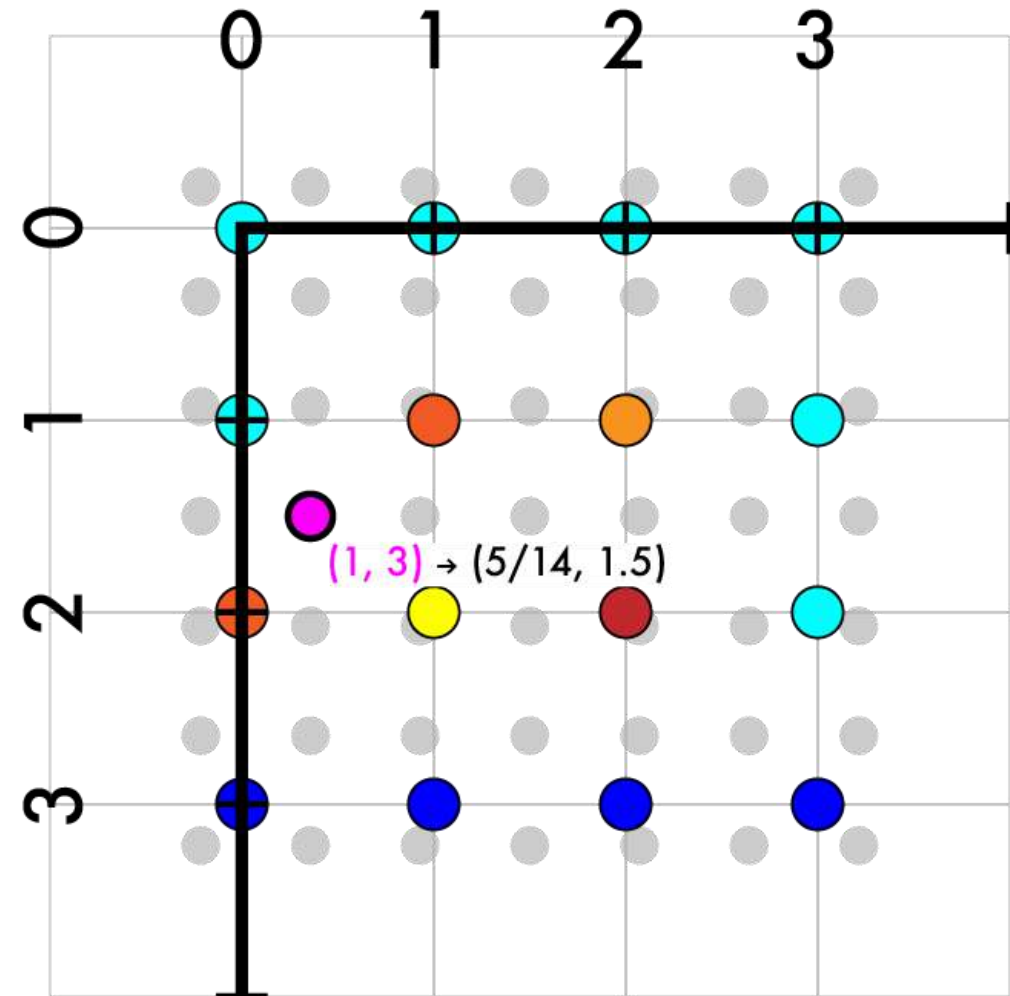
## Resize 4x4 -> 7x7

- Create our new image
- Match up coordinates
  - $4/7 X - 3/14 = Y$
- Iterate over new pts
  - Map to old coords
  - $(1, 3) \rightarrow (5/14, 21/14)$



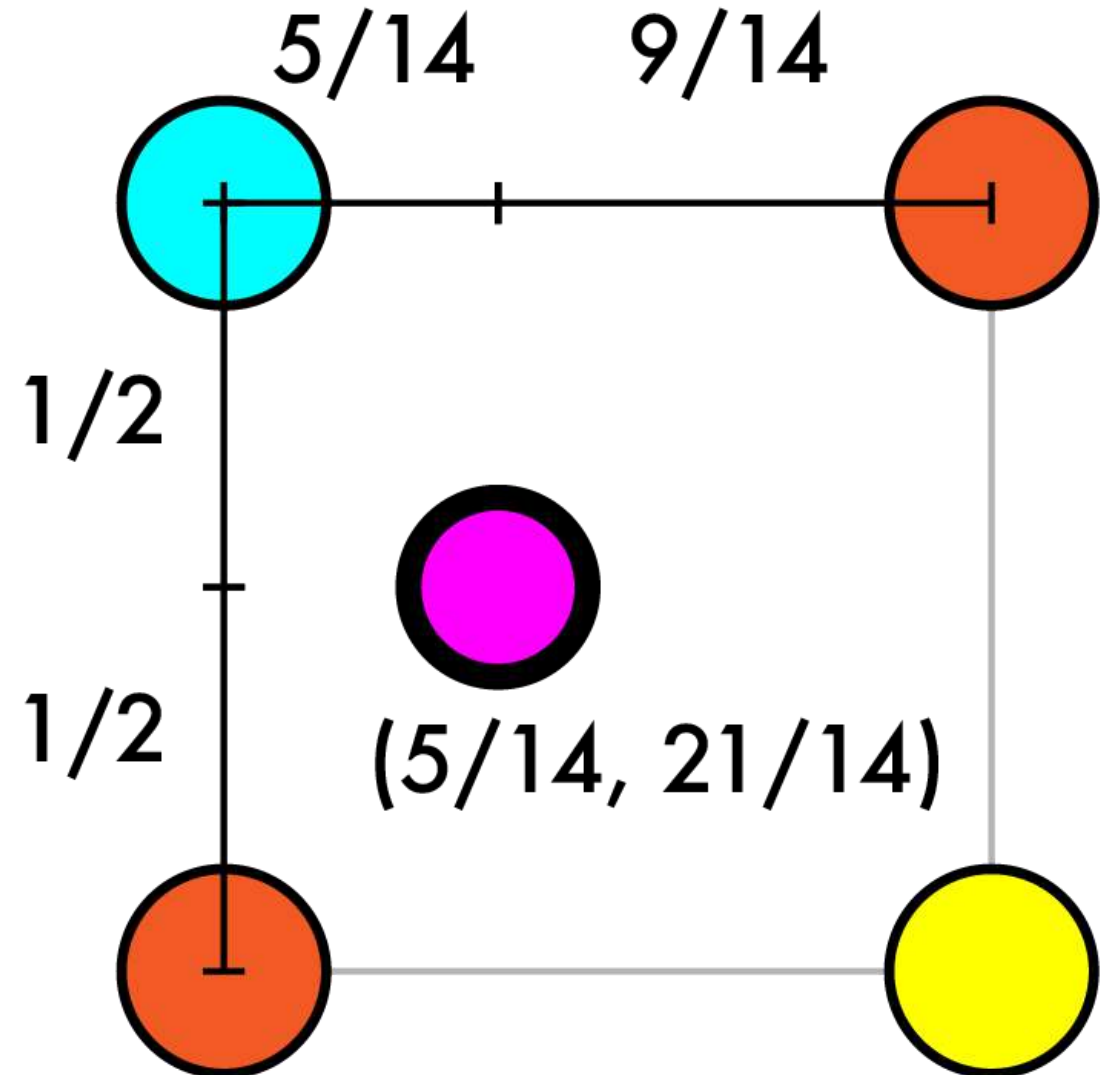
## Resize 4x4 -> 7x7

- Create our new image
- Match up coordinates
  - $\frac{4}{7} X - \frac{3}{14} = Y$
- Iterate over new pts
  - Map to old coords
  - $(1, 3) \rightarrow (5/14, 21/14)$
  - Interpolate old values



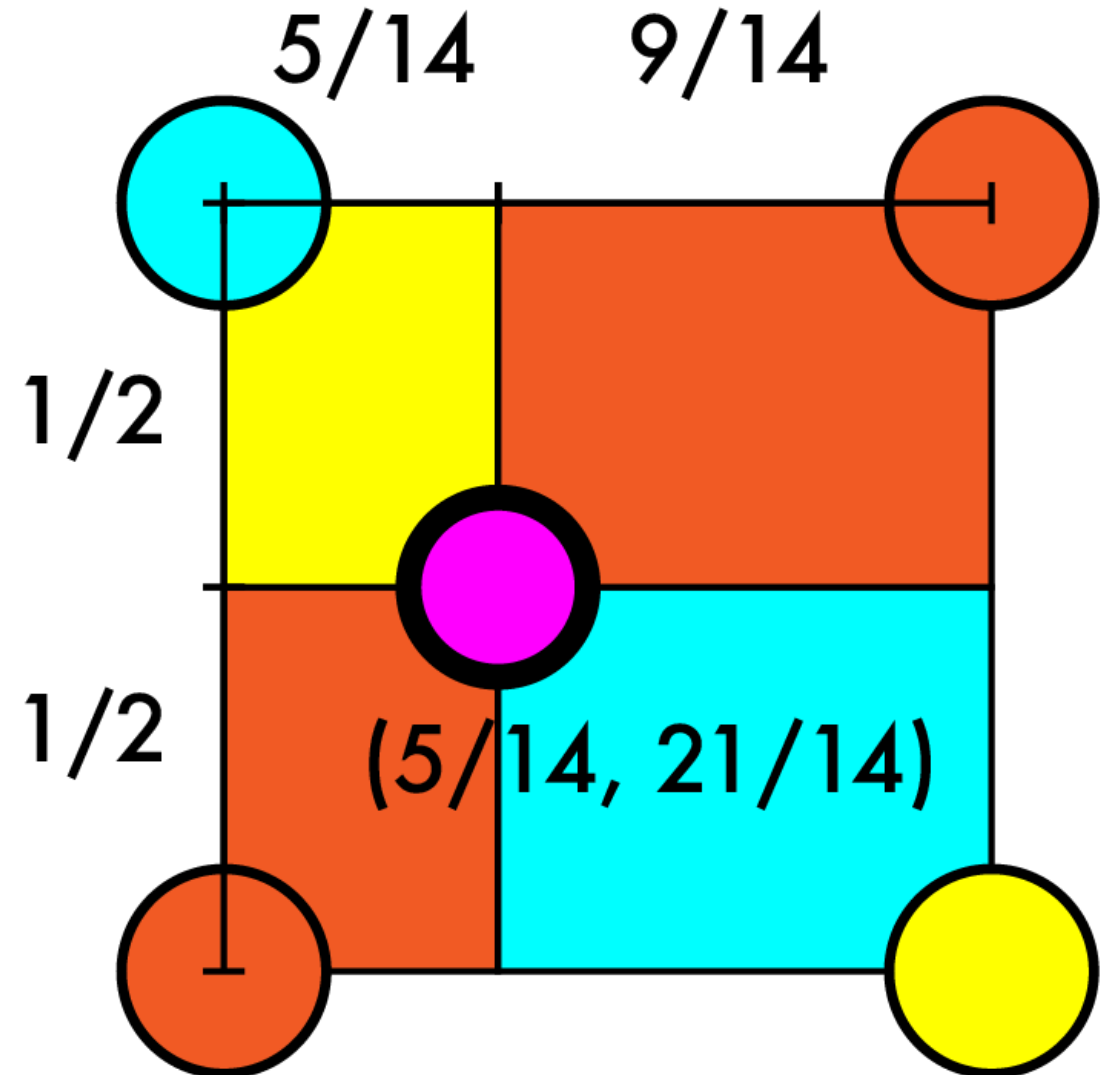
## Resize 4x4 -> 7x7

- Create our new image
- Match up coordinates
  - $4/7 X - 3/14 = Y$
- Iterate over new pts
  - Map to old coords
  - $(1, 3) \rightarrow (5/14, 21/14)$
  - Interpolate old values



## Resize 4x4 -> 7x7

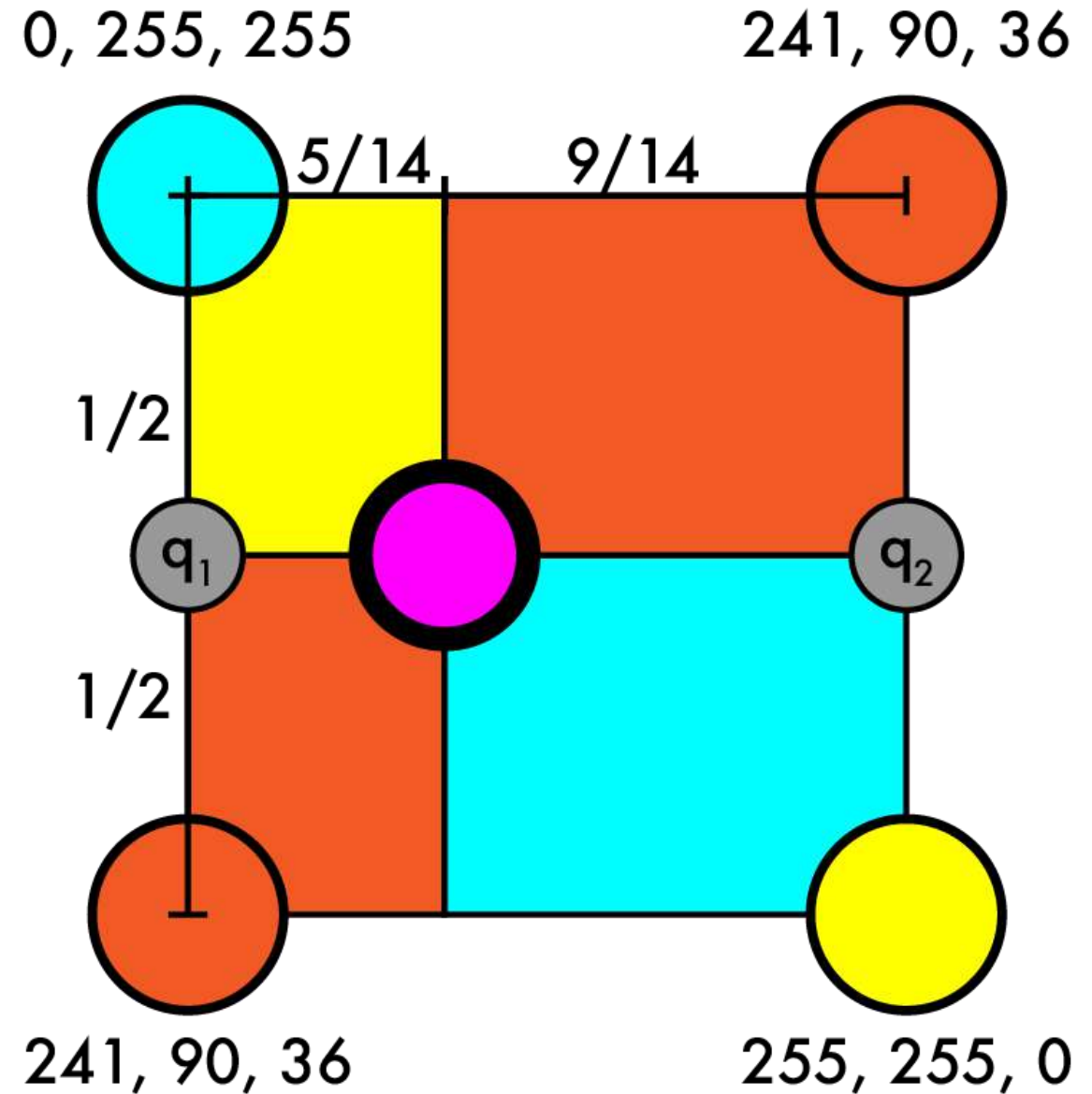
- Create our new image
- Match up coordinates
  - $4/7 X - 3/14 = Y$
- Iterate over new pts
  - Map to old coords
  - $(1, 3) \rightarrow (5/14, 21/14)$
  - Interpolate old values
    - Size of opposite rects





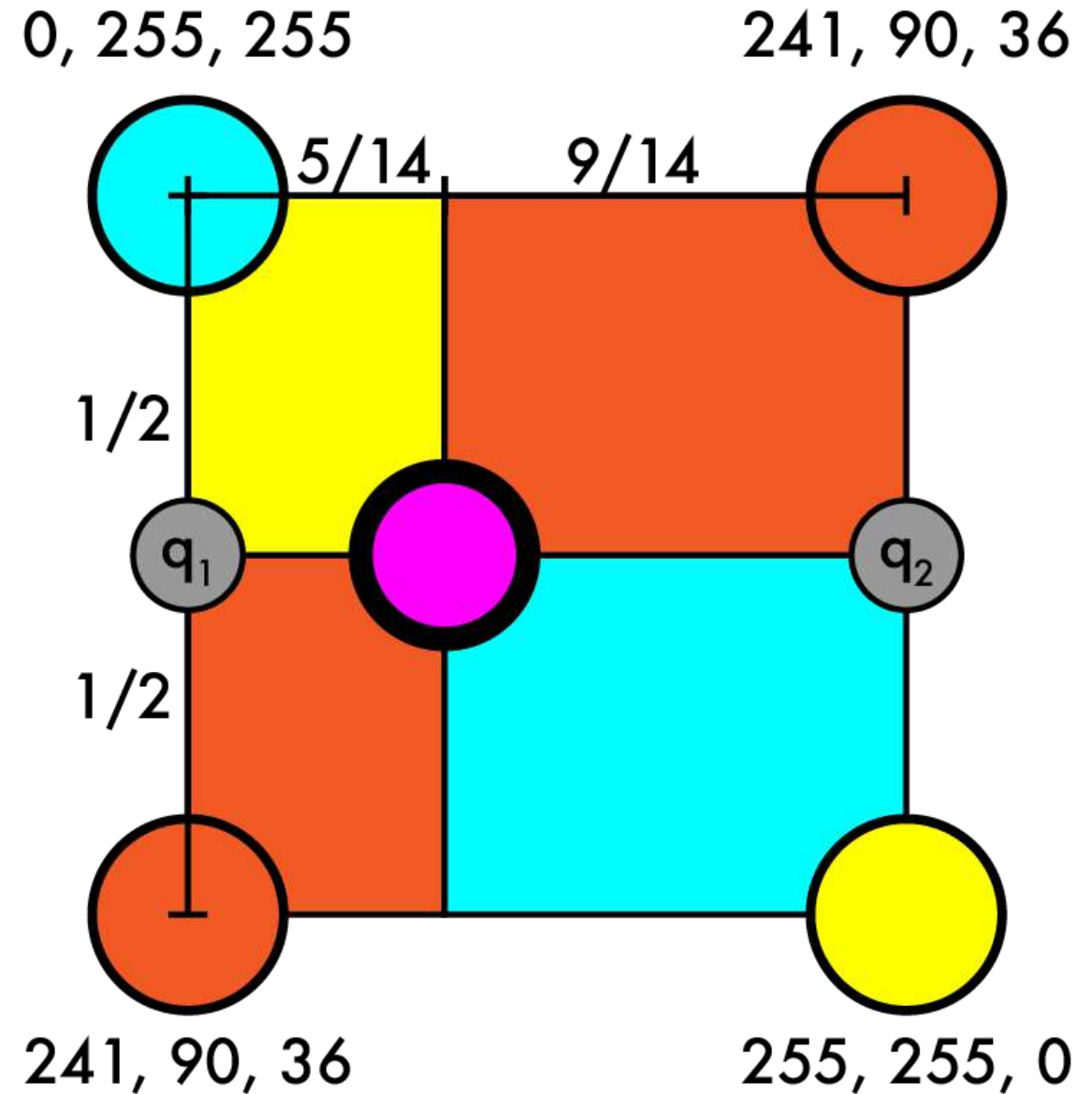
## Resize 4x4 -> 7x7

- Create our new image
- Match up coordinates
  - $4/7 X - 3/14 = Y$
- Iterate over new pts
  - Map to old coords
  - $(1, 3) \rightarrow (5/14, 21/14)$
  - Interpolate old values
    - Size of opposite rects
    - OR find  $q_1$  and  $q_2$ , then interpolate between them



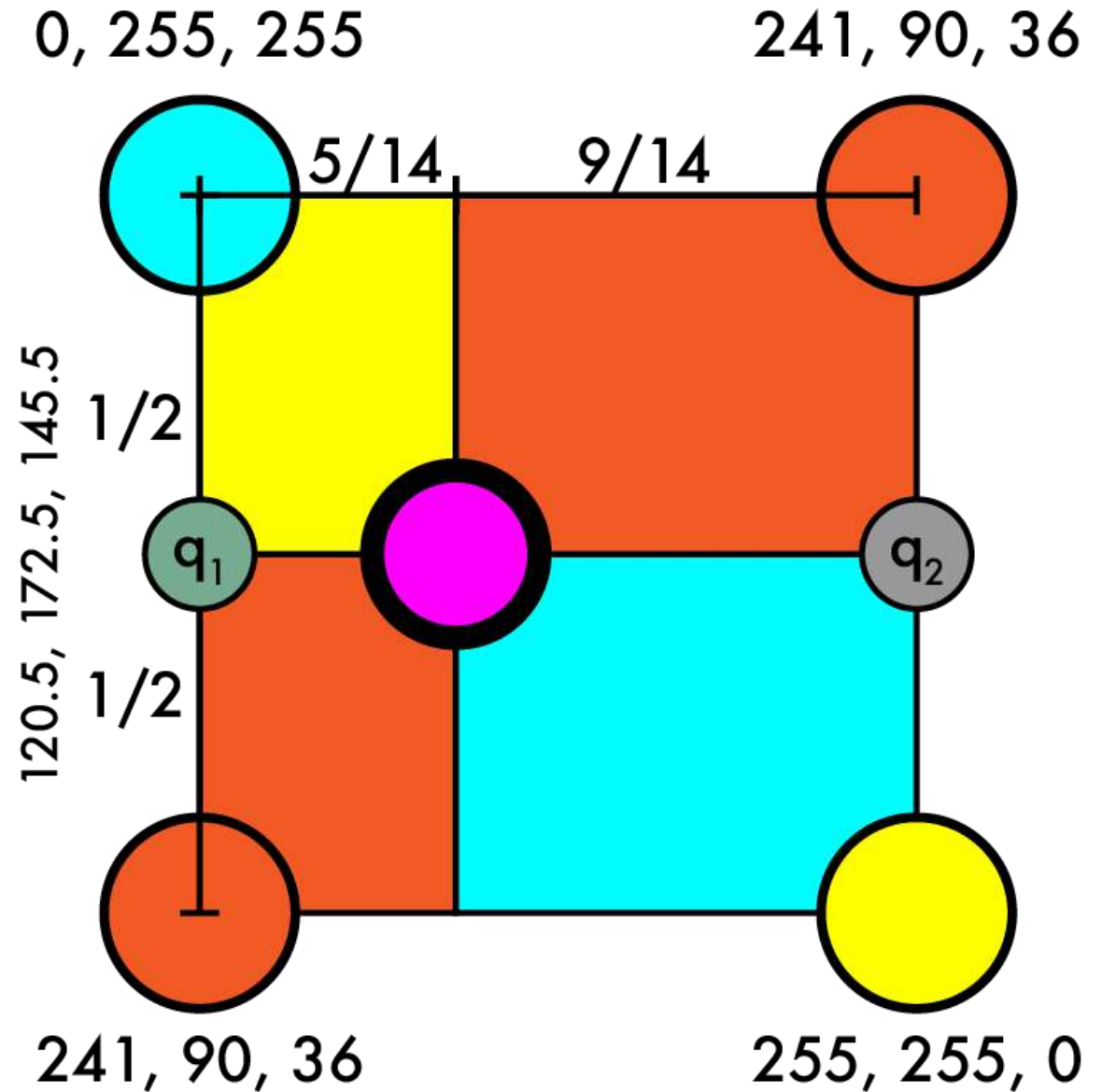
## Resize 4x4 -> 7x7

- Create our new image
- Match up coordinates
  - $4/7 X - 3/14 = Y$
- Iterate over new pts
  - Map to old coords
  - $(1, 3) \rightarrow (5/14, 21/14)$
  - Interpolate old values
    - $q1 = r1, g1, b1$
    - $r1 = .5 * 0 + .5 * 241$
    - $g1 = .5 * 255 + .5 * 90$
    - $b1 = .5 * 255 + .5 * 36$



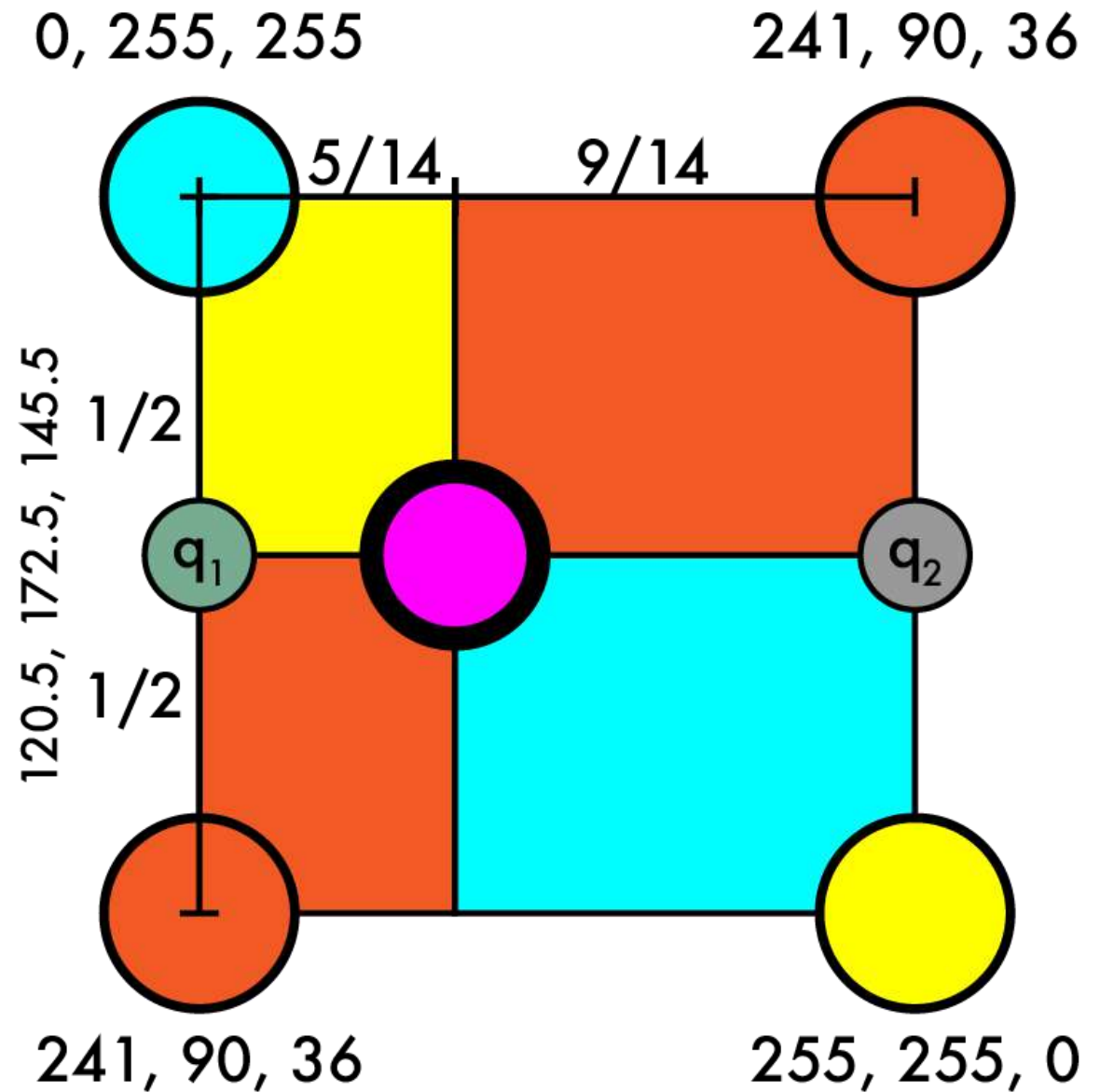
## Resize 4x4 -> 7x7

- Create our new image
- Match up coordinates
  - $4/7 X - 3/14 = Y$
- Iterate over new pts
  - Map to old coords
  - $(1, 3) \rightarrow (5/14, 21/14)$
  - Interpolate old values
    - $q1 = (120.5, 172.5, 145.5)$
    - $q2 = r2, g2, b2$



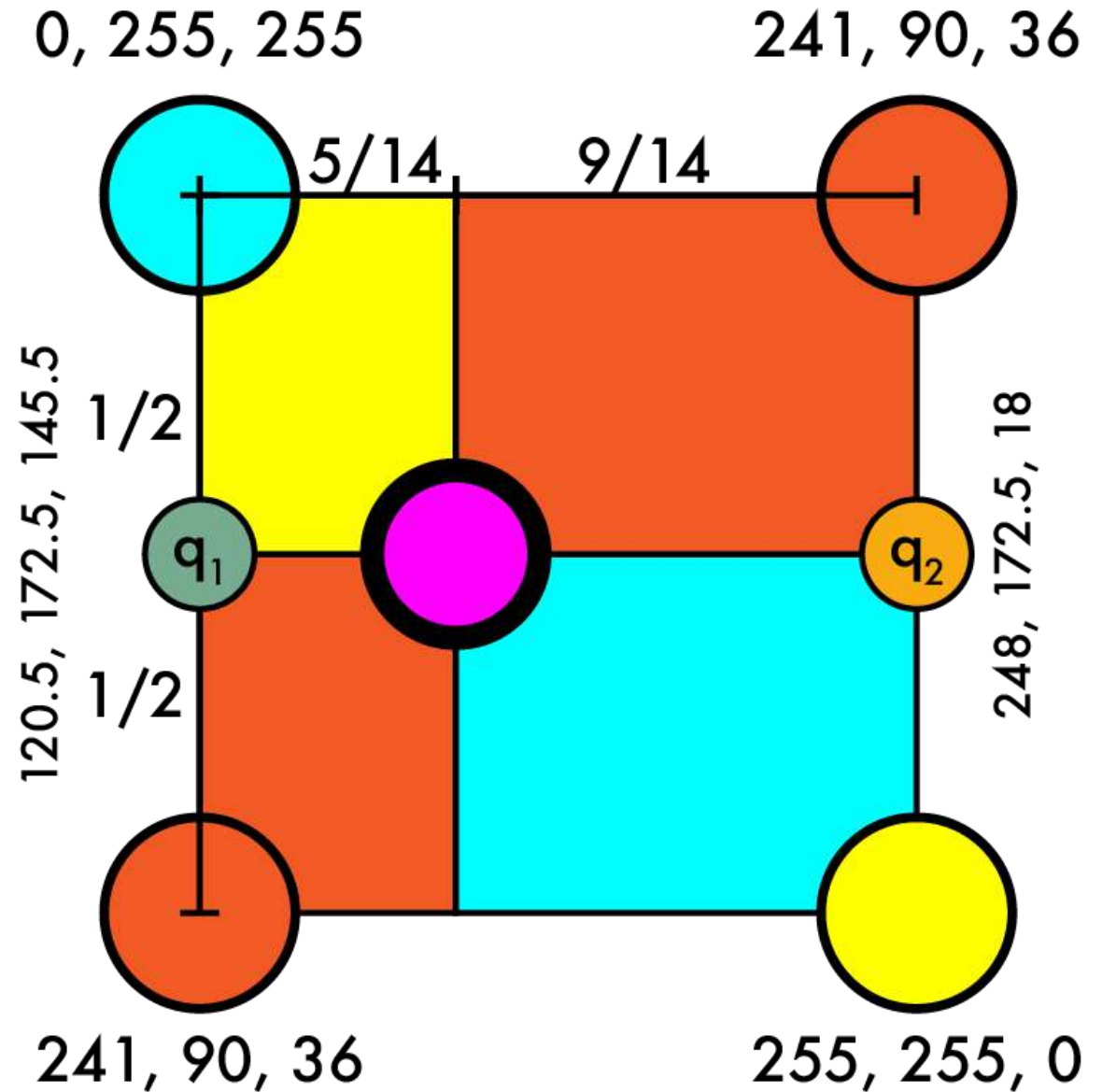
## Resize 4x4 -> 7x7

- Create our new image
- Match up coordinates
  - $4/7 X - 3/14 = Y$
- Iterate over new pts
  - Map to old coords
  - $(1, 3) \rightarrow (5/14, 21/14)$
  - Interpolate old values
    - $q1 = (120.5, 172.5, 145.5)$
    - $q2 = r2, g2, b2$
    - $r2 = .5 * 241 + .5 * 255$
    - $g2 = .5 * 90 + .5 * 255$
    - $b2 = .5 * 36 + .5 * 0$



## Resize 4x4 -> 7x7

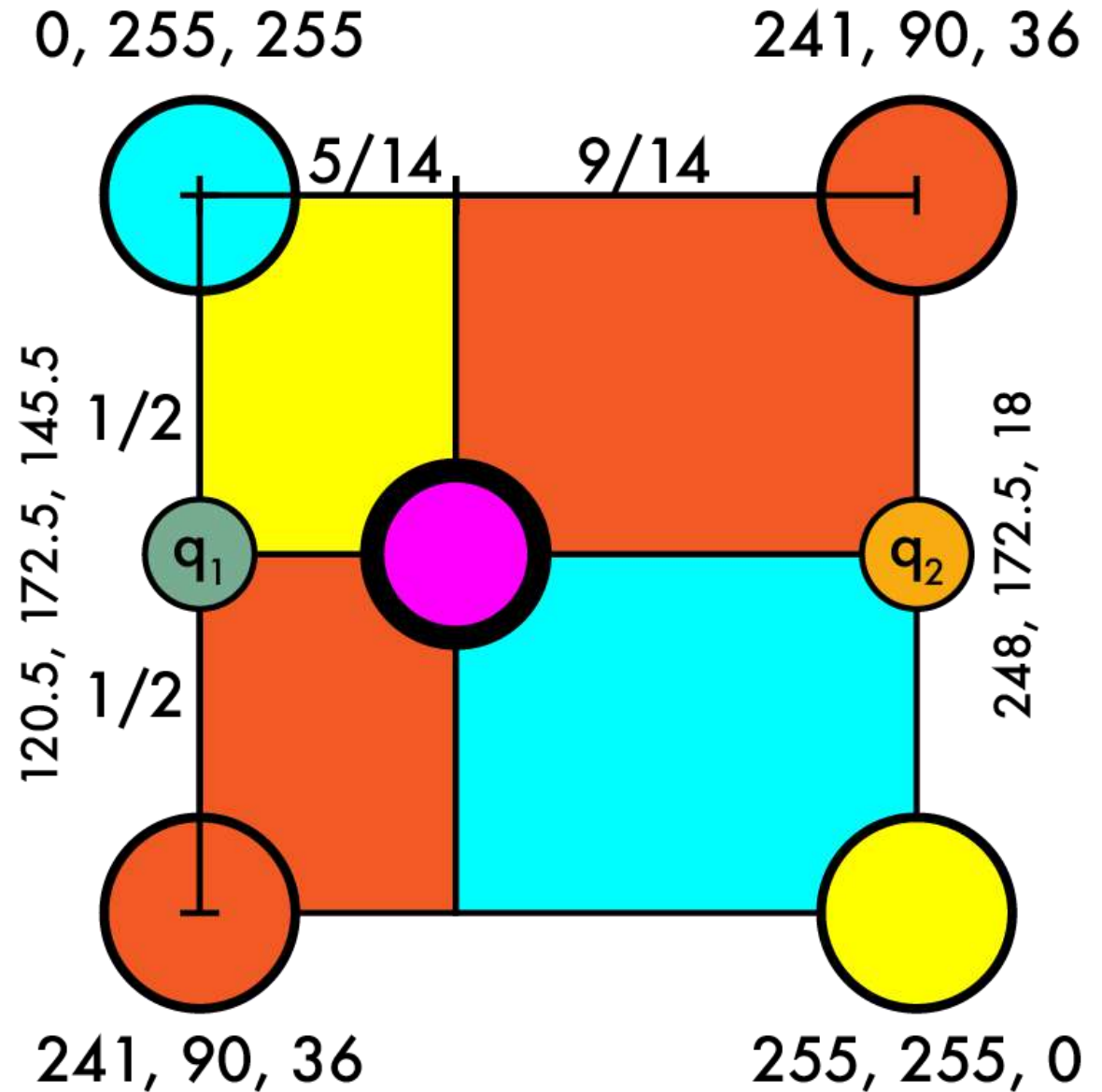
- Create our new image
- Match up coordinates
  - $4/7 X - 3/14 = Y$
- Iterate over new pts
  - Map to old coords
  - $(1, 3) \rightarrow (5/14, 21/14)$
  - Interpolate old values
    - $q_1 = (120.5, 172.5, 145.5)$
    - $q_2 = (248, 172.5, 18)$





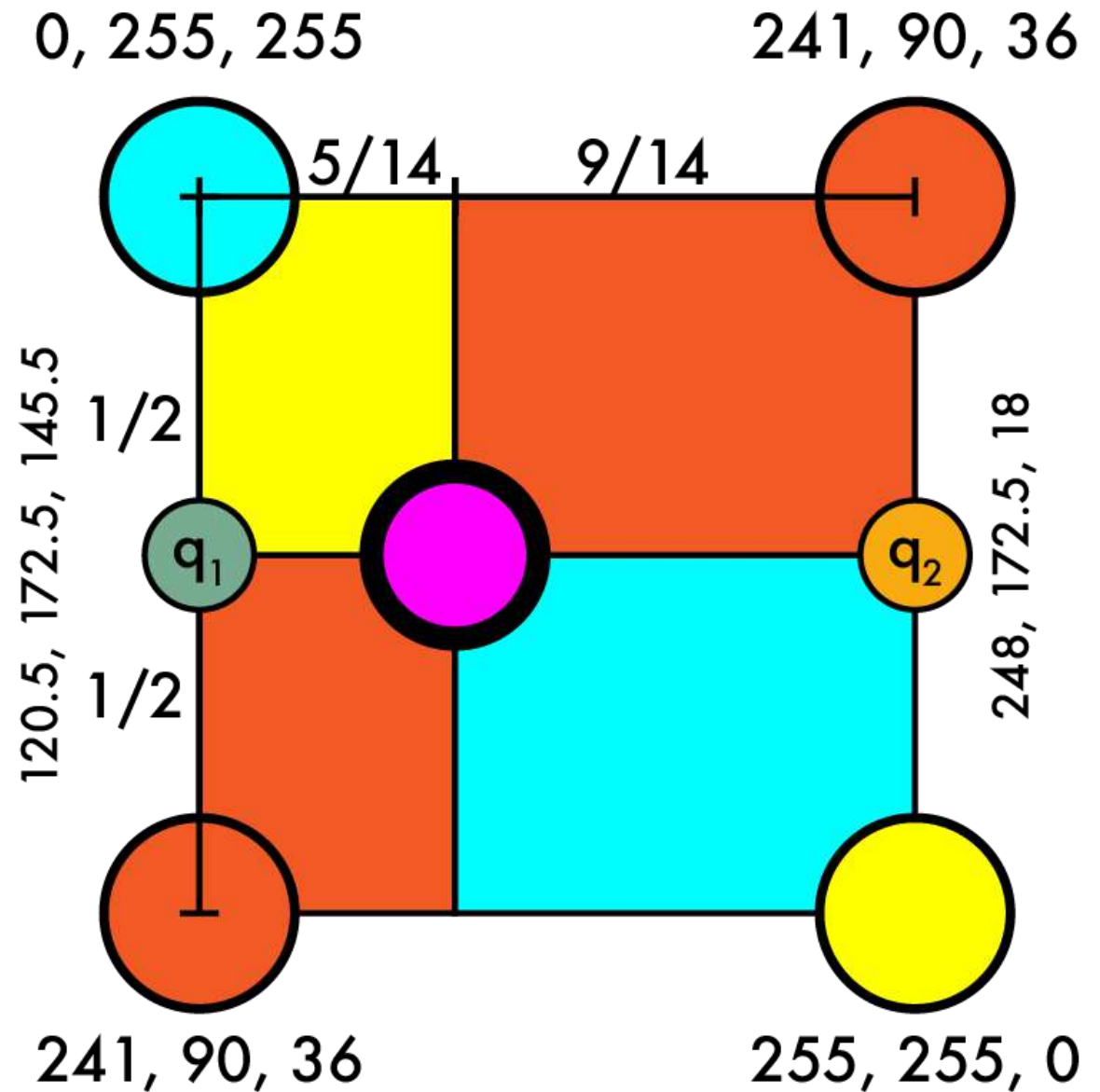
## Resize 4x4 -> 7x7

- Create our new image
- Match up coordinates
  - $4/7 X - 3/14 = Y$
- Iterate over new pts
  - Map to old coords
  - $(1, 3) \rightarrow (5/14, 21/14)$
  - Interpolate old values
    - $q_1 = (120.5, 172.5, 145.5)$
    - $q_2 = (248, 172.5, 18)$
    - $q = r, g, b$
    - $q = 9/14 * q_1 + 5/14 * q_2$



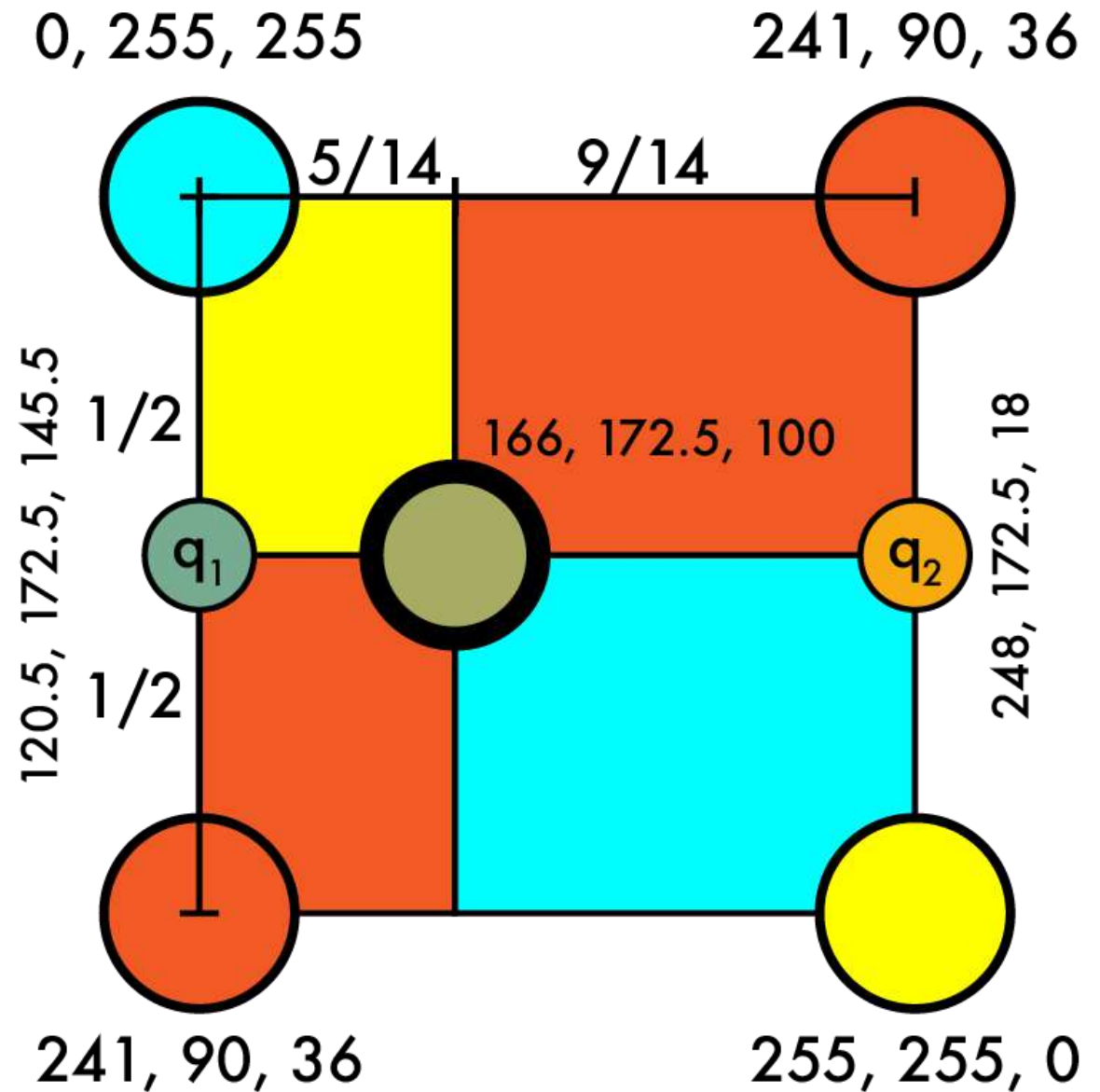
## Resize 4x4 -> 7x7

- Create our new image
- Match up coordinates
  - $4/7 X - 3/14 = Y$
- Iterate over new pts
  - Map to old coords
  - $(1, 3) \rightarrow (5/14, 21/14)$
  - Interpolate old values
    - $q_1 = (120.5, 172.5, 145.5)$
    - $q_2 = (248, 172.5, 18)$
    - $q = r, g, b$
    - $q = 9/14 * q_1 + 5/14 * q_2$
    - $r = 9/14 * 120.5 + 5/14 * 248$
    - $g = 9/14 * 172.5 + 5/14 * 172.5$
    - $b = 9/14 * 145.5 + 5/14 * 18$



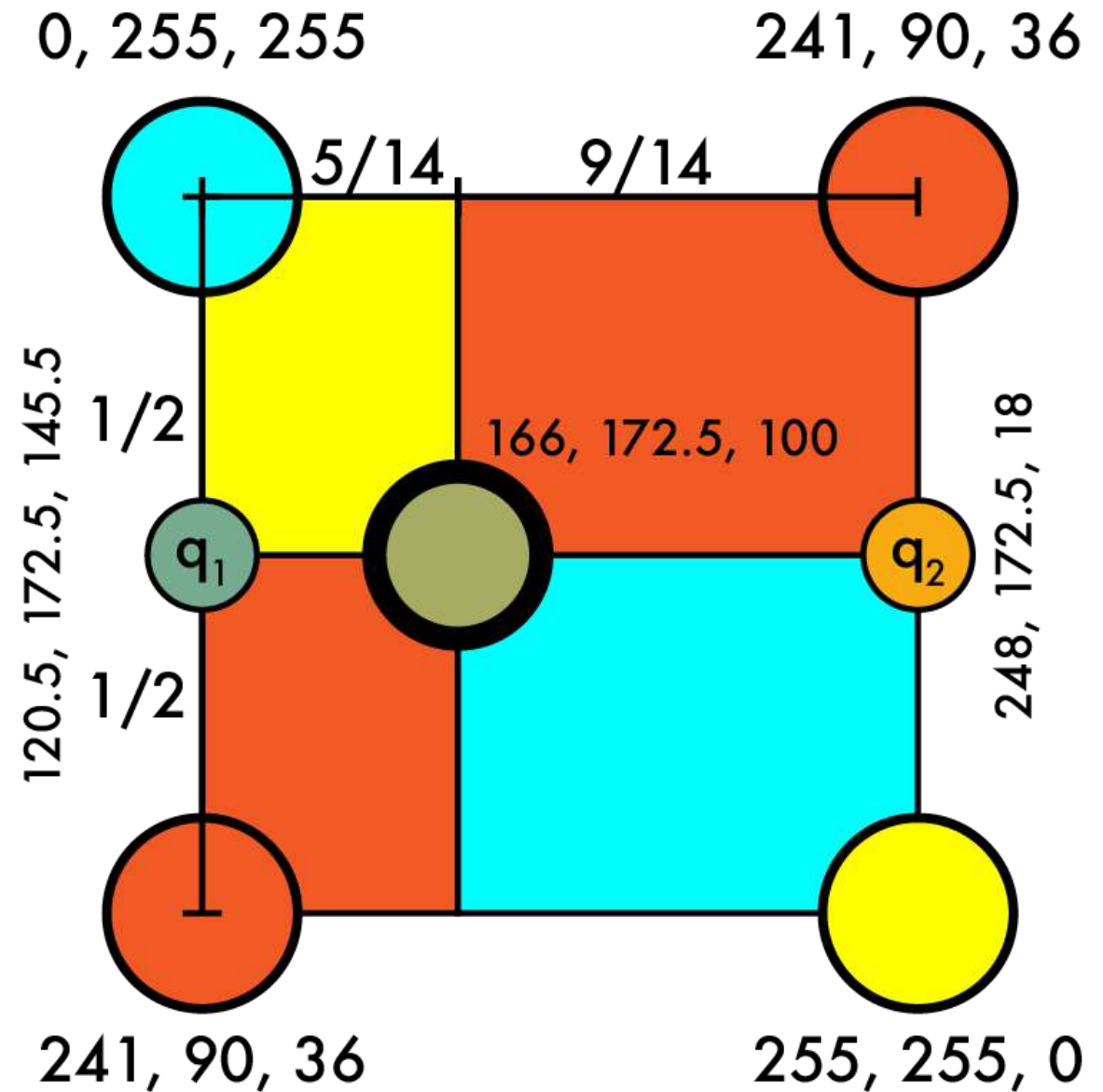
## Resize 4x4 -> 7x7

- Create our new image
- Match up coordinates
  - $4/7 X - 3/14 = Y$
- Iterate over new pts
  - Map to old coords
  - $(1, 3) \rightarrow (5/14, 21/14)$
  - Interpolate old values
    - $q_1 = (120.5, 172.5, 145.5)$
    - $q_2 = (248, 172.5, 18)$
    - $q = (166, 172.5, 100)$



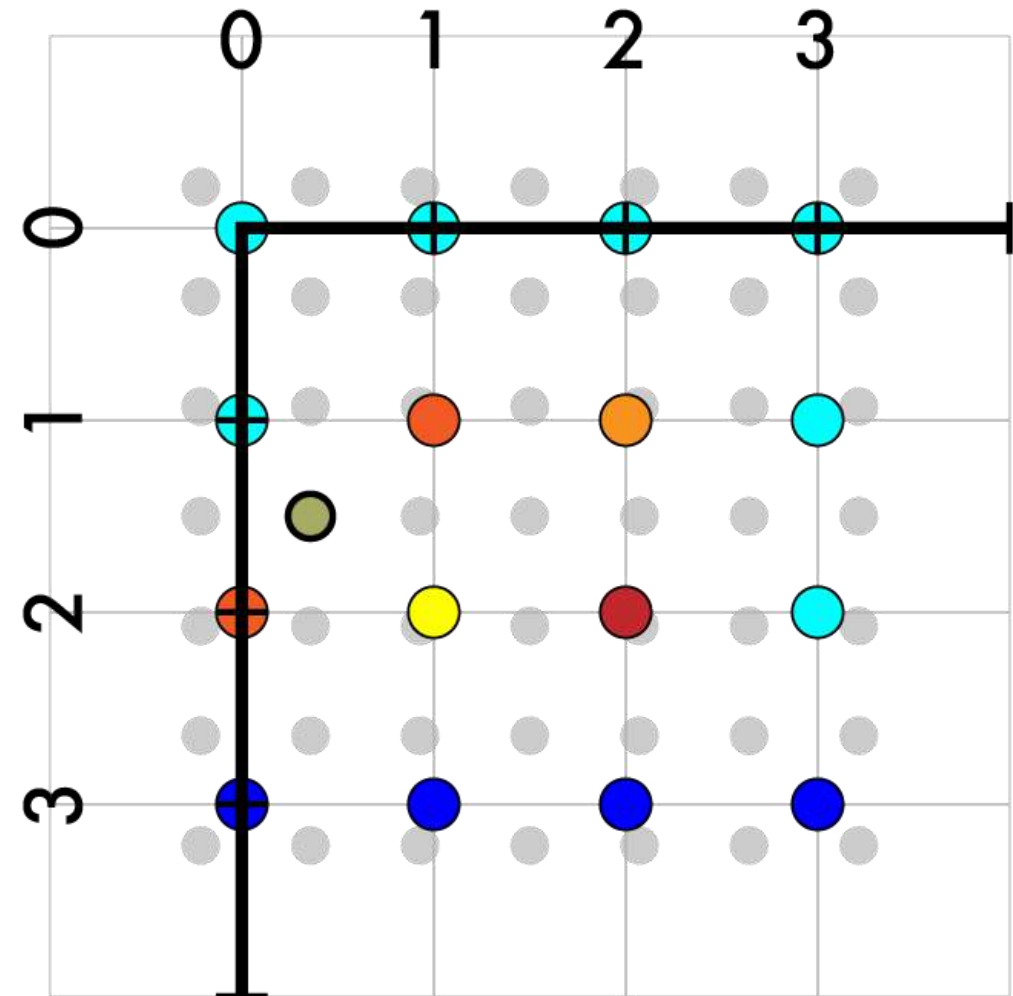
## Resize 4x4 -> 7x7

- Create our new image
- Match up coordinates
  - $4/7 X - 3/14 = Y$
- Iterate over new pts
  - Map to old coords
  - $(1, 3) \rightarrow (5/14, 21/14)$
  - Interpolate old values
    - $q = (166, 172.5, 100)$



## Resize 4x4 -> 7x7

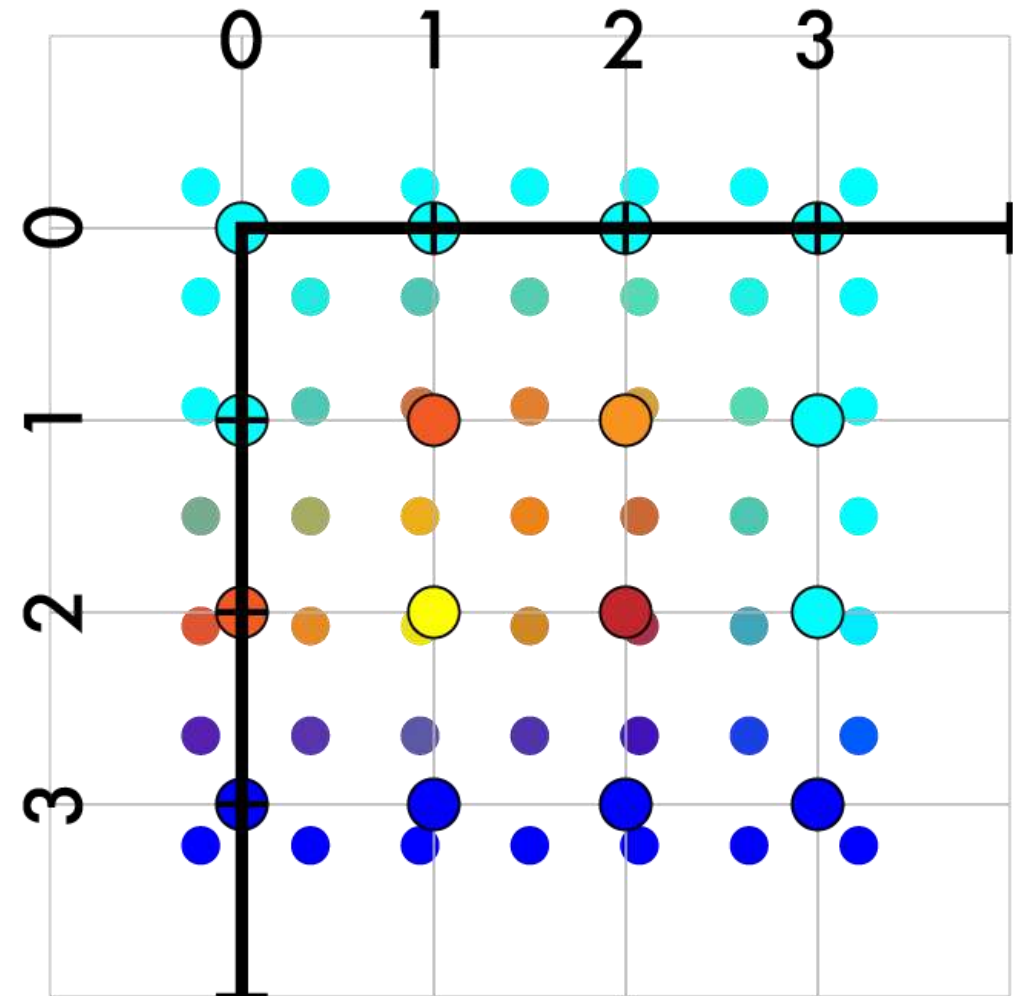
- Create our new image
- Match up coordinates
  - $4/7 X - 3/14 = Y$
- Iterate over new pts
  - Map to old coords
  - $(1, 3) \rightarrow (5/14, 21/14)$
  - Interpolate old values
    - $q = (166, 172.5, 100)$





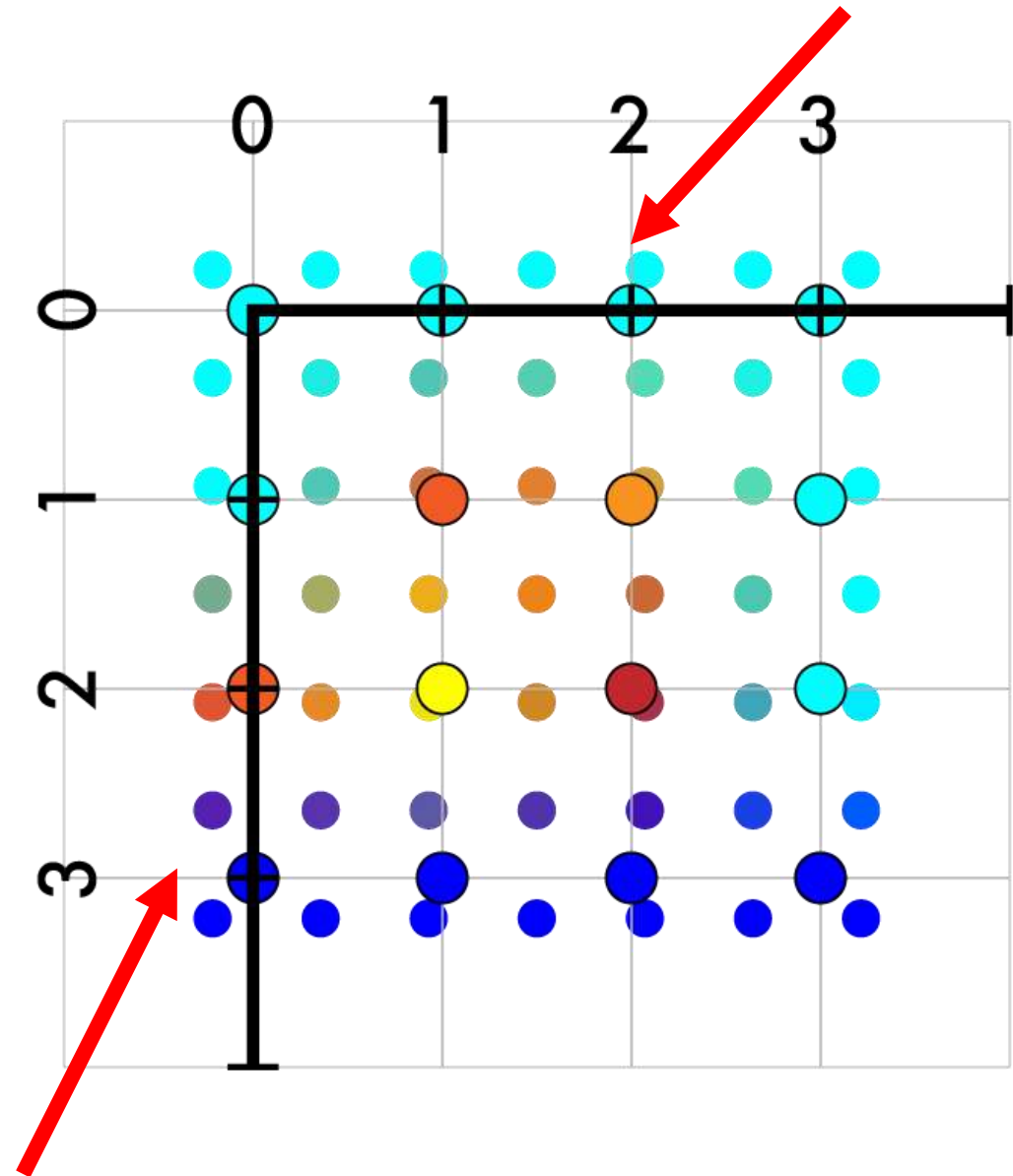
## Resize 4x4 -> 7x7

- Create our new image
- Match up coordinates
  - $\frac{4}{7} X - \frac{3}{14} = Y$
- Iterate over new pts
  - Map to old coords
  - $(1, 3) \rightarrow (5/14, 21/14)$
  - Interpolate old values
    - $q = (166, 172.5, 100)$
- Fill in the rest



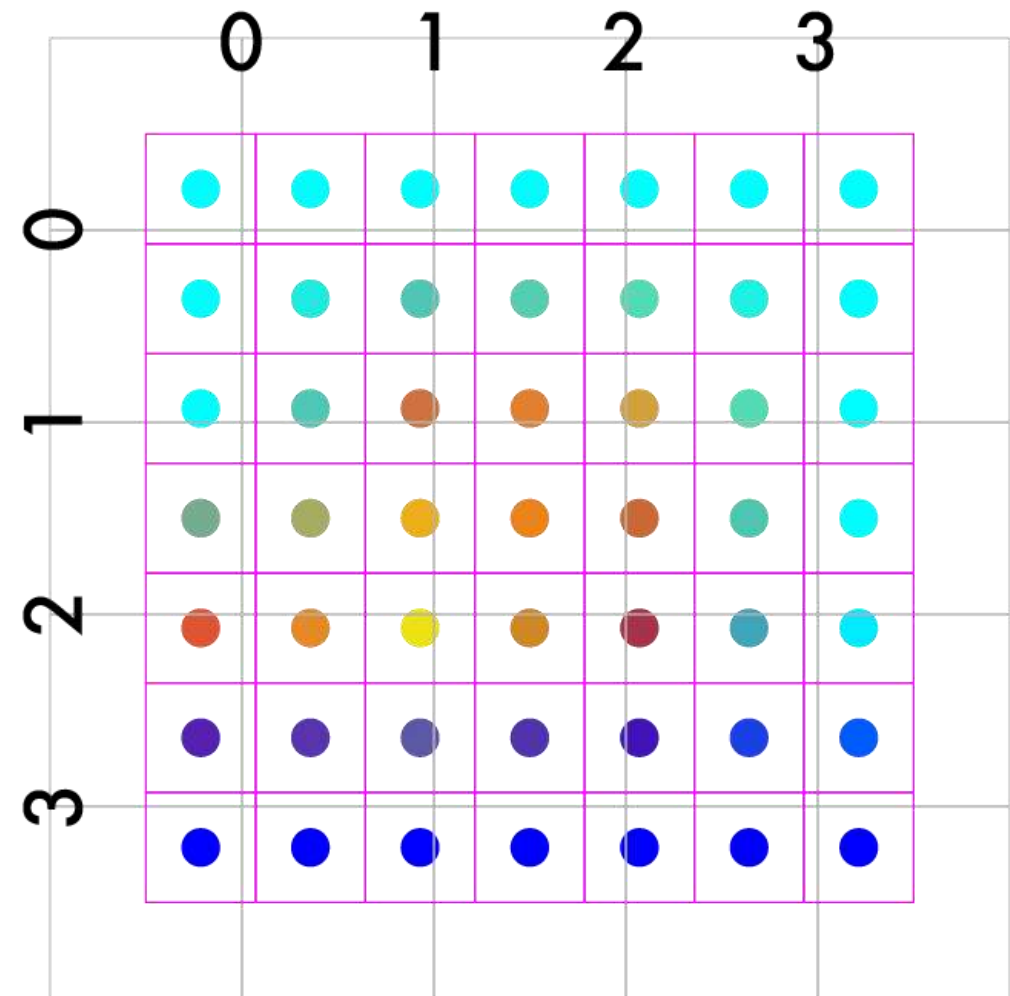
## Resize 4x4 -> 7x7

- Create our new image
- Match up coordinates
  - $4/7 X - 3/14 = Y$
- Iterate over new pts
  - Map to old coords
  - $(1, 3) \rightarrow (5/14, 21/14)$
  - Interpolate old values
    - $q = (166, 172.5, 100)$
- Fill in the rest
  - On outer edges use padding!



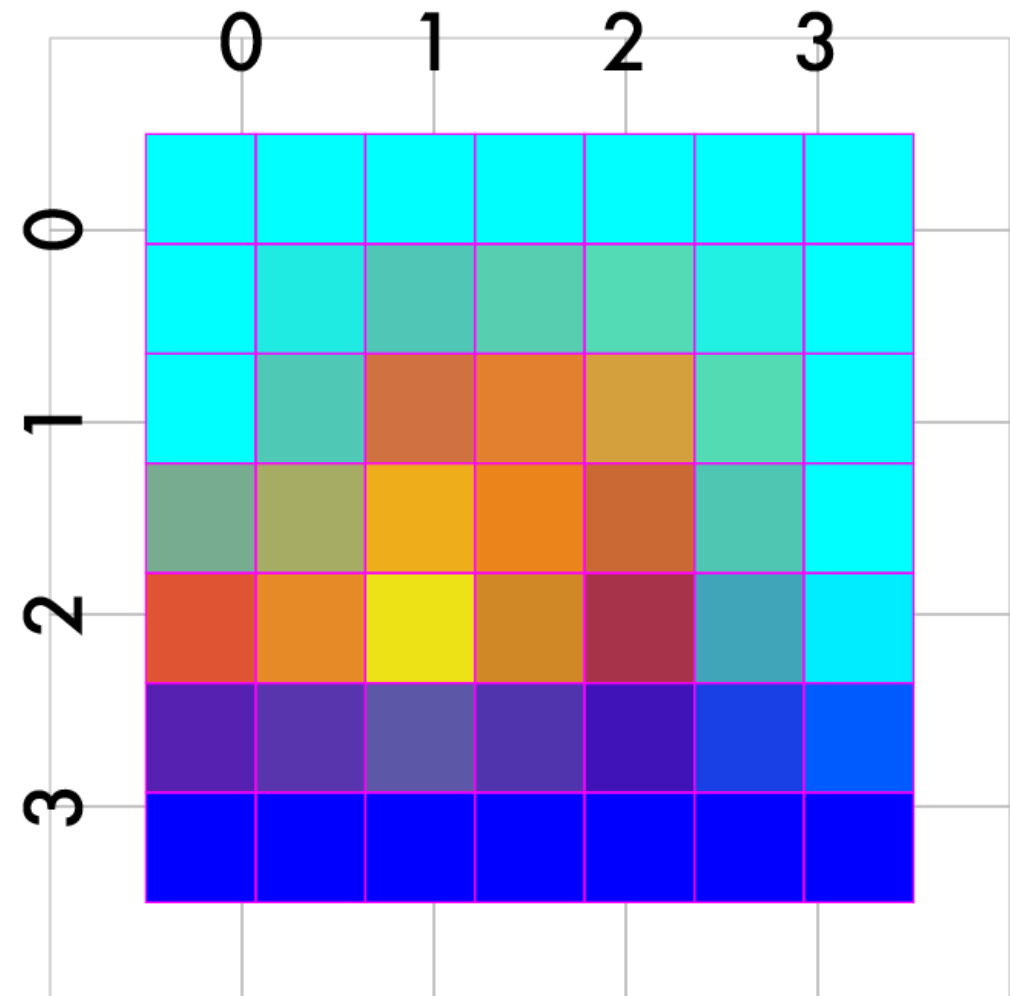
## Resize 4x4 -> 7x7

- Create our new image
- Match up coordinates
  - $4/7 X - 3/14 = Y$
- Iterate over new pts
  - Map to old coords
  - $(1, 3) \rightarrow (5/14, 21/14)$
  - Interpolate old values
    - $q = (166, 172.5, 100)$
- Fill in the rest

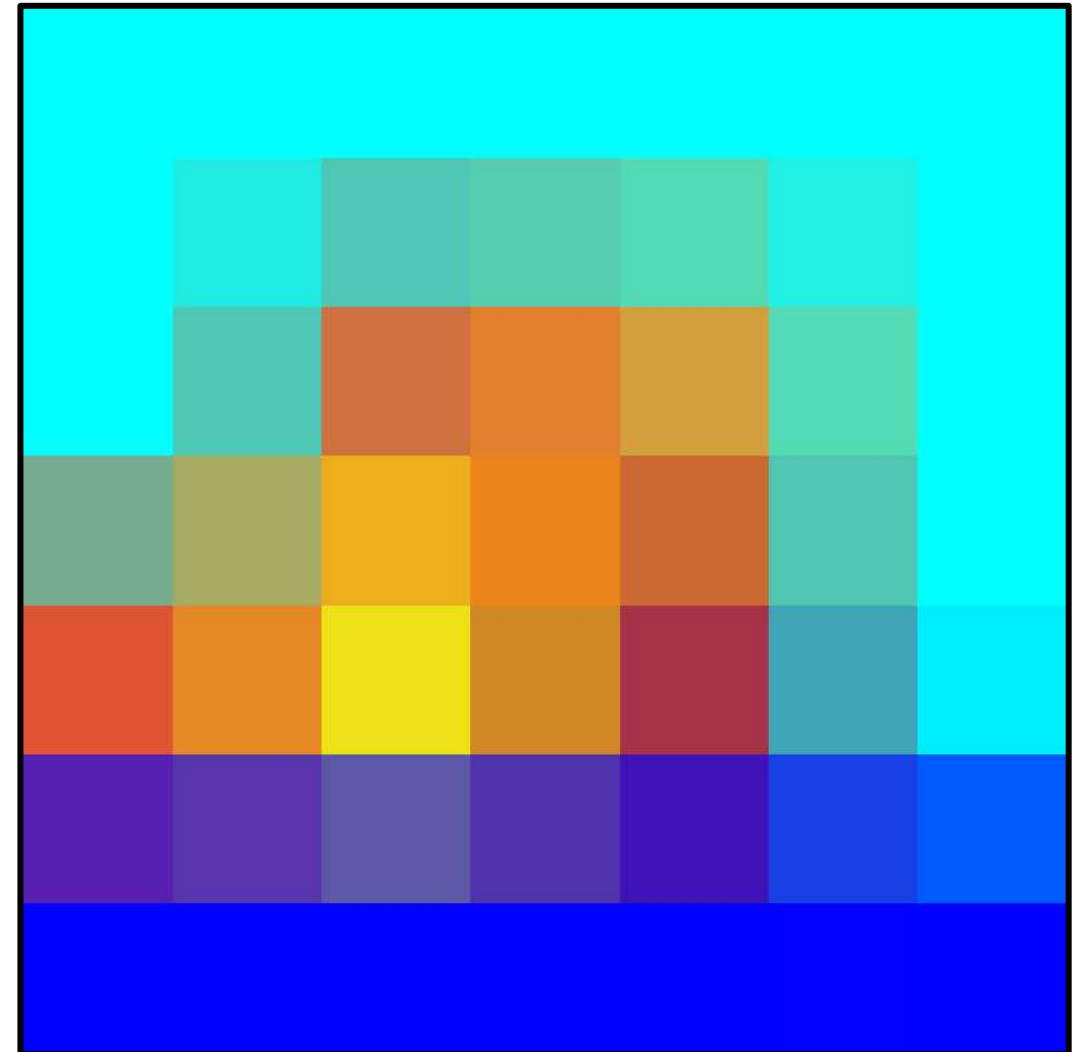
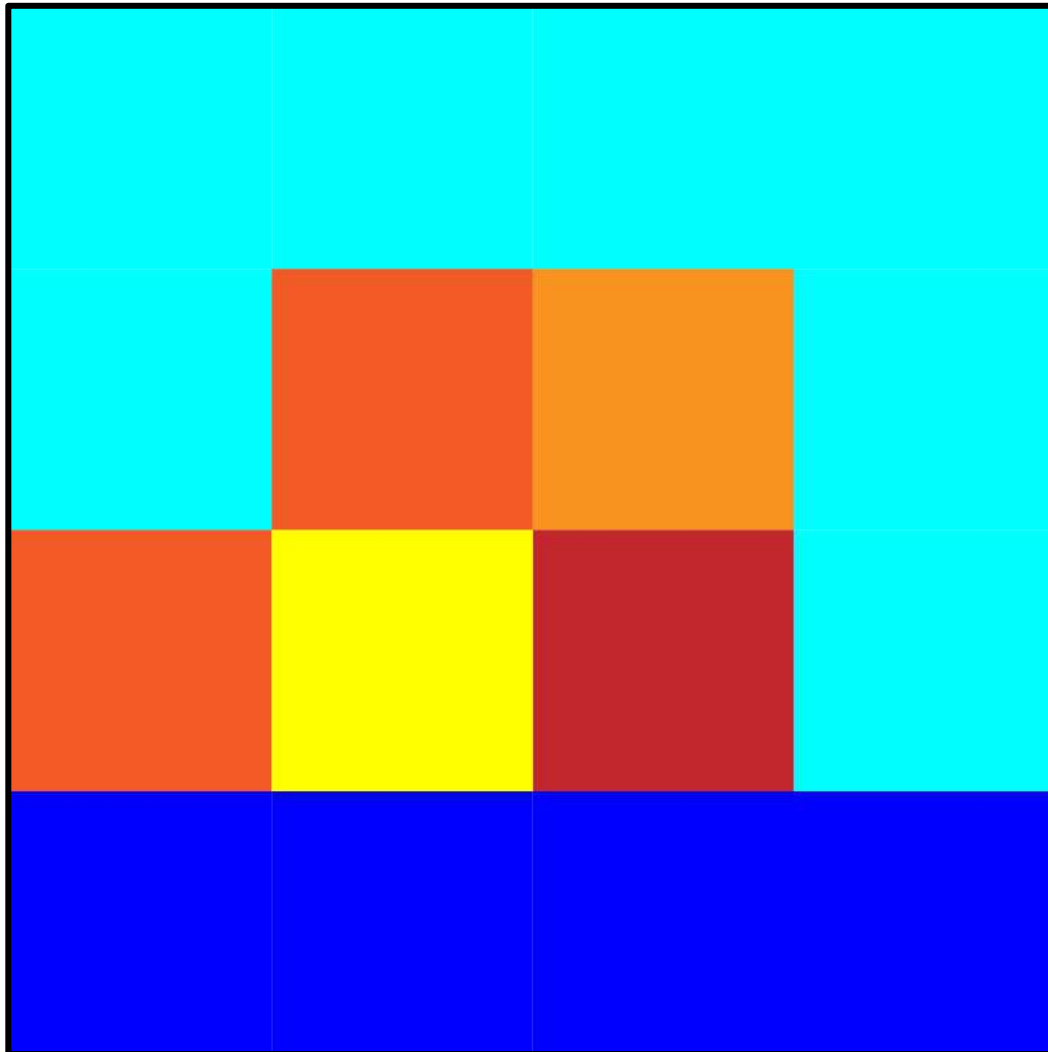


## Resize 4x4 -> 7x7

- Create our new image
- Match up coordinates
  - $4/7 X - 3/14 = Y$
- Iterate over new pts
  - Map to old coords
  - (1, 3) -> (5/14, 21/14)
  - Interpolate old values
    - $q = (166, 172.5, 100)$
- Fill in the rest

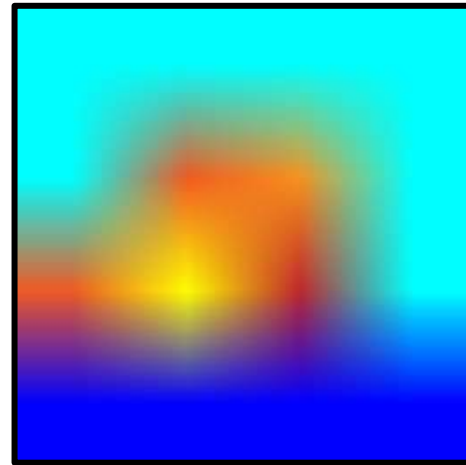
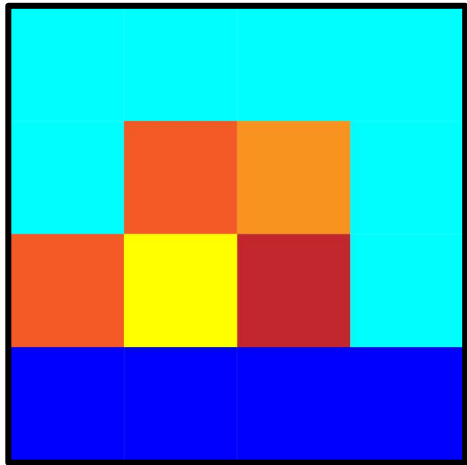


## We did it!

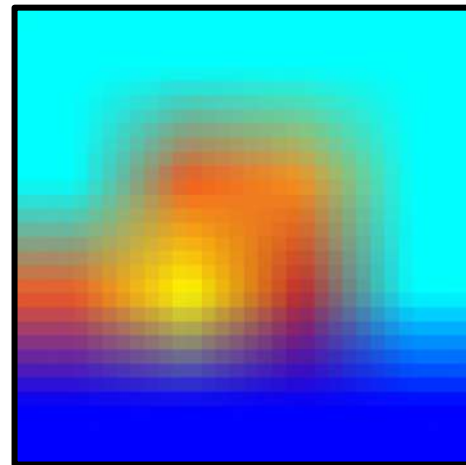
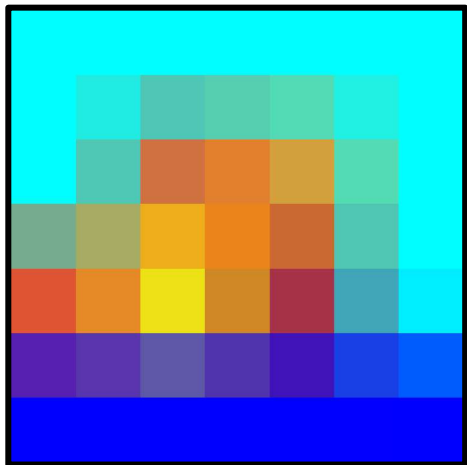




## Different scales



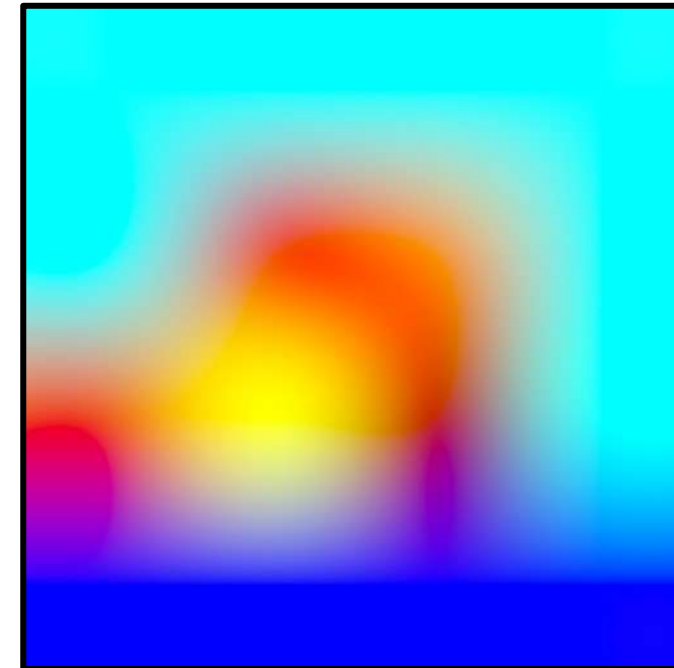
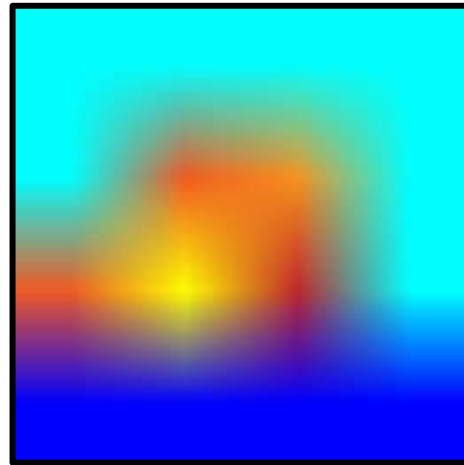
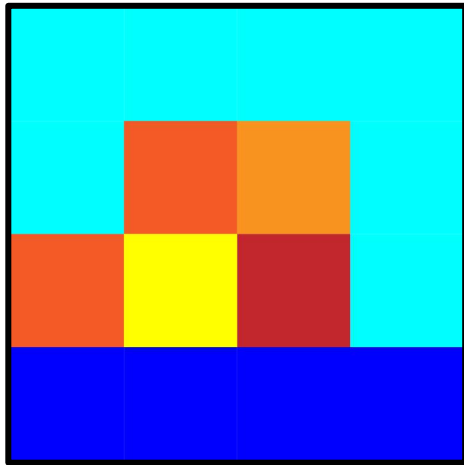
256x256



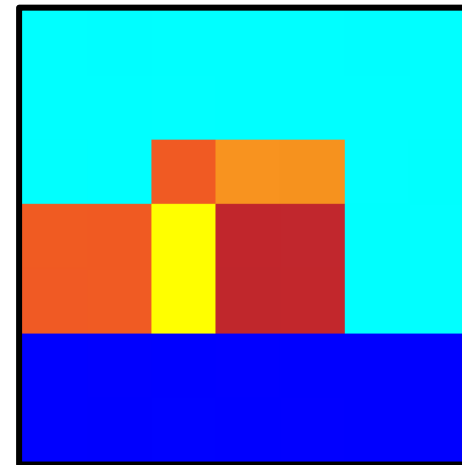
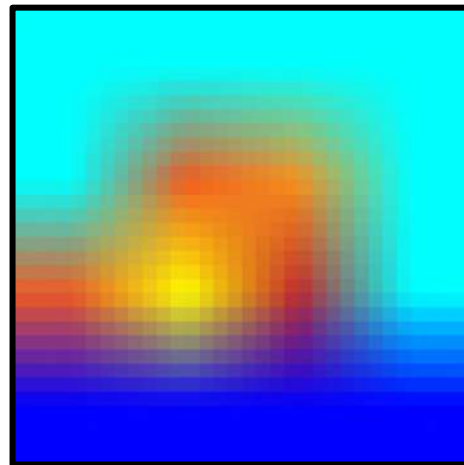
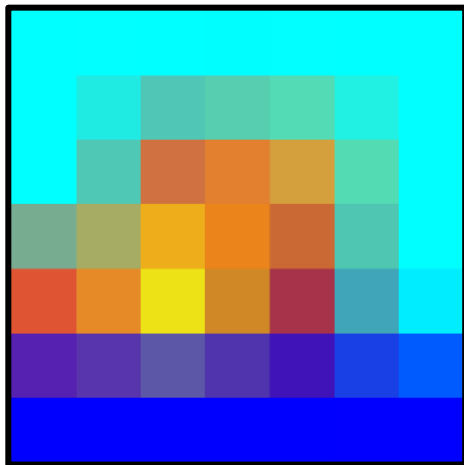
7x7

32x32

## Different methods



Bicubic



7x7  
NN

# Today's Agenda

- Image basics
  - What is an image – addressing pixels
  - Image as a function – image coordinates
- Image interpolation
  - Nearest neighbor
  - Bilinear
  - Bicubic
- Image resizing
  - Enlarge
  - Shrink

## Want to make image smaller



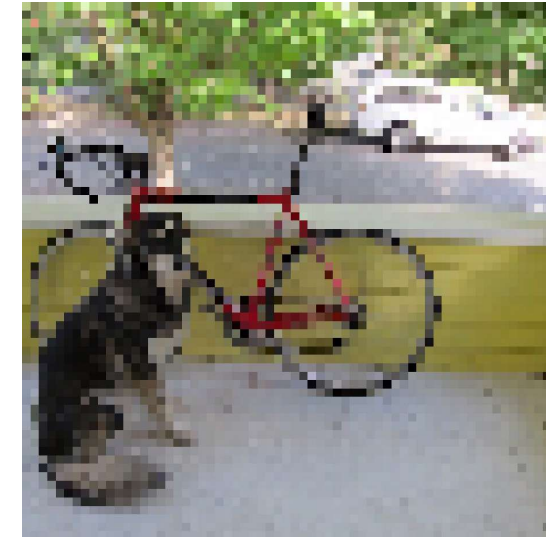


## 448x448 -> 64x64

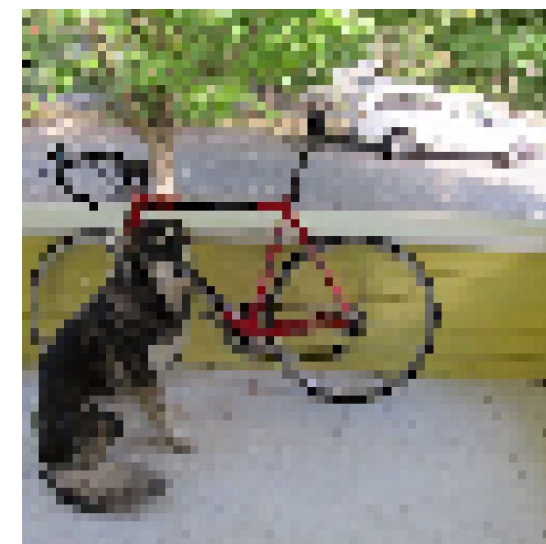




## 448x448 -> 64x64



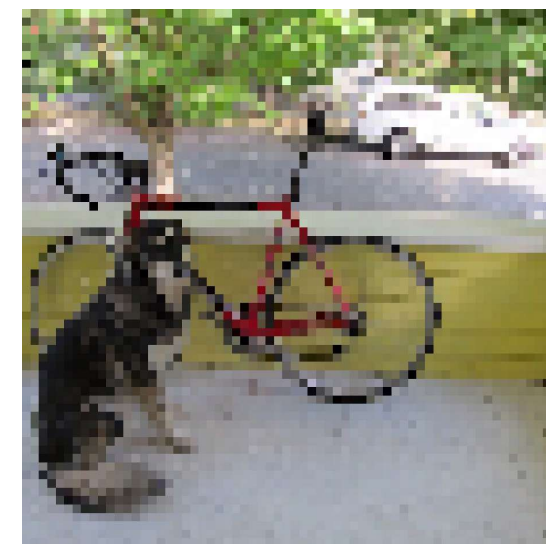
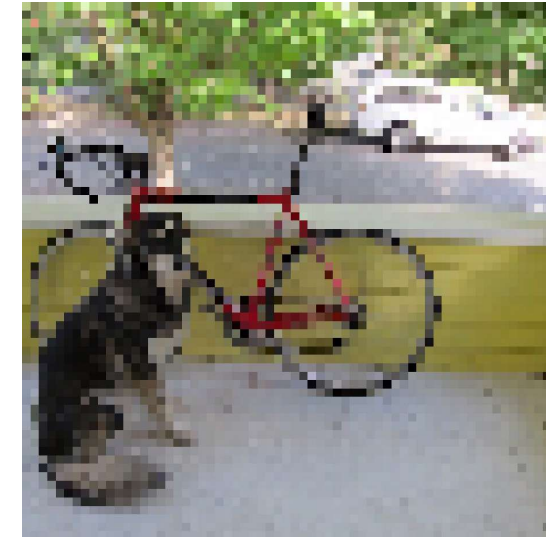
NN



Bilinear



## 448x448 -> 64x64



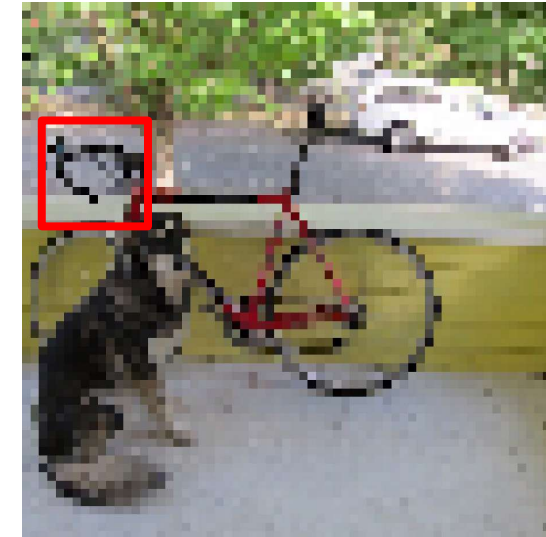


## 448x448 -> 64x64



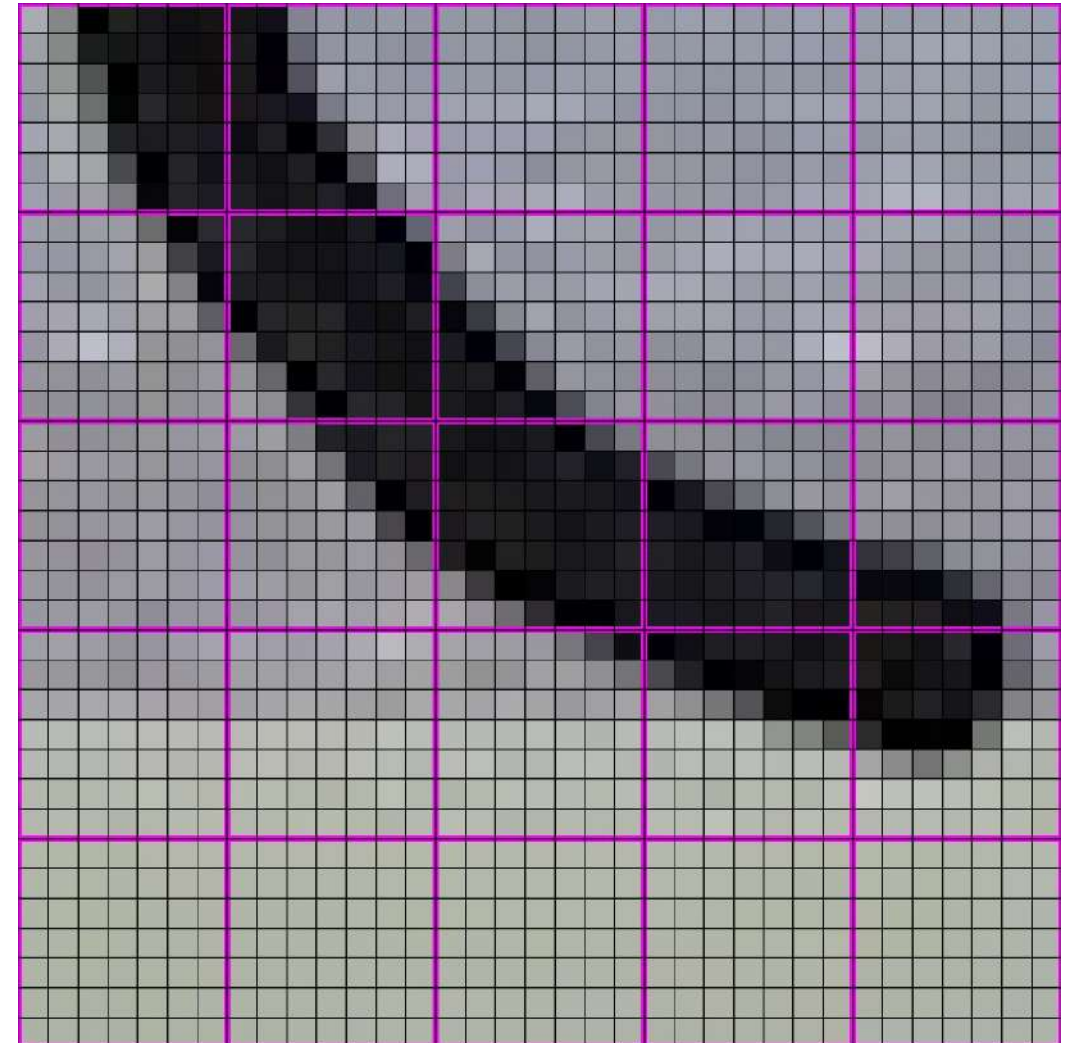


## 448x448 -> 64x64



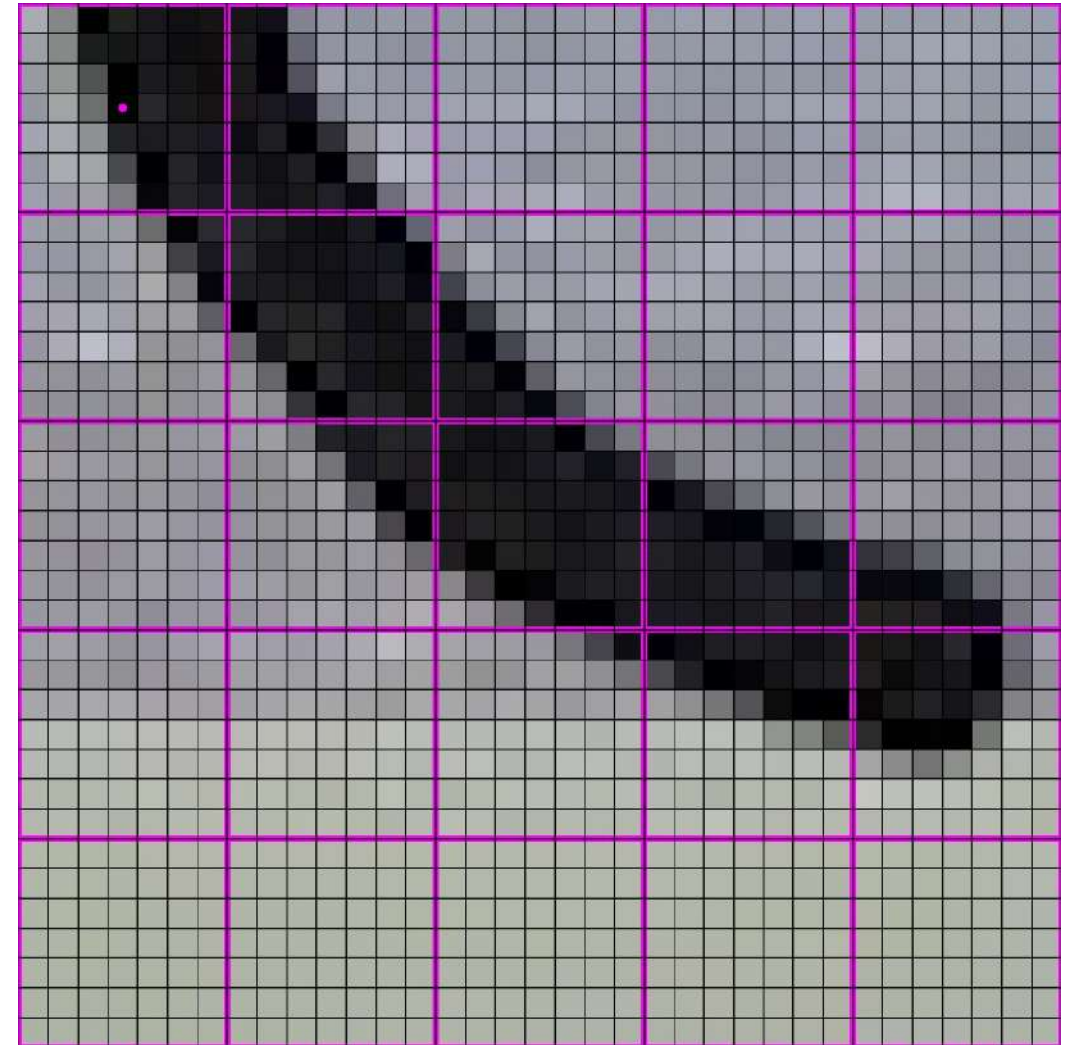


## 448x448 -> 64x64



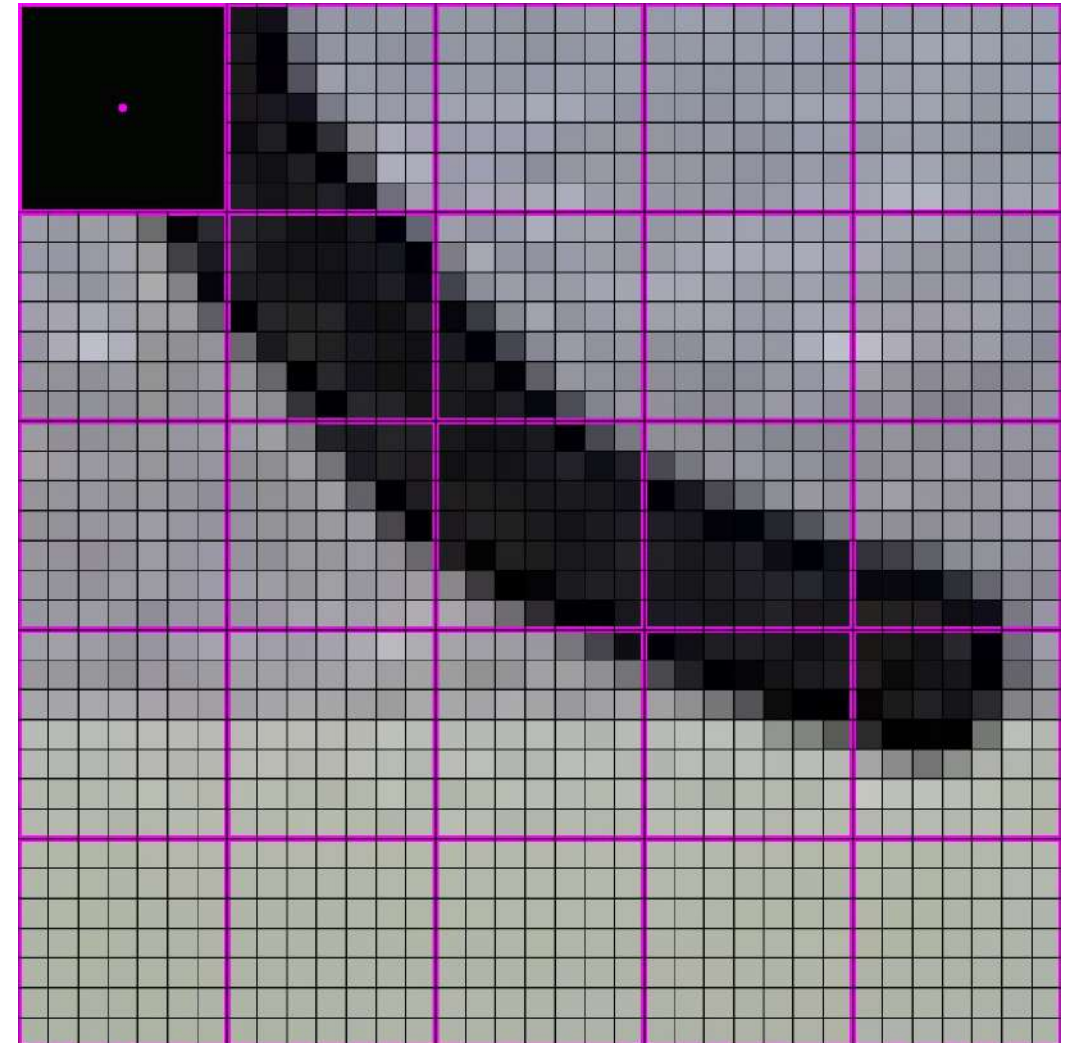


## 448x448 -> 64x64



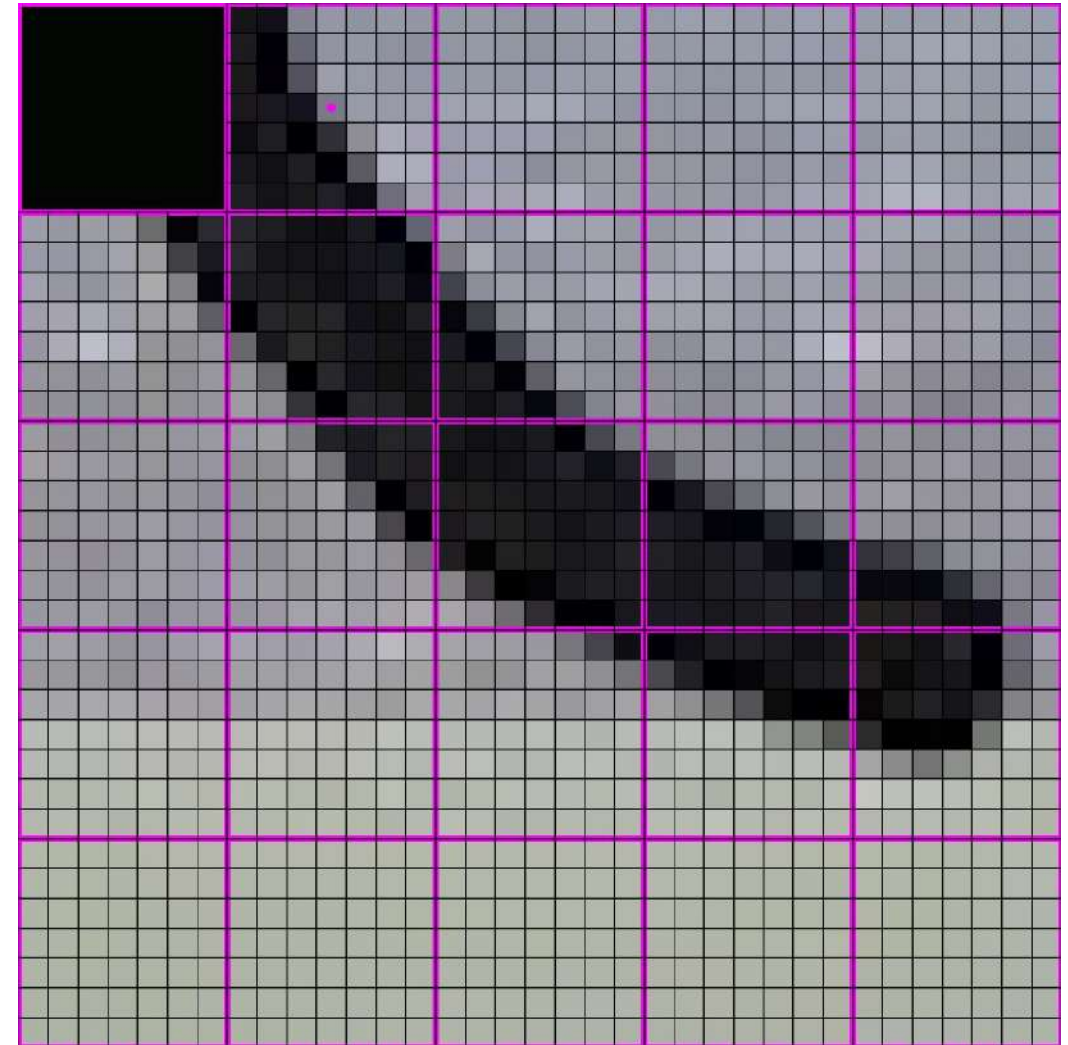


## 448x448 -> 64x64



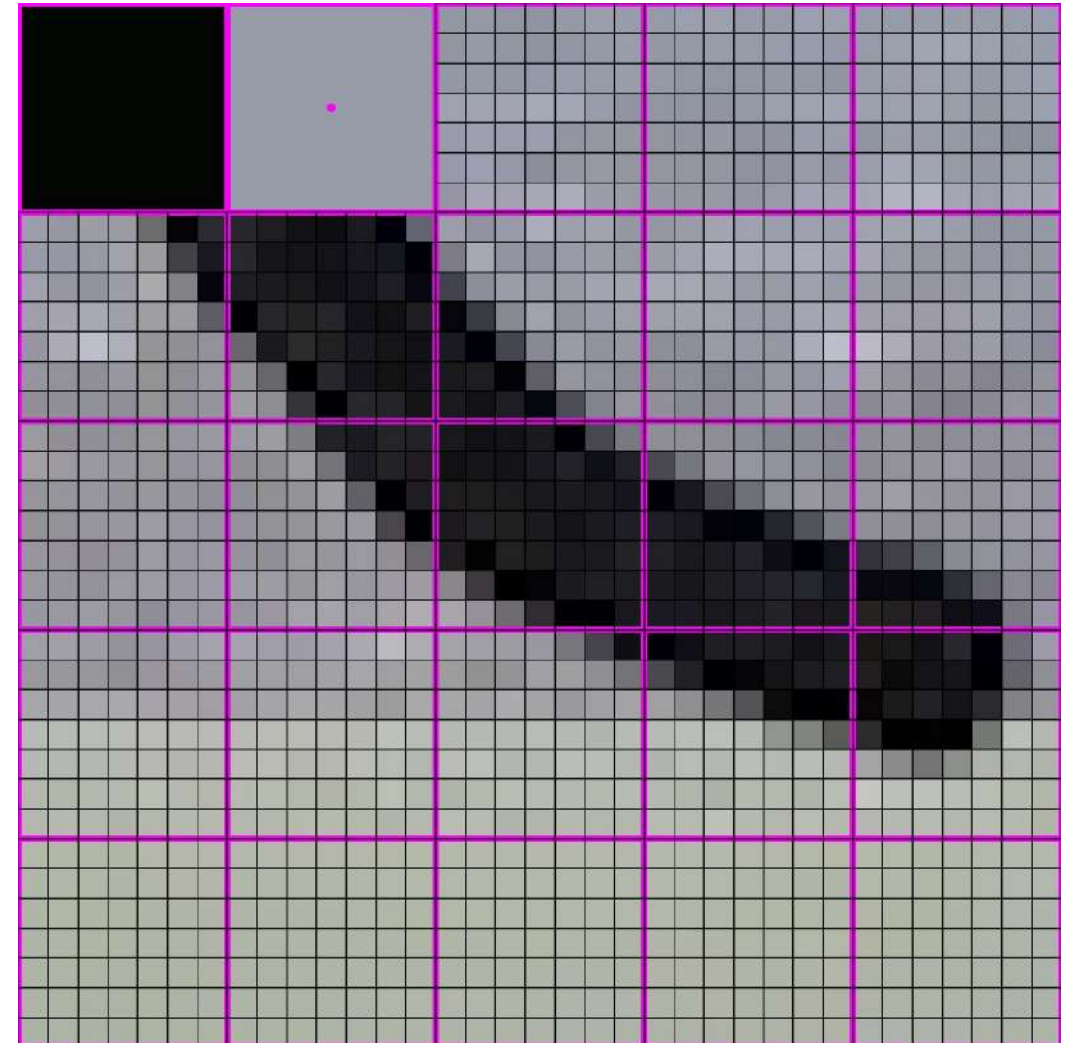


## 448x448 -> 64x64



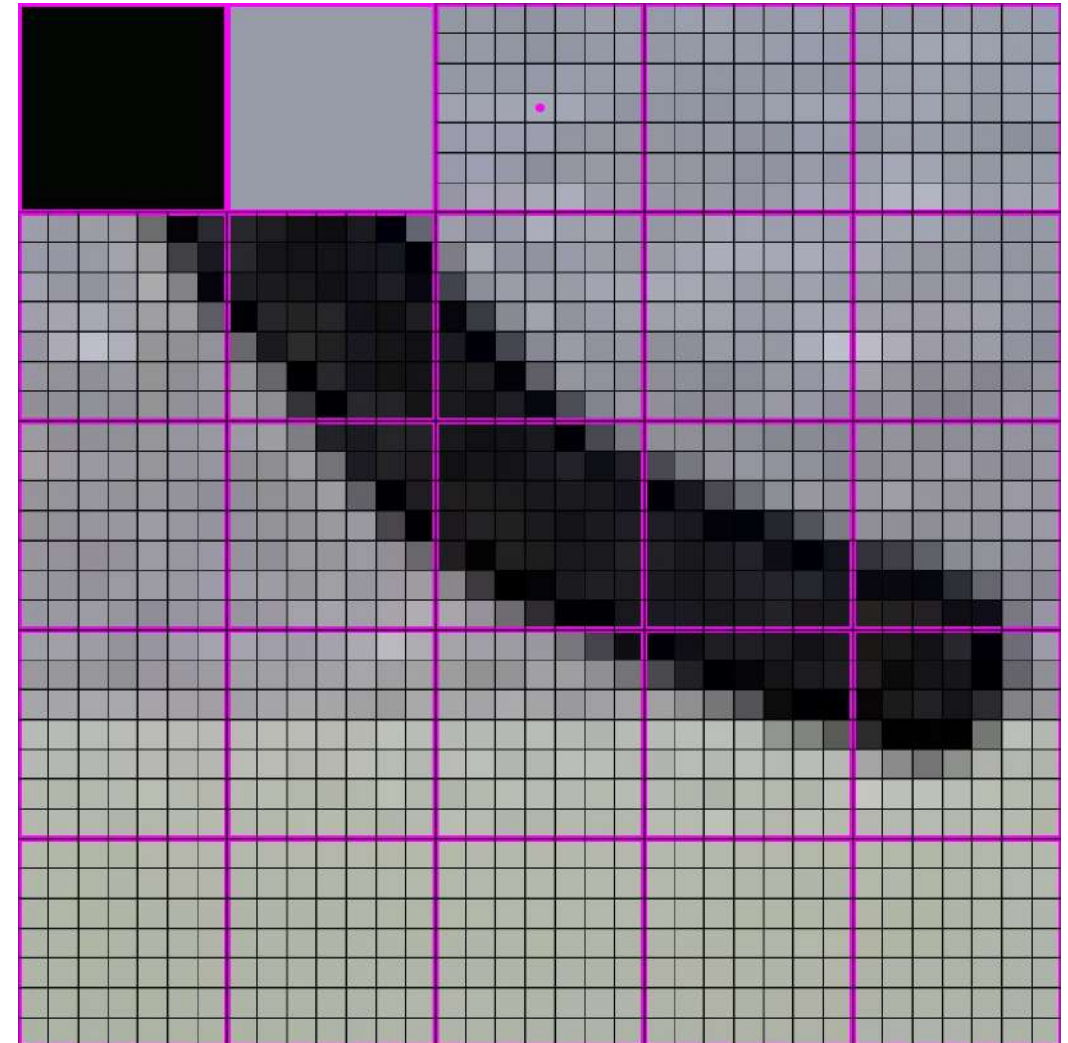


## 448x448 -> 64x64



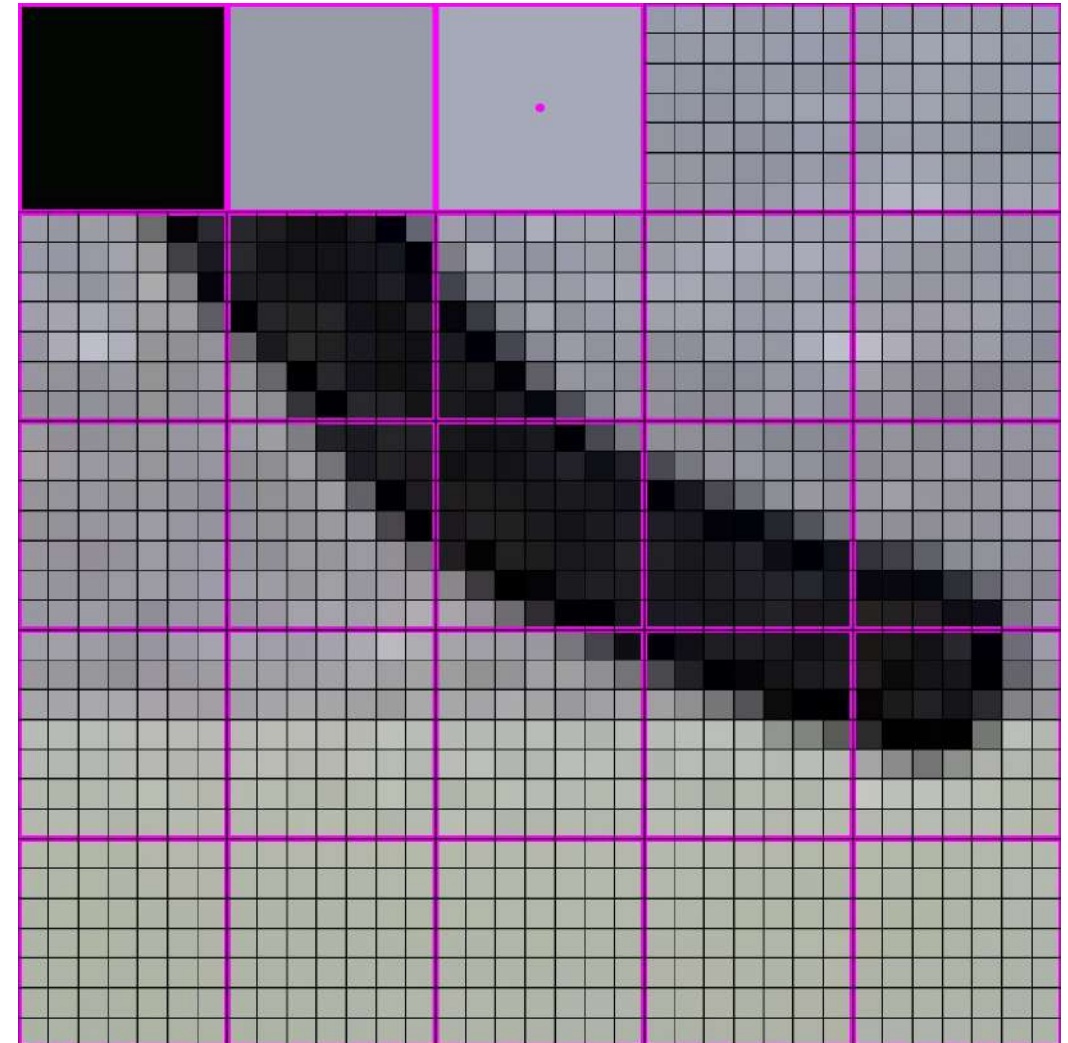


## 448x448 -> 64x64



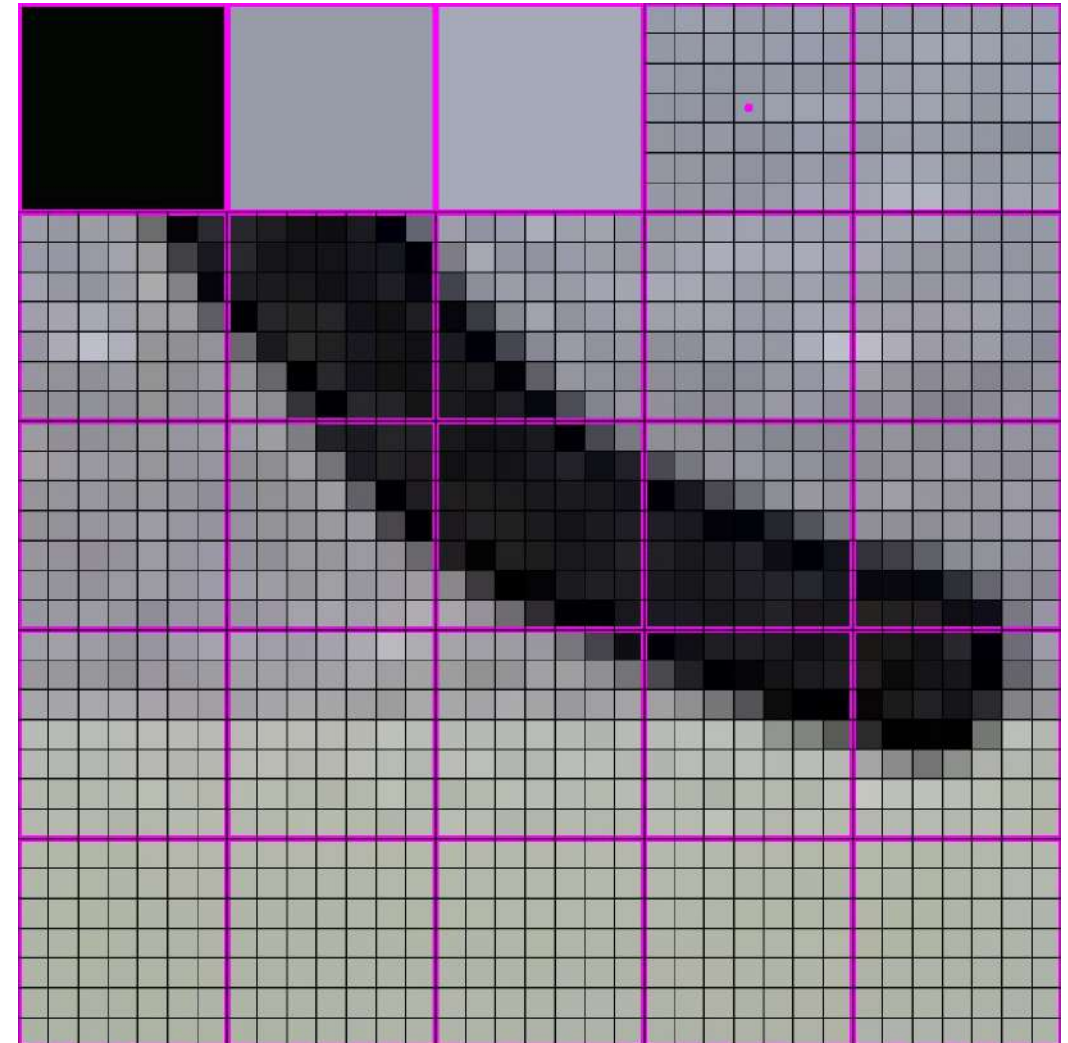


## 448x448 -> 64x64



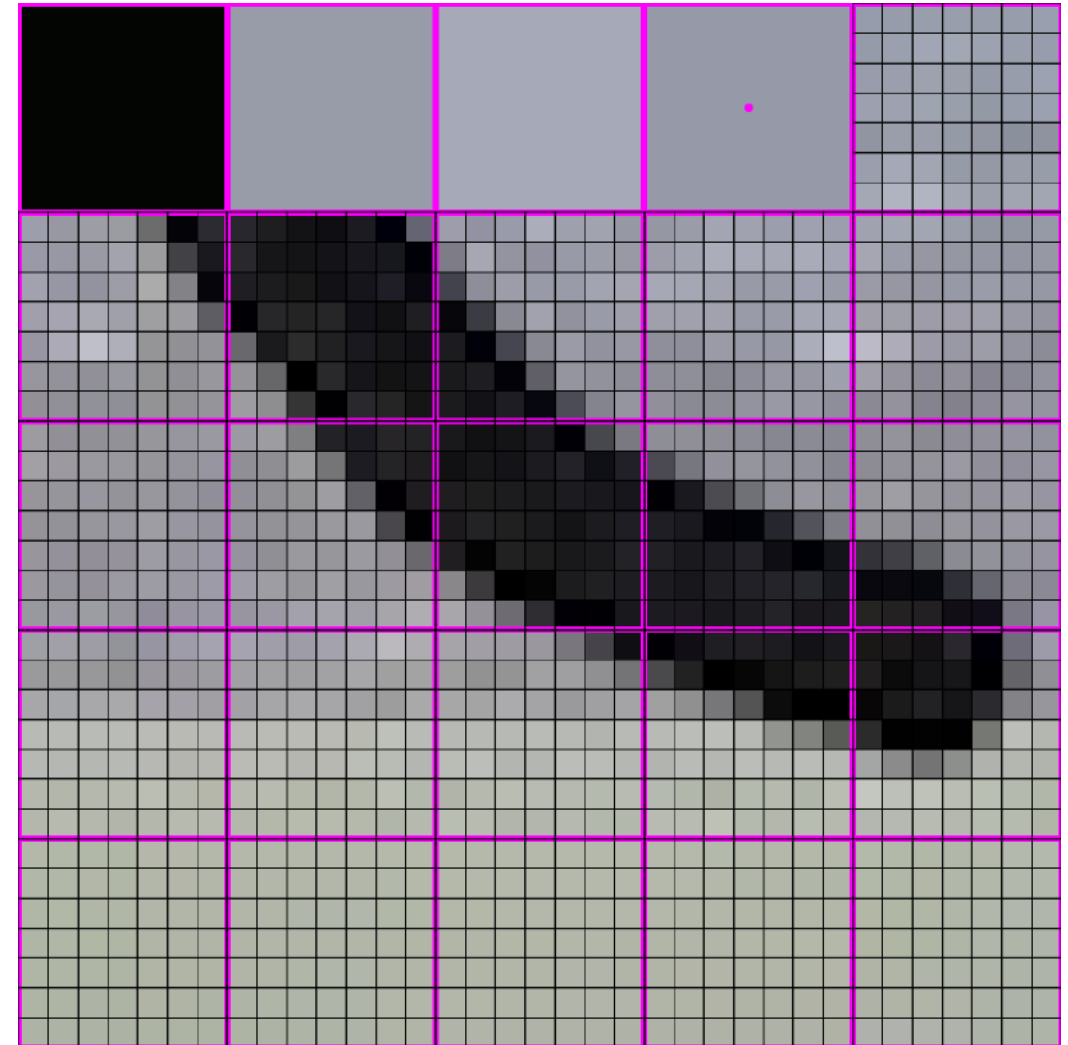


## 448x448 -> 64x64



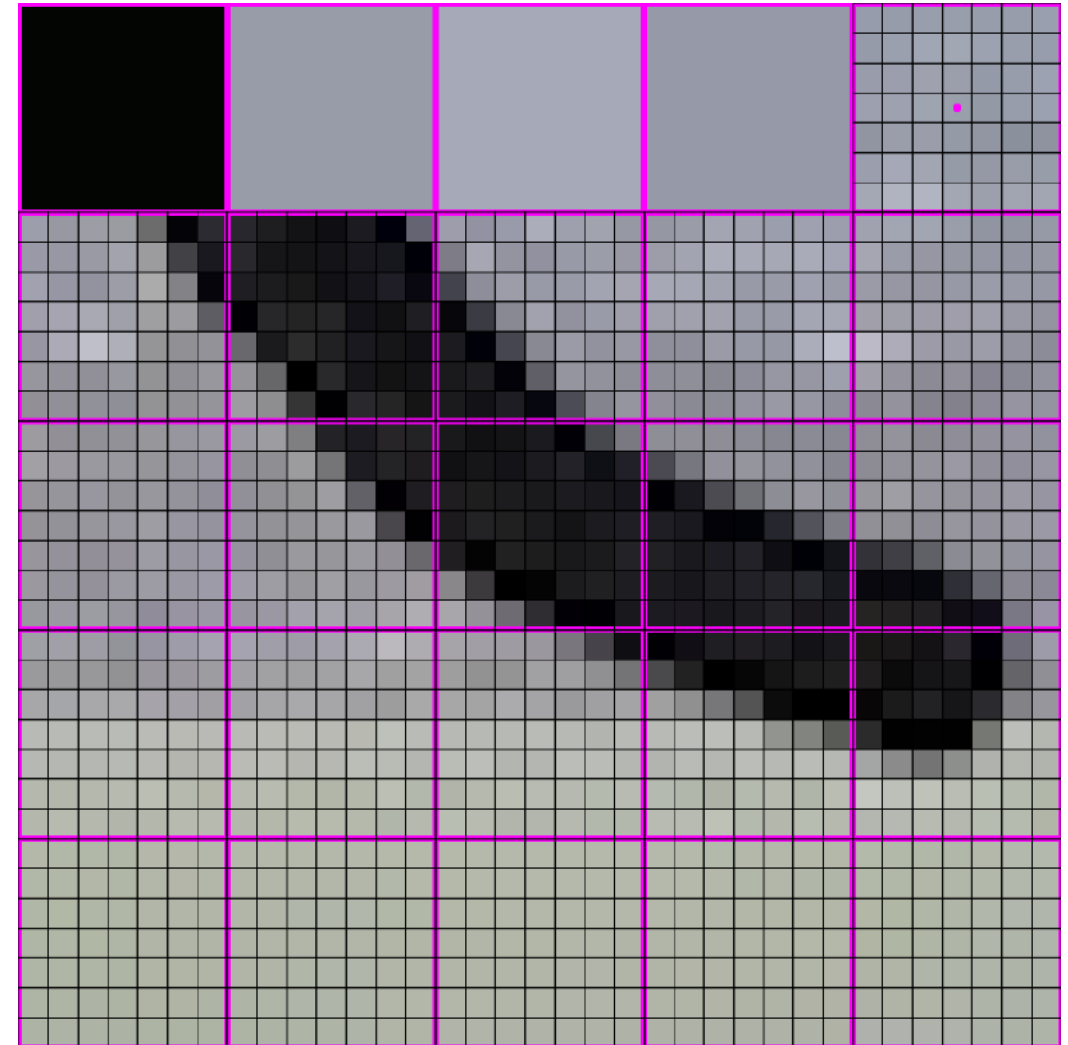


## 448x448 -> 64x64



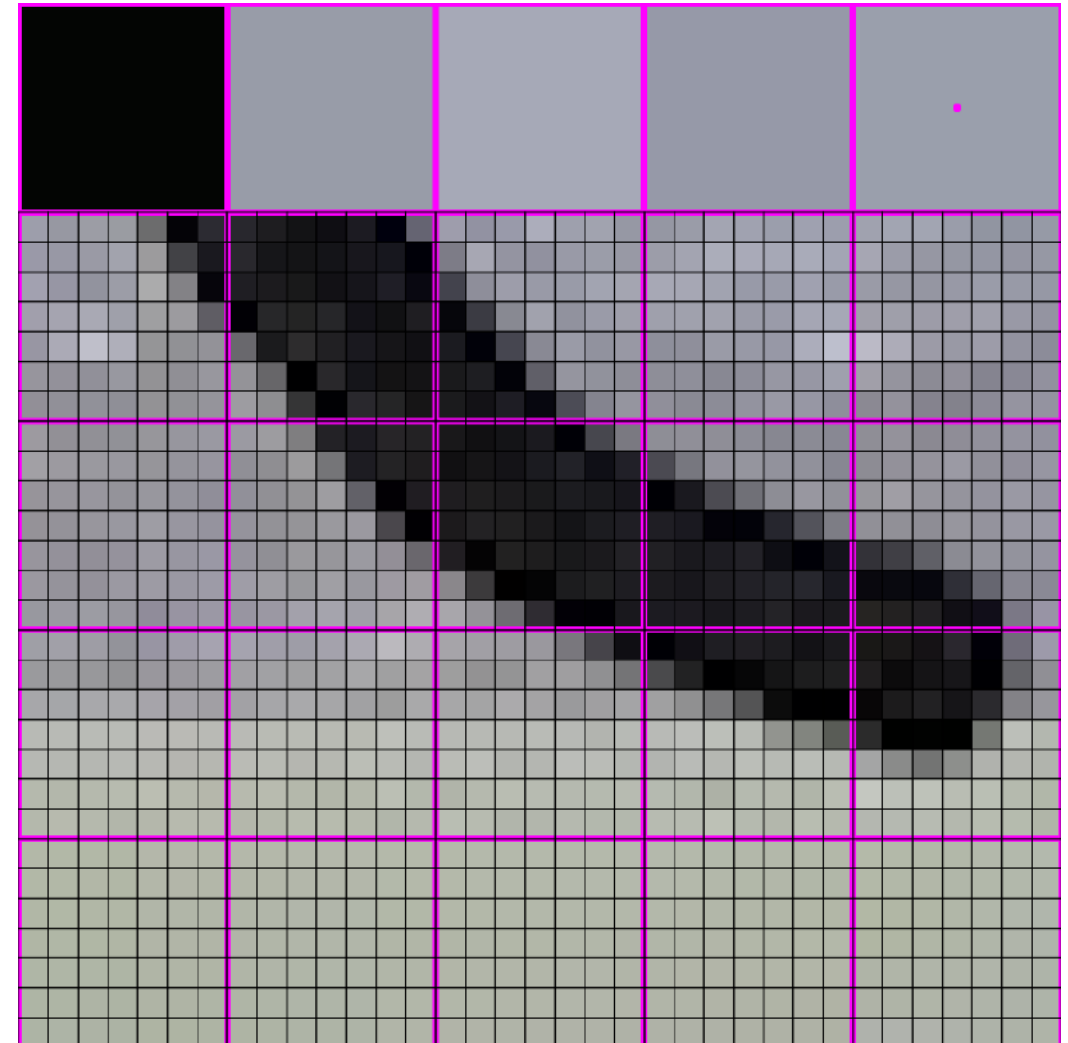


## 448x448 -> 64x64



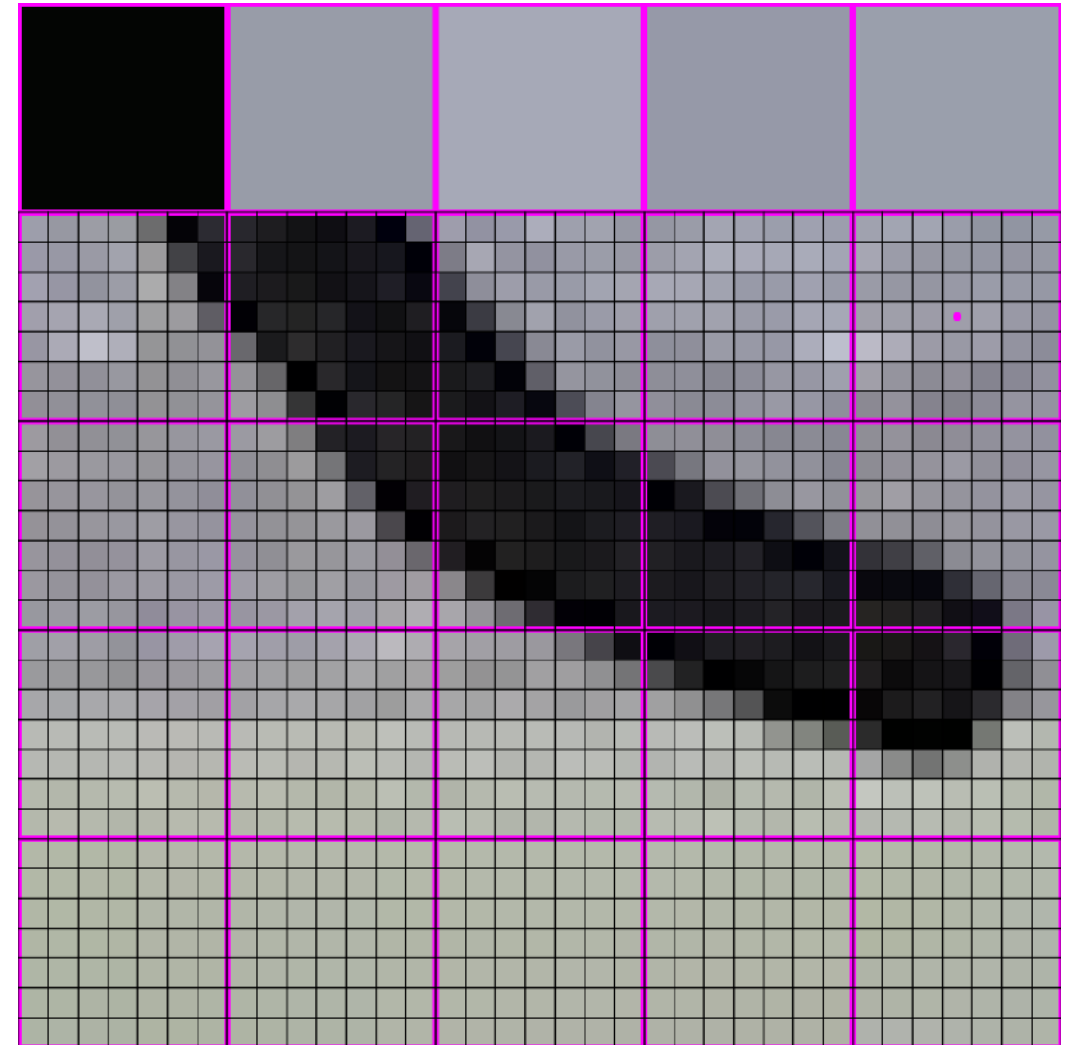


## 448x448 -> 64x64



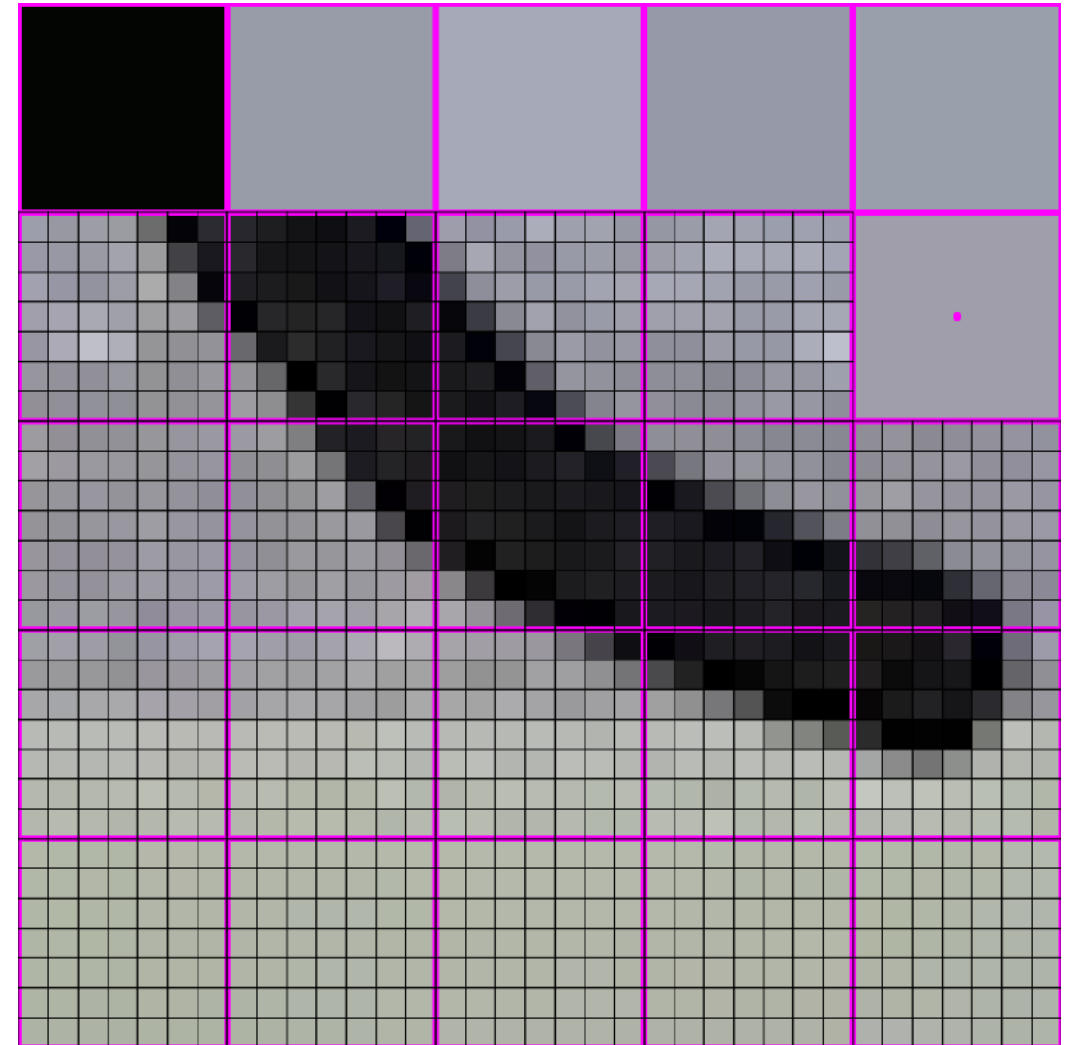


## 448x448 -> 64x64



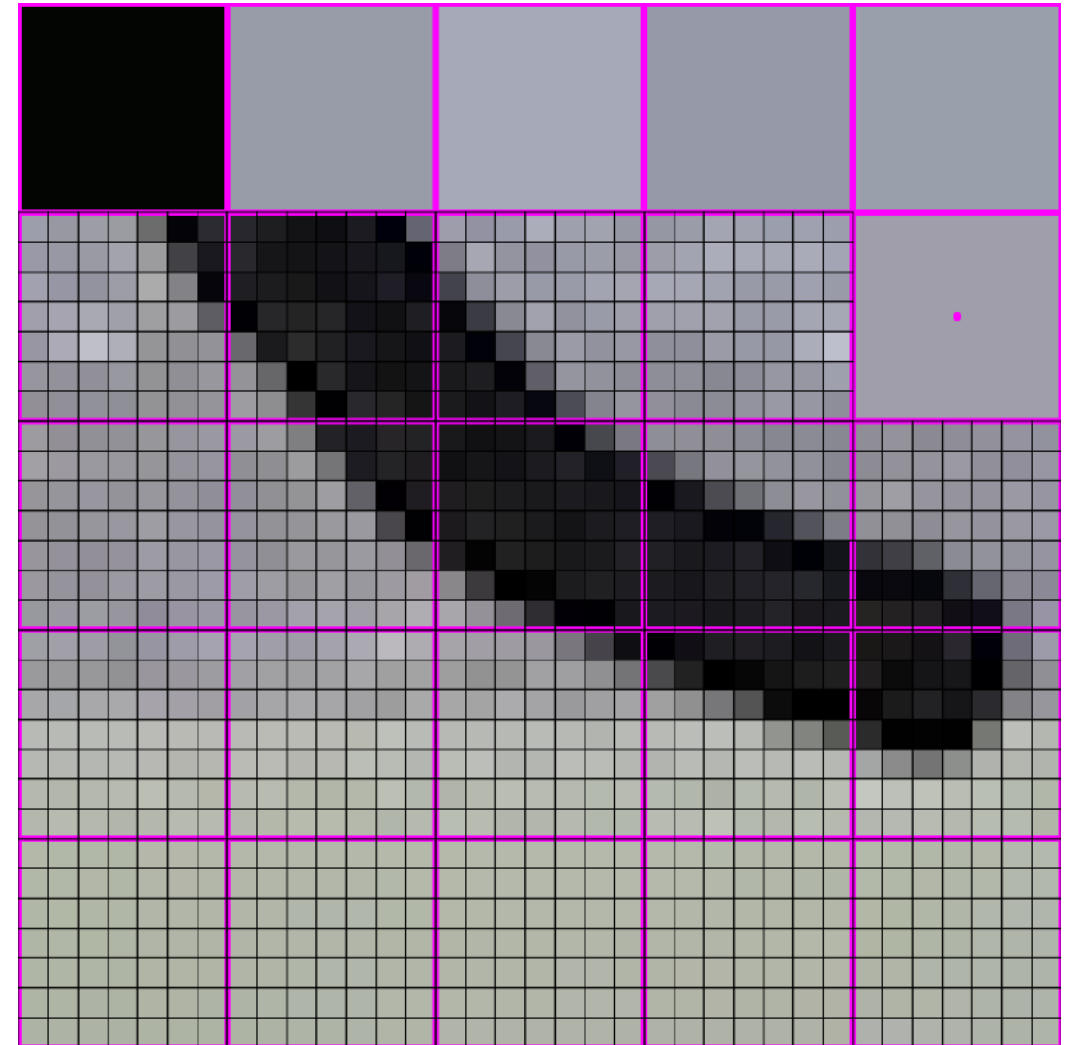


## 448x448 -> 64x64



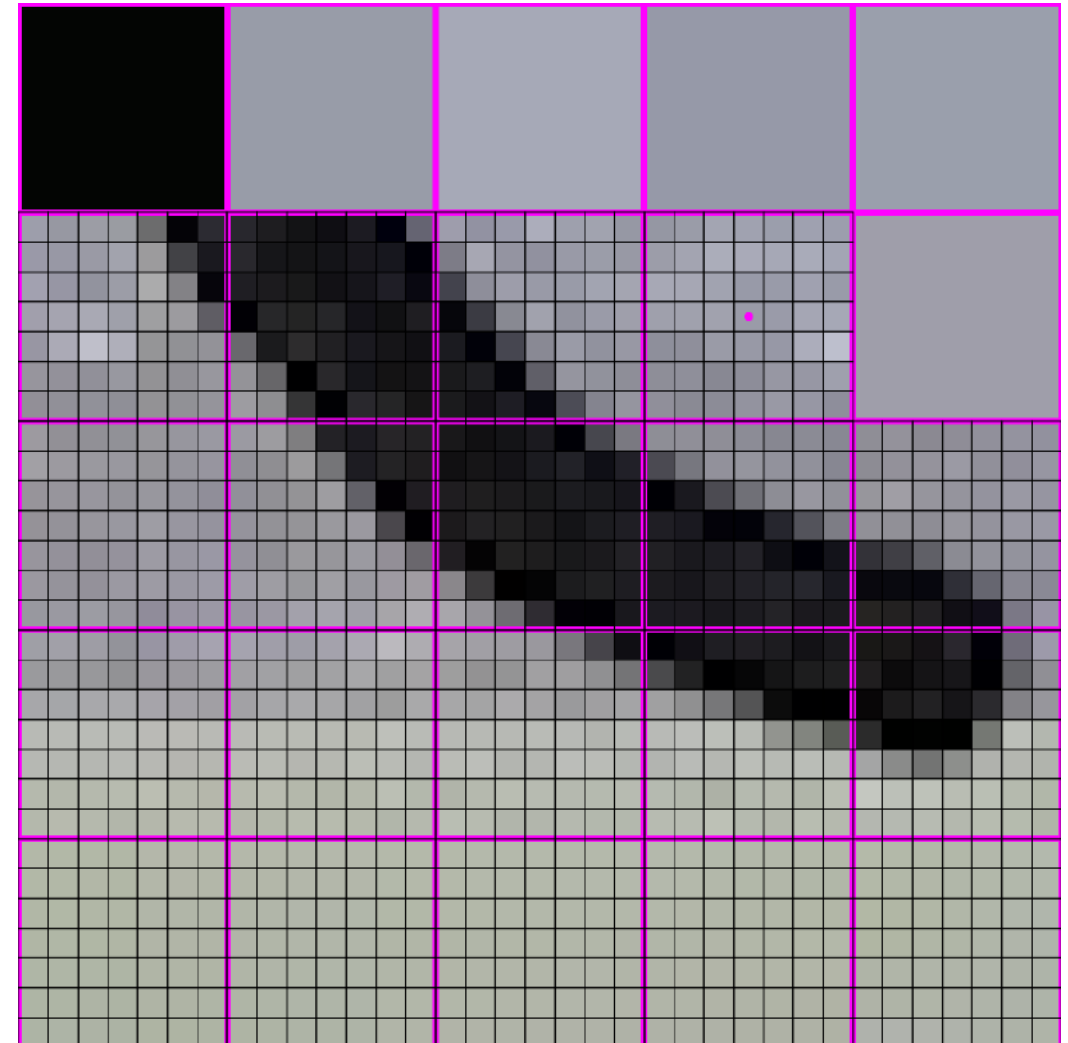


## 448x448 -> 64x64



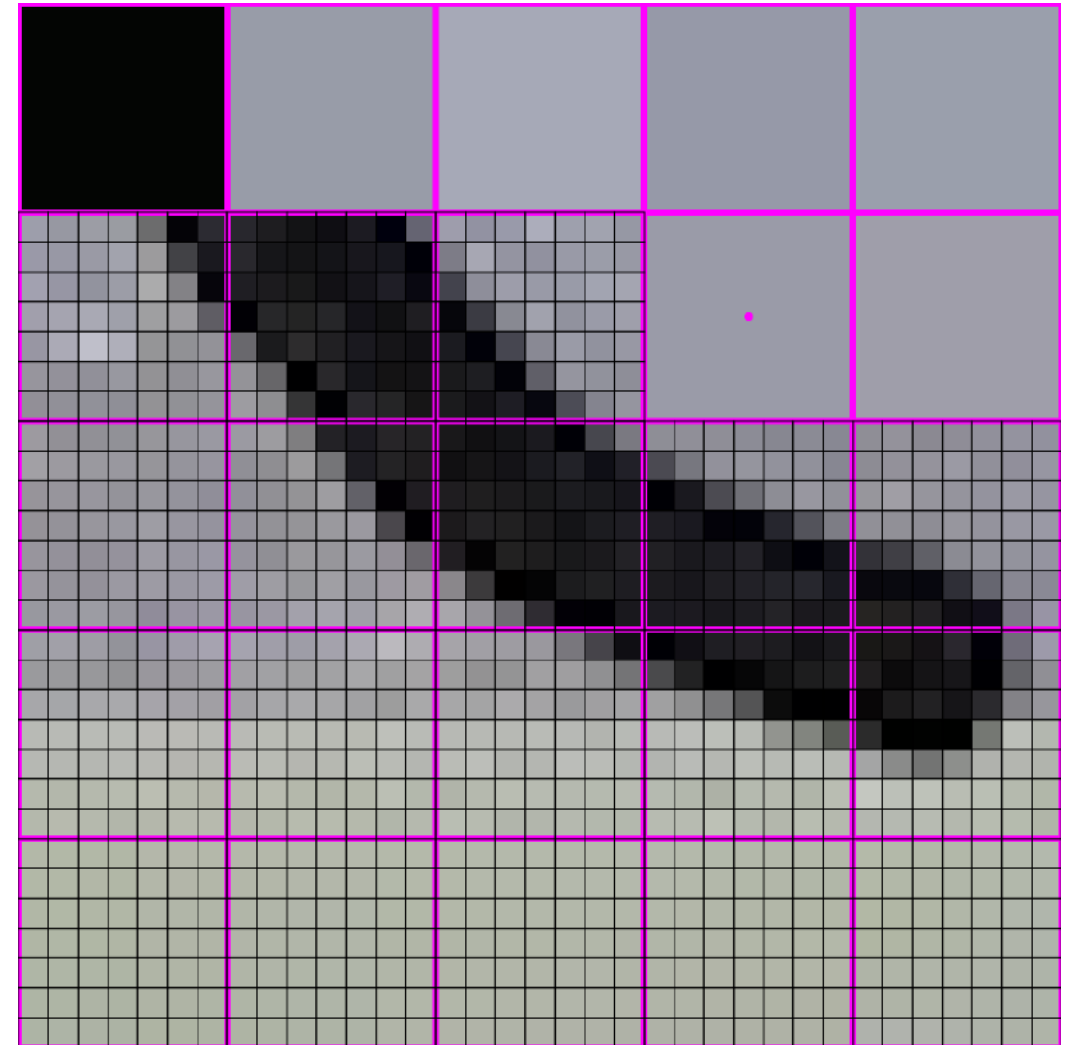


## 448x448 -> 64x64



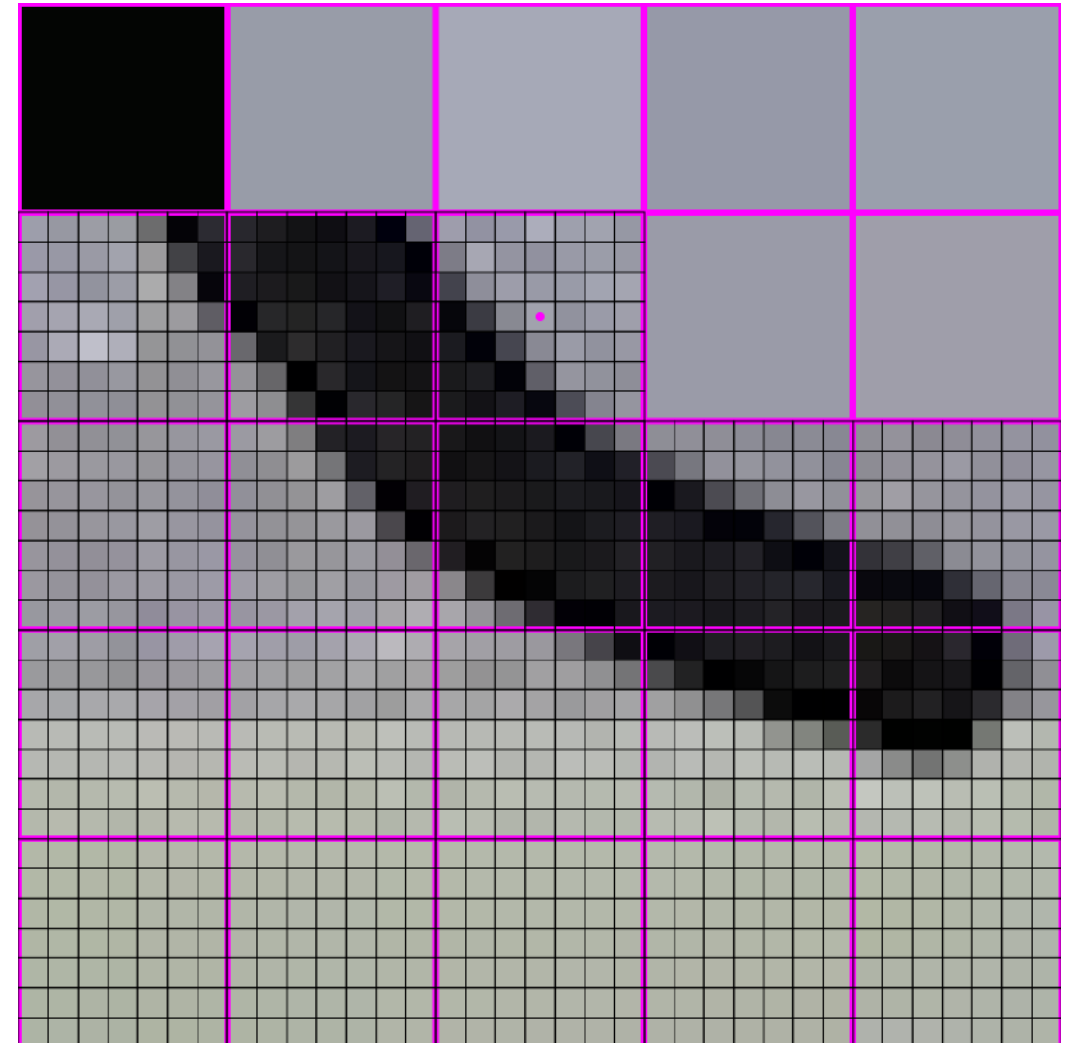


## 448x448 -> 64x64



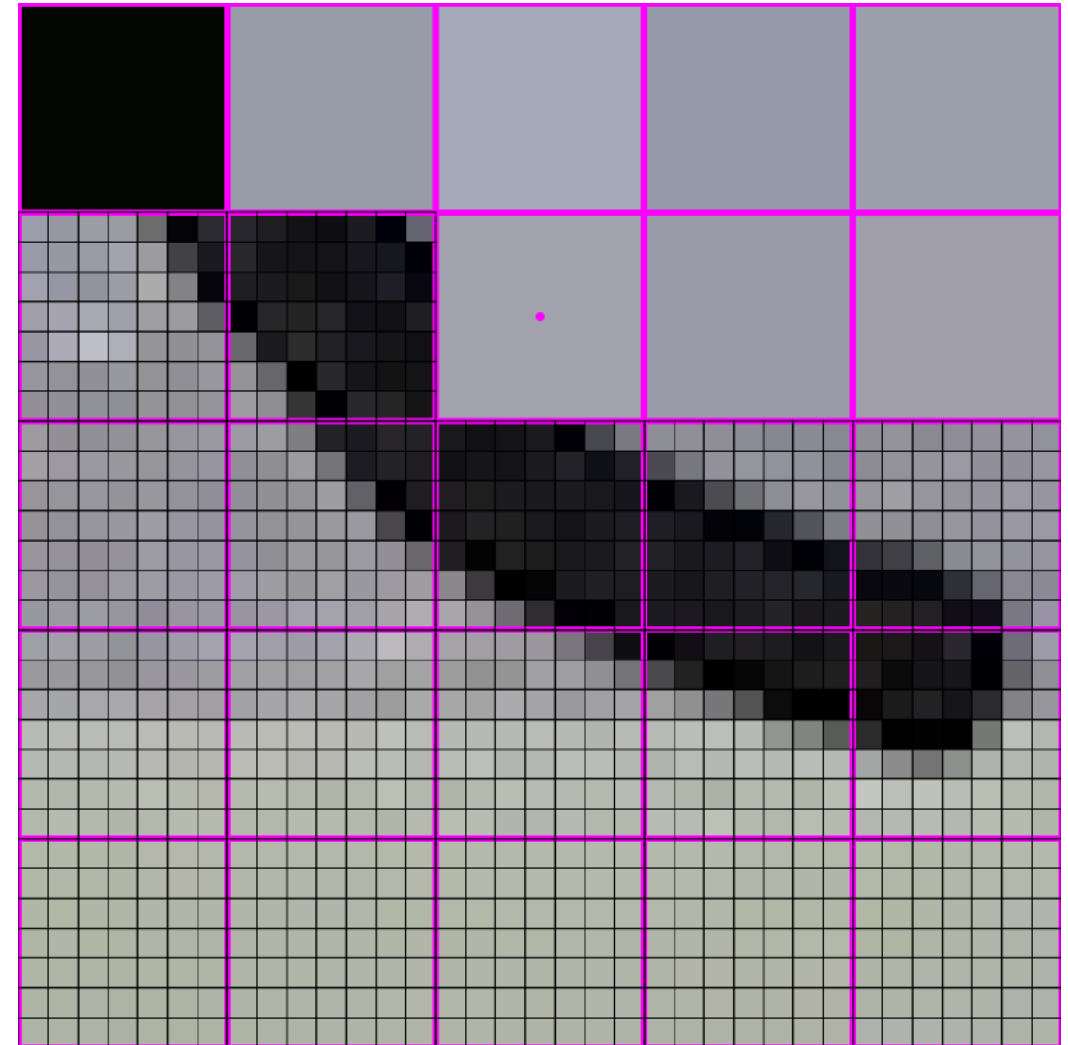


## 448x448 -> 64x64



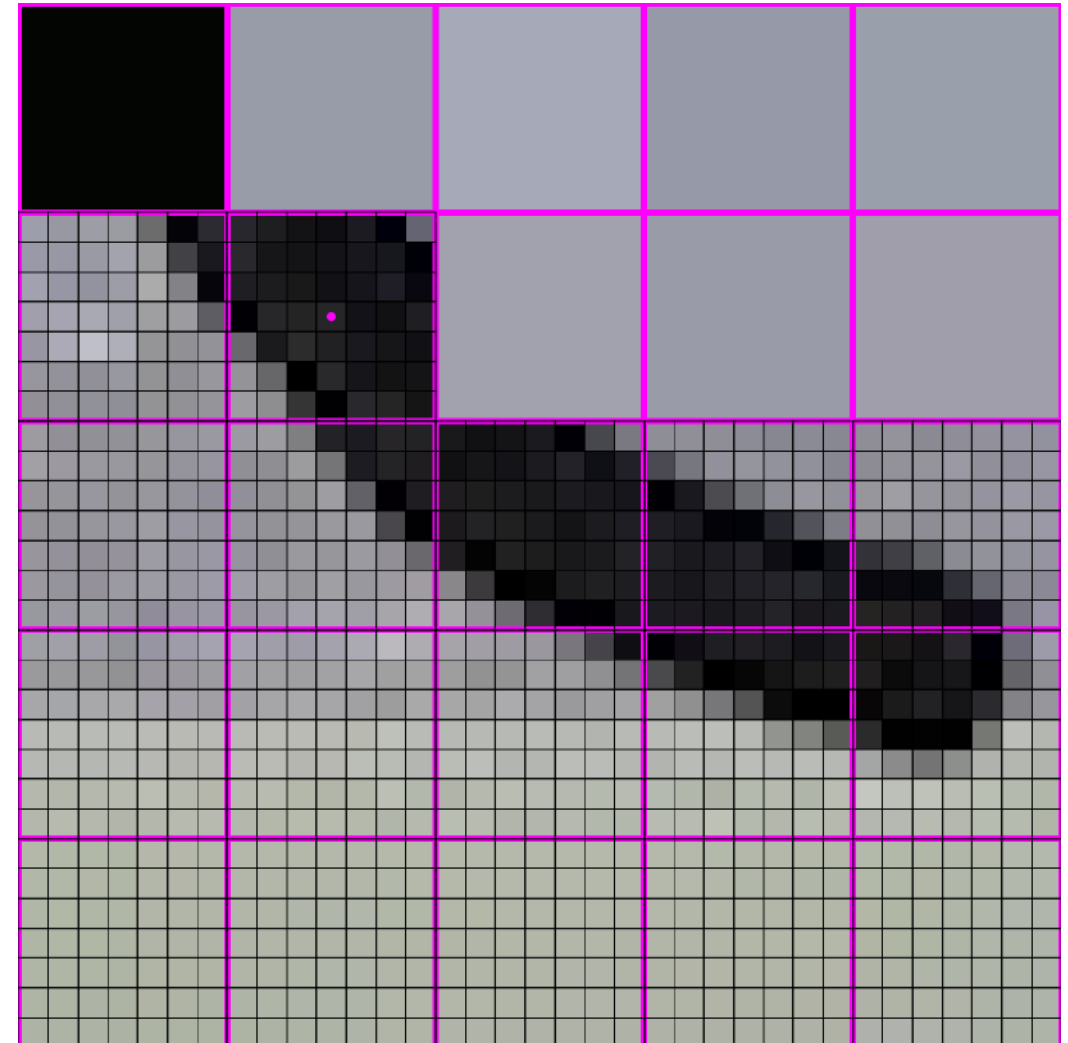


## 448x448 -> 64x64



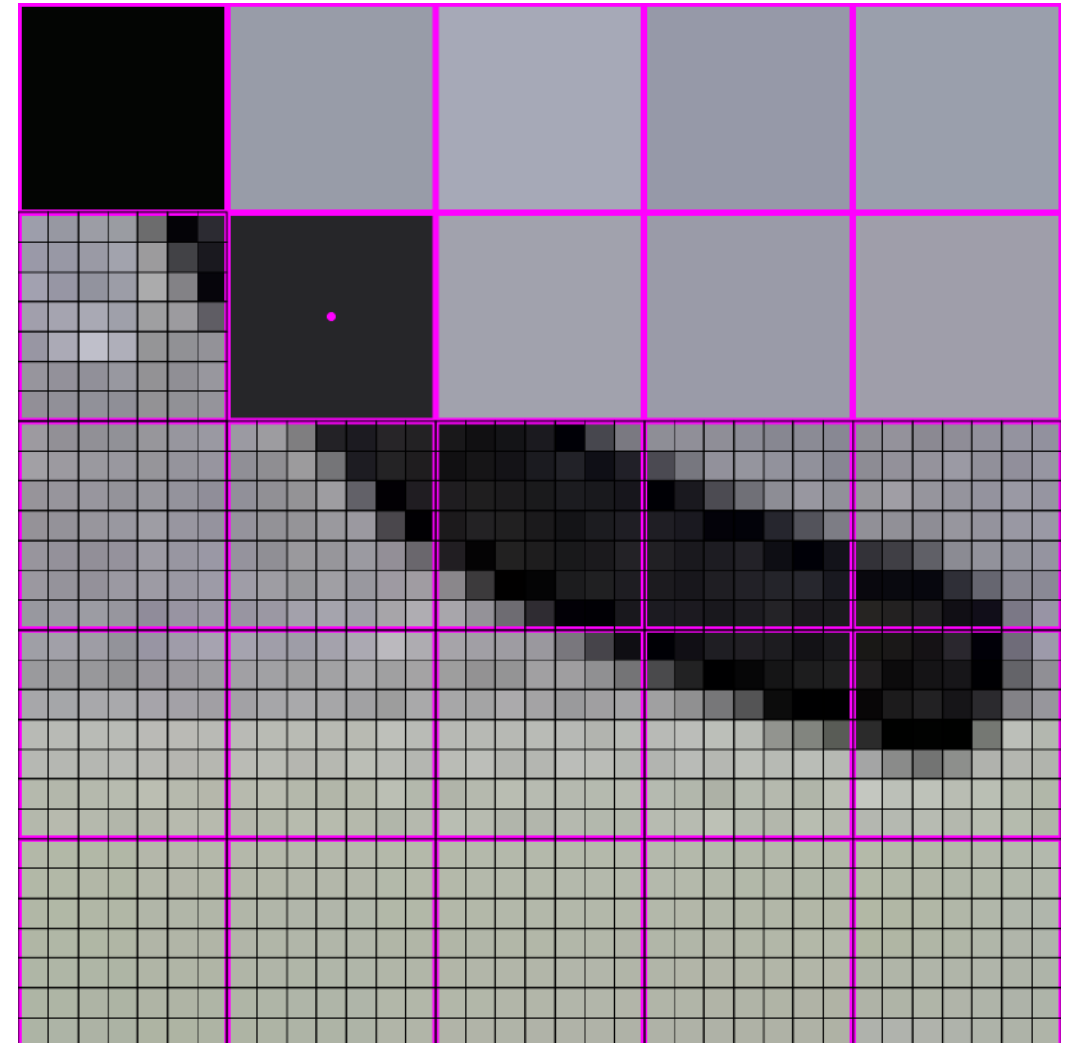


## 448x448 -> 64x64



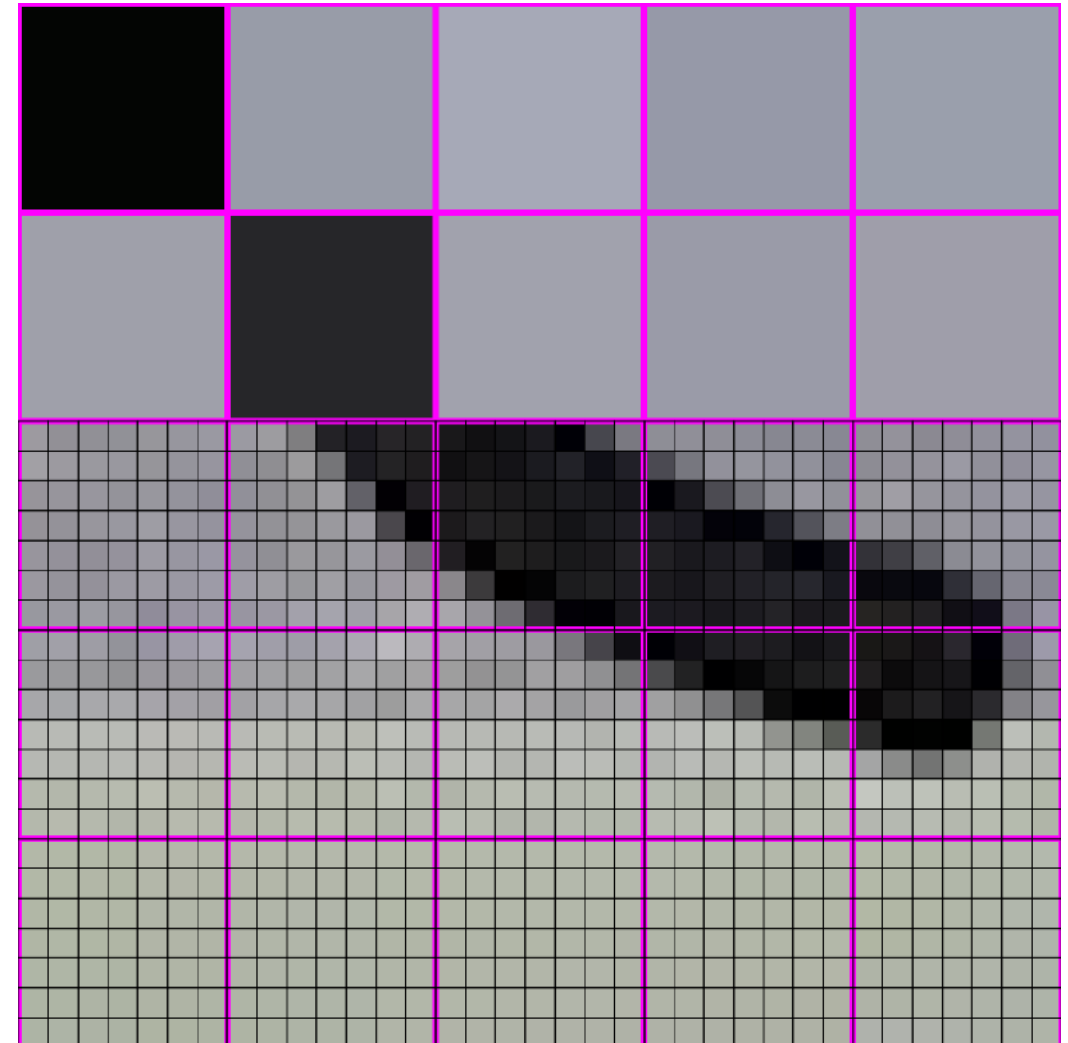


## 448x448 -> 64x64



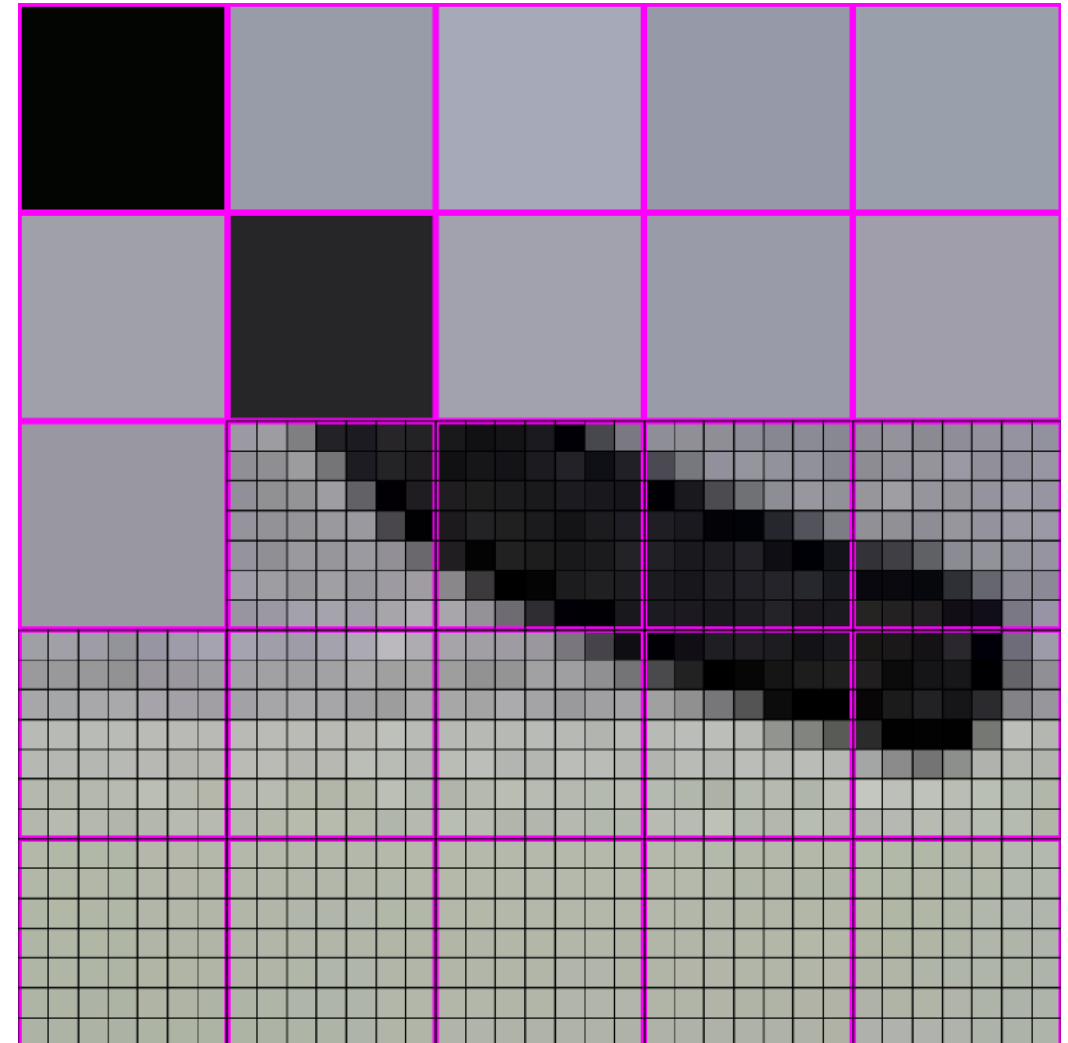


## 448x448 -> 64x64



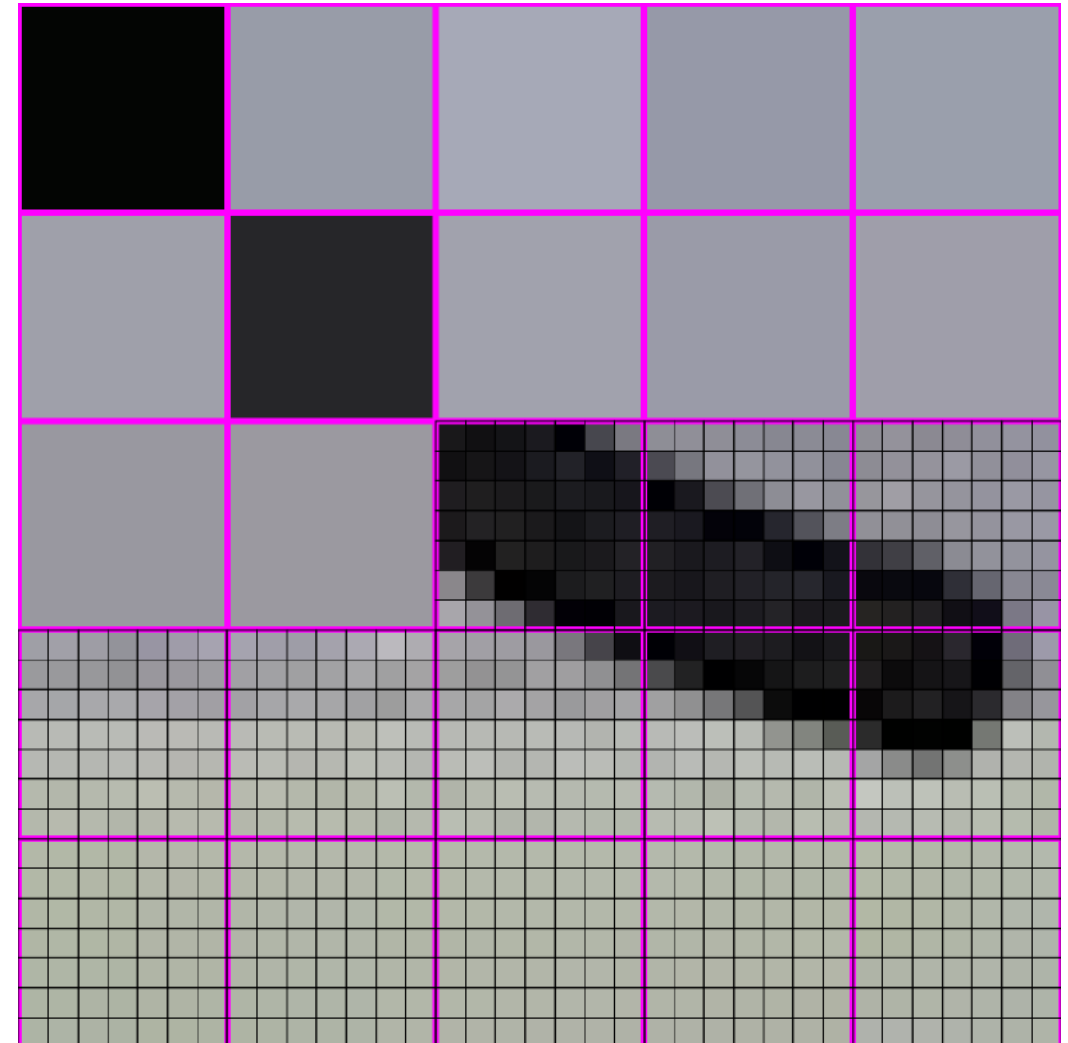


## 448x448 -> 64x64



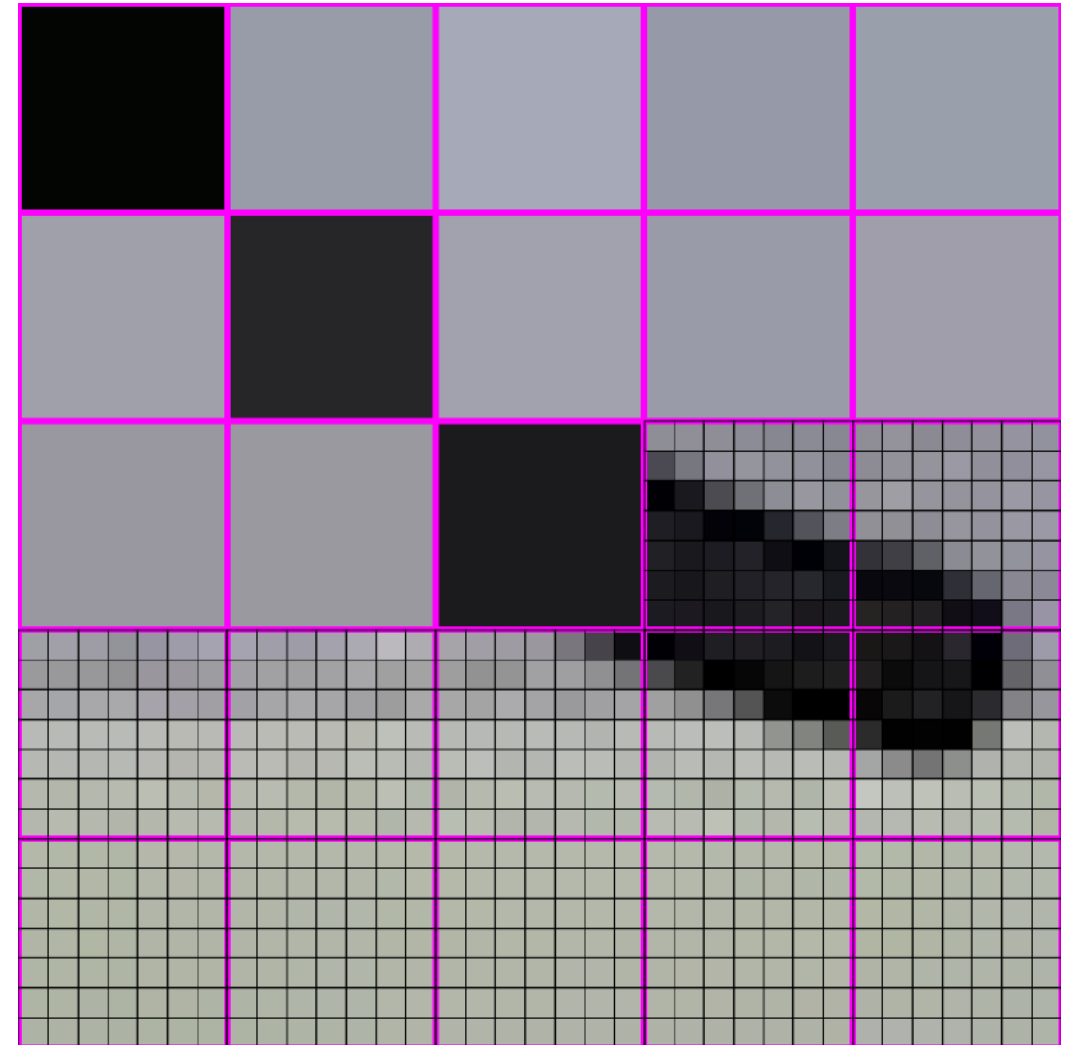


## 448x448 -> 64x64



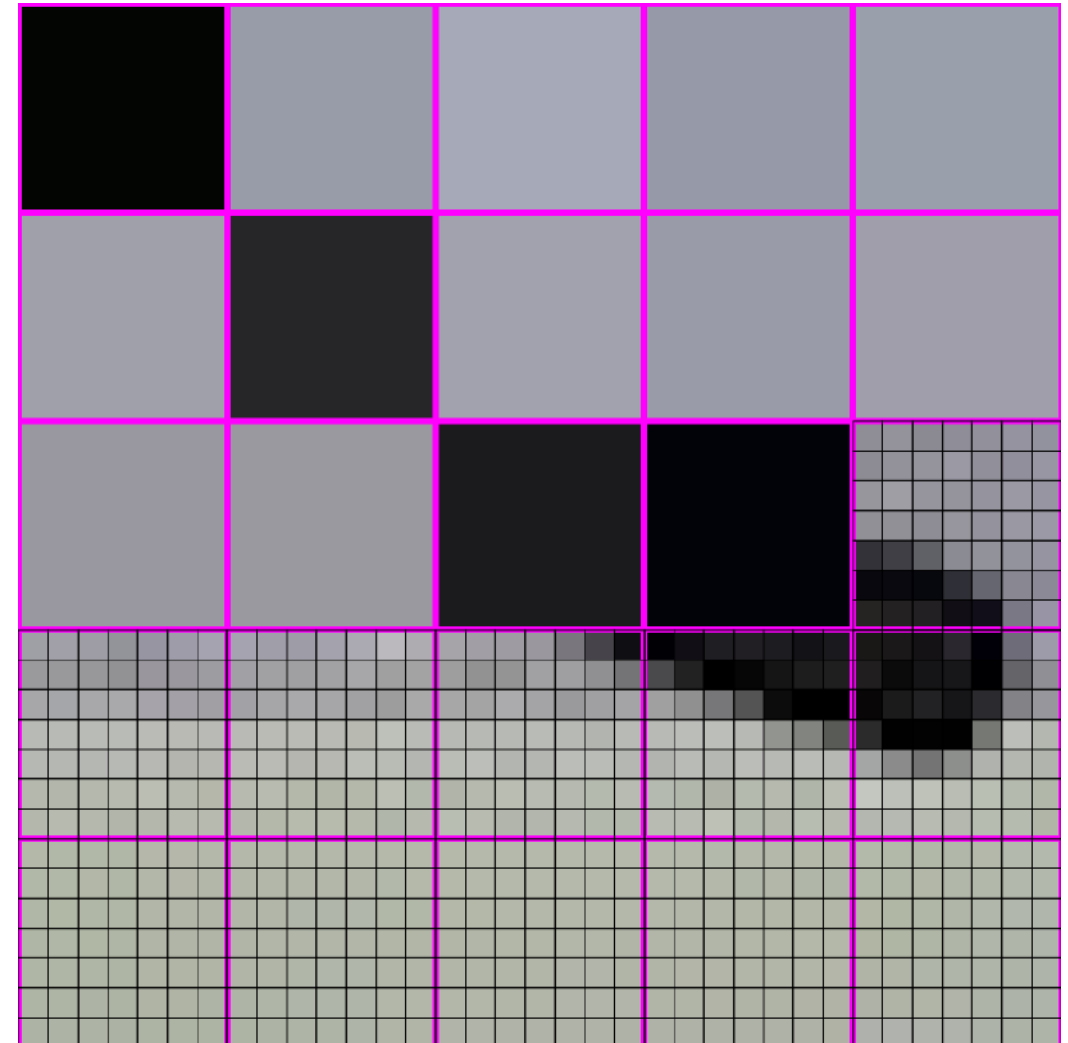


## 448x448 -> 64x64



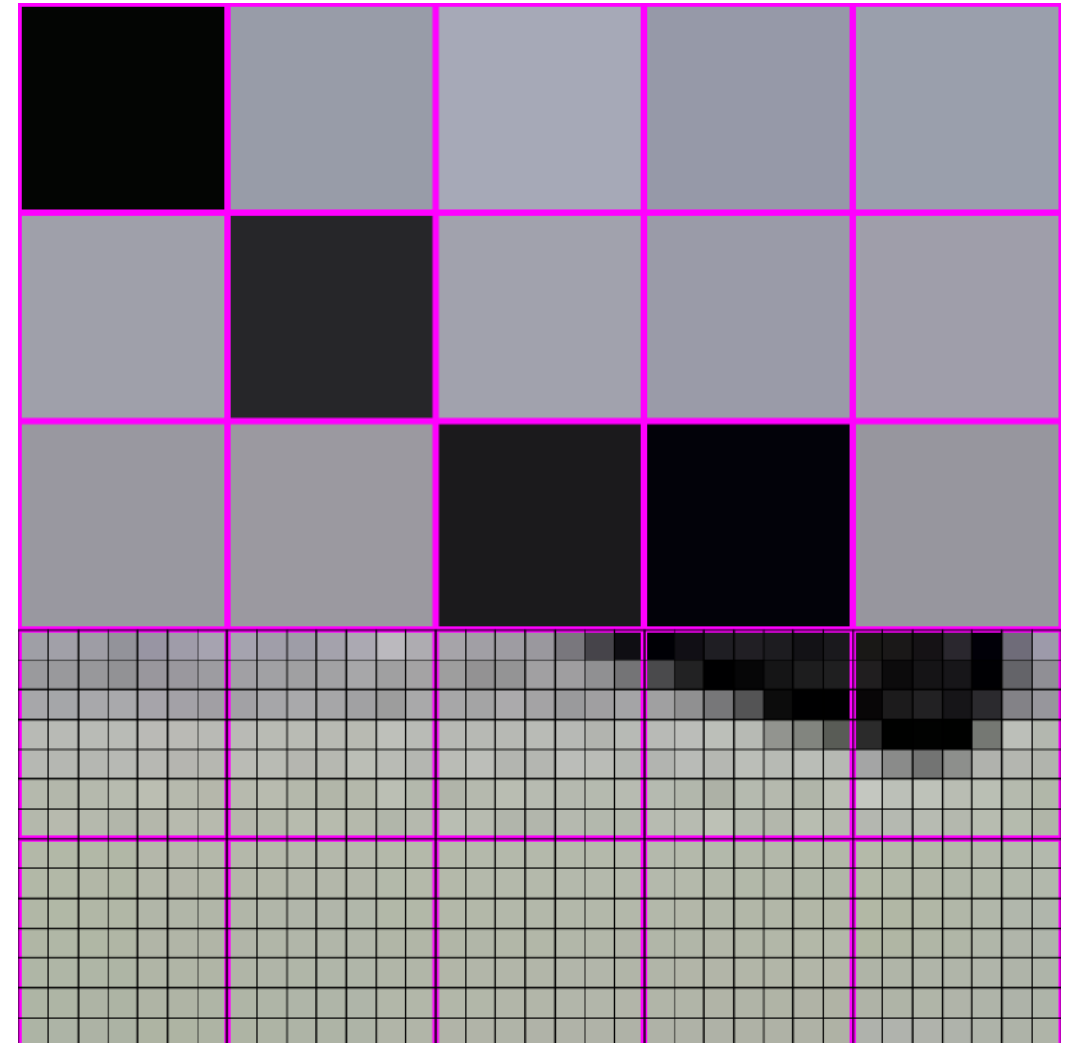


## 448x448 -> 64x64



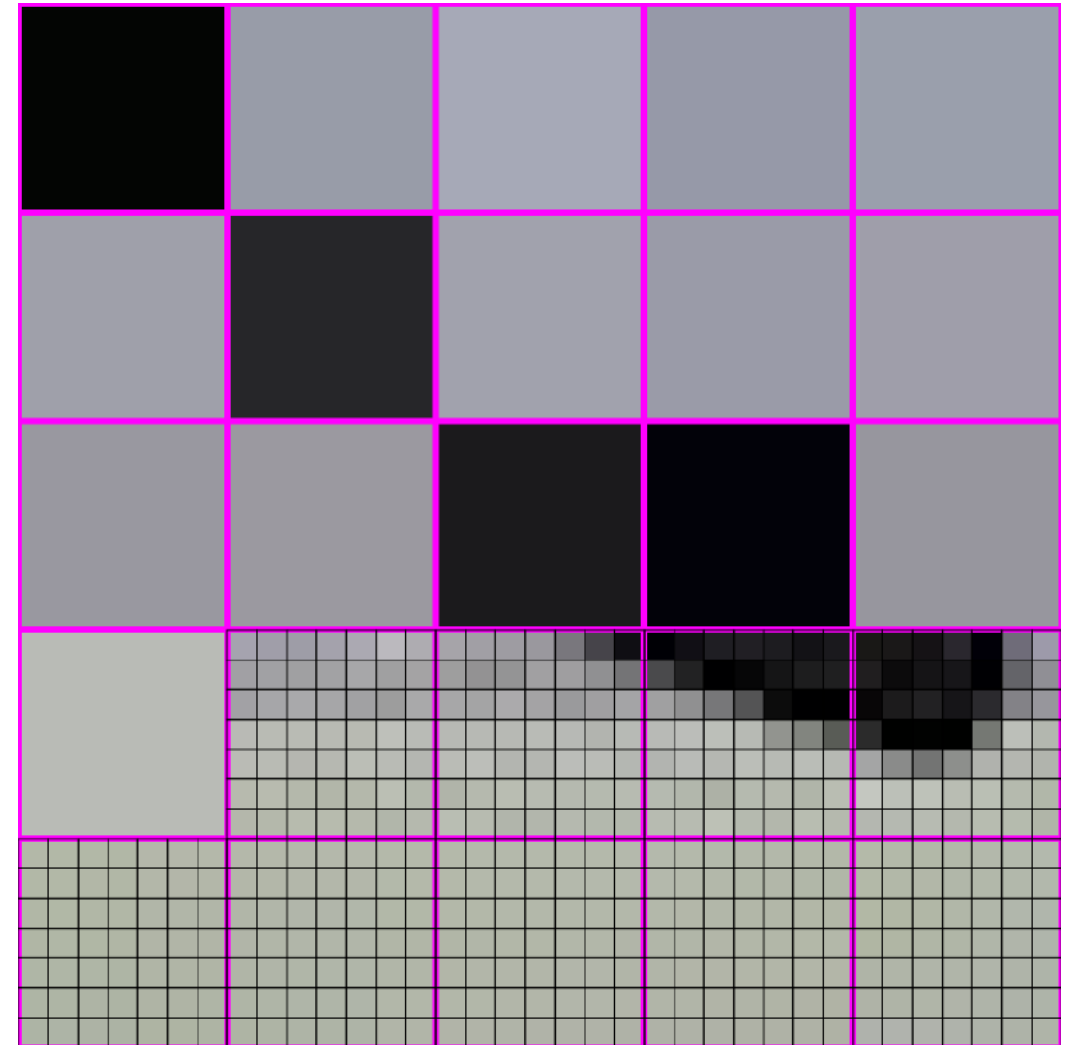


## 448x448 -> 64x64



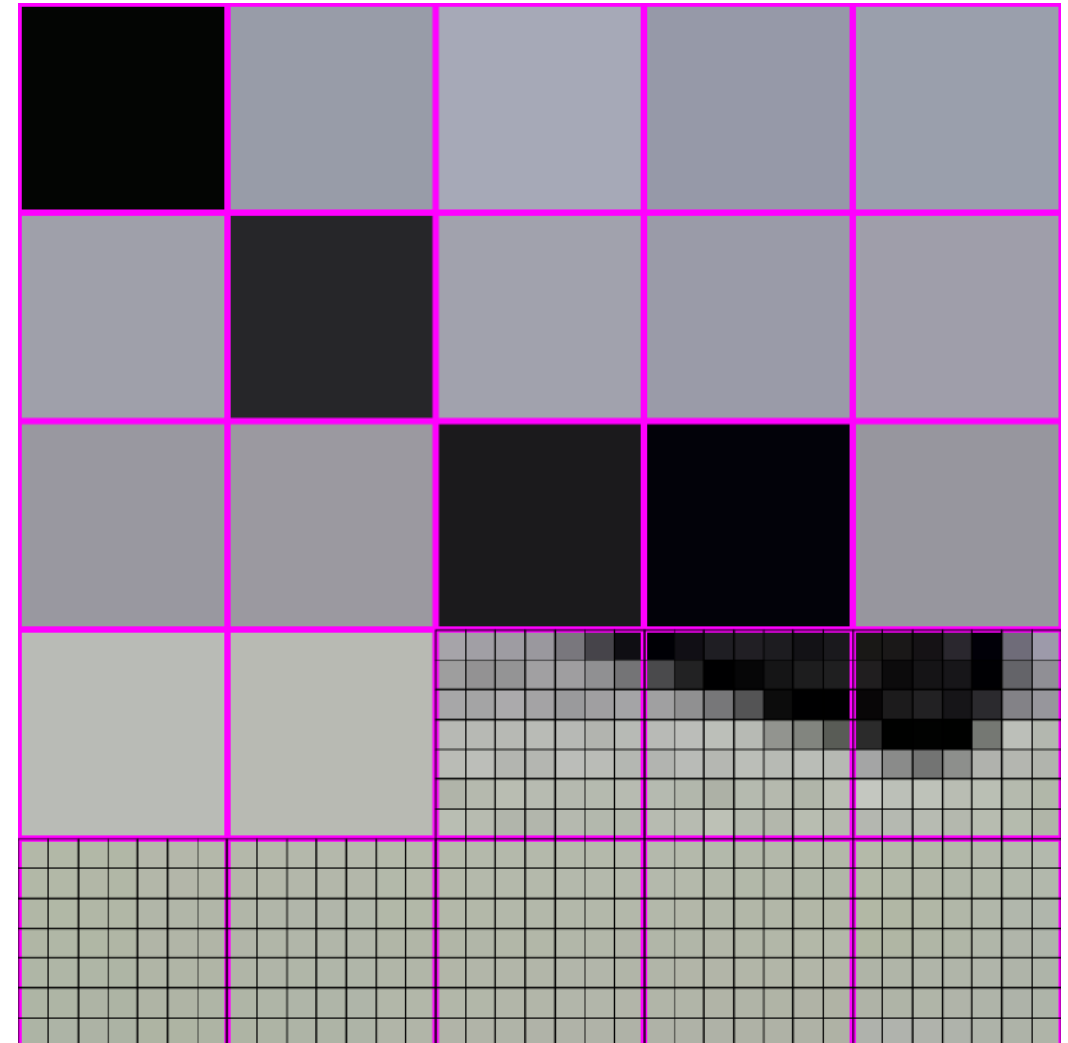


## 448x448 -> 64x64



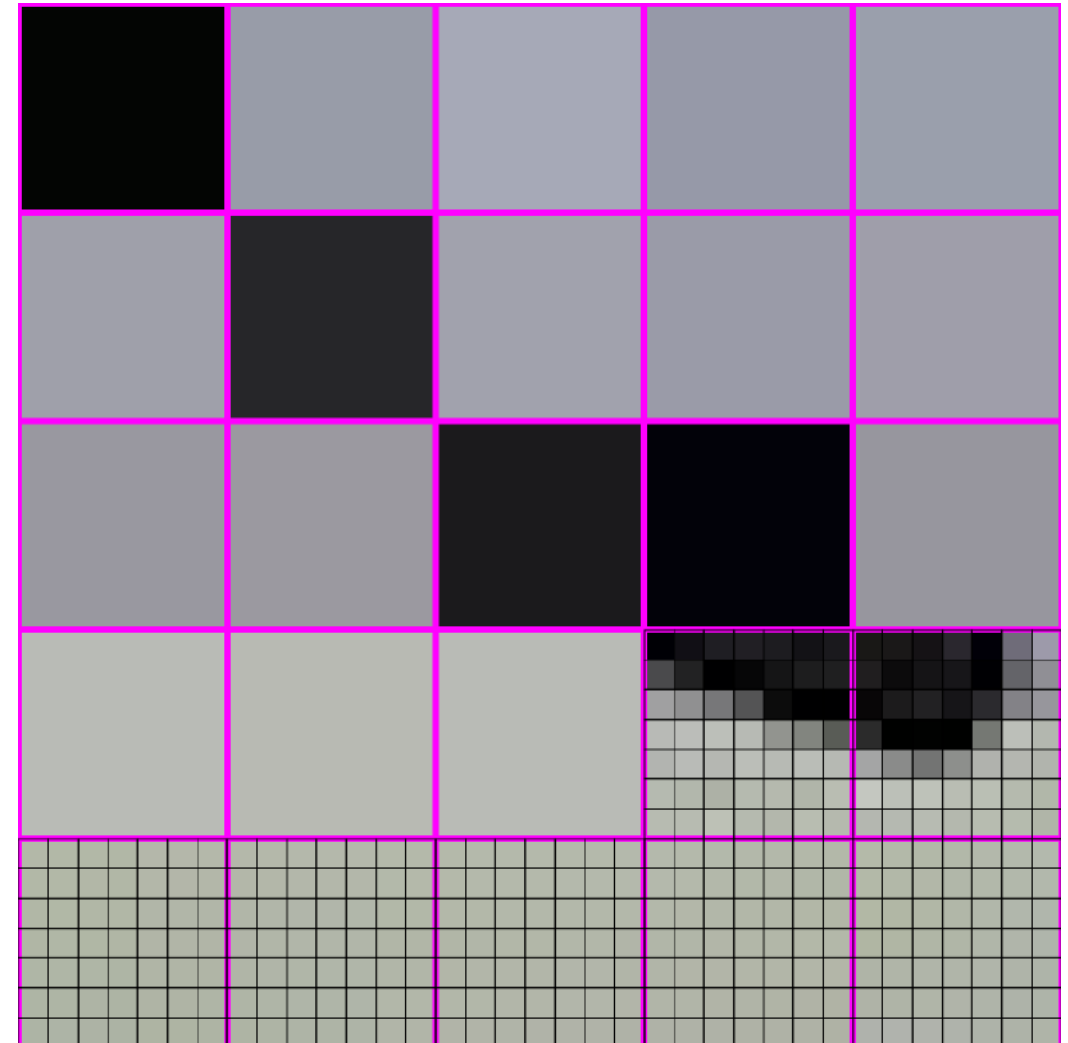


## 448x448 -> 64x64



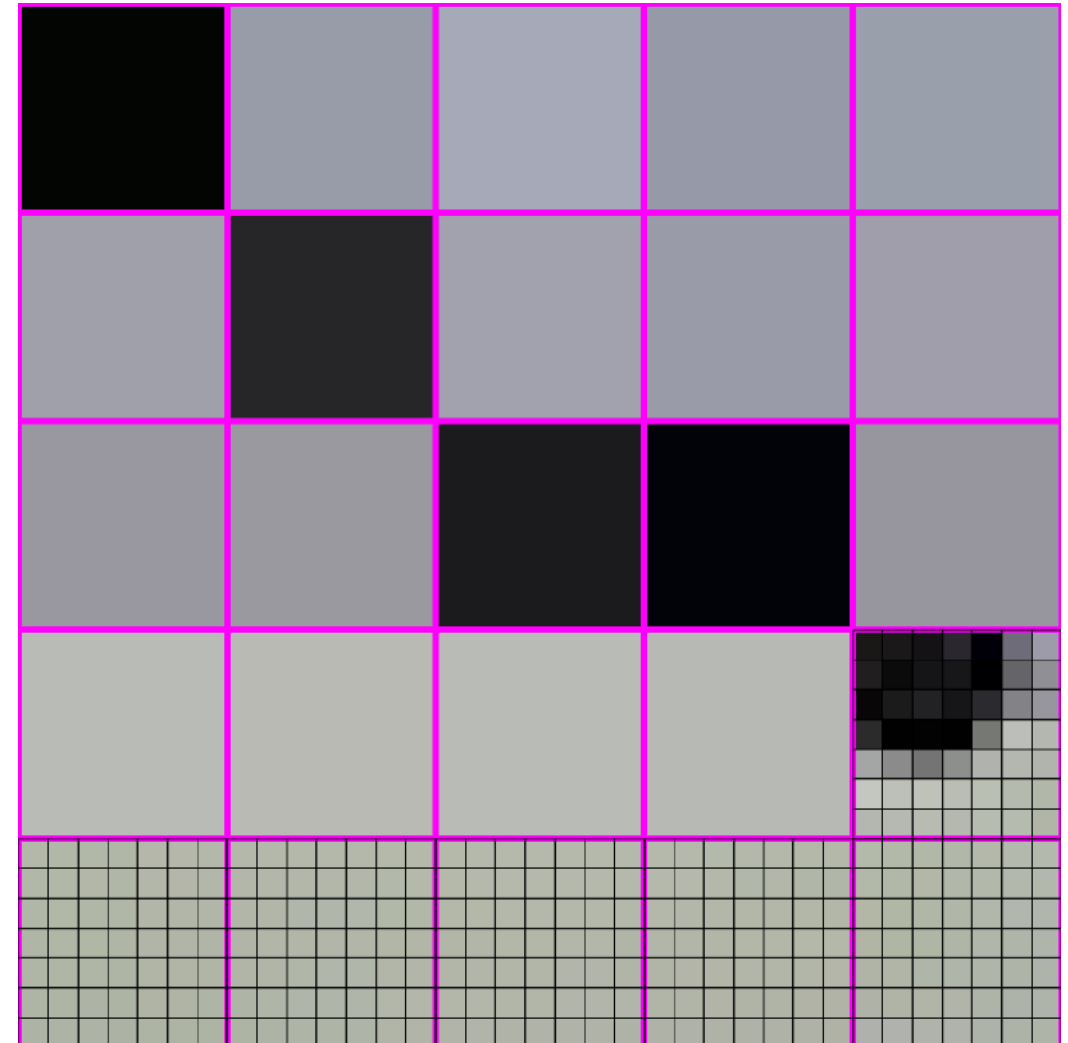


## 448x448 -> 64x64



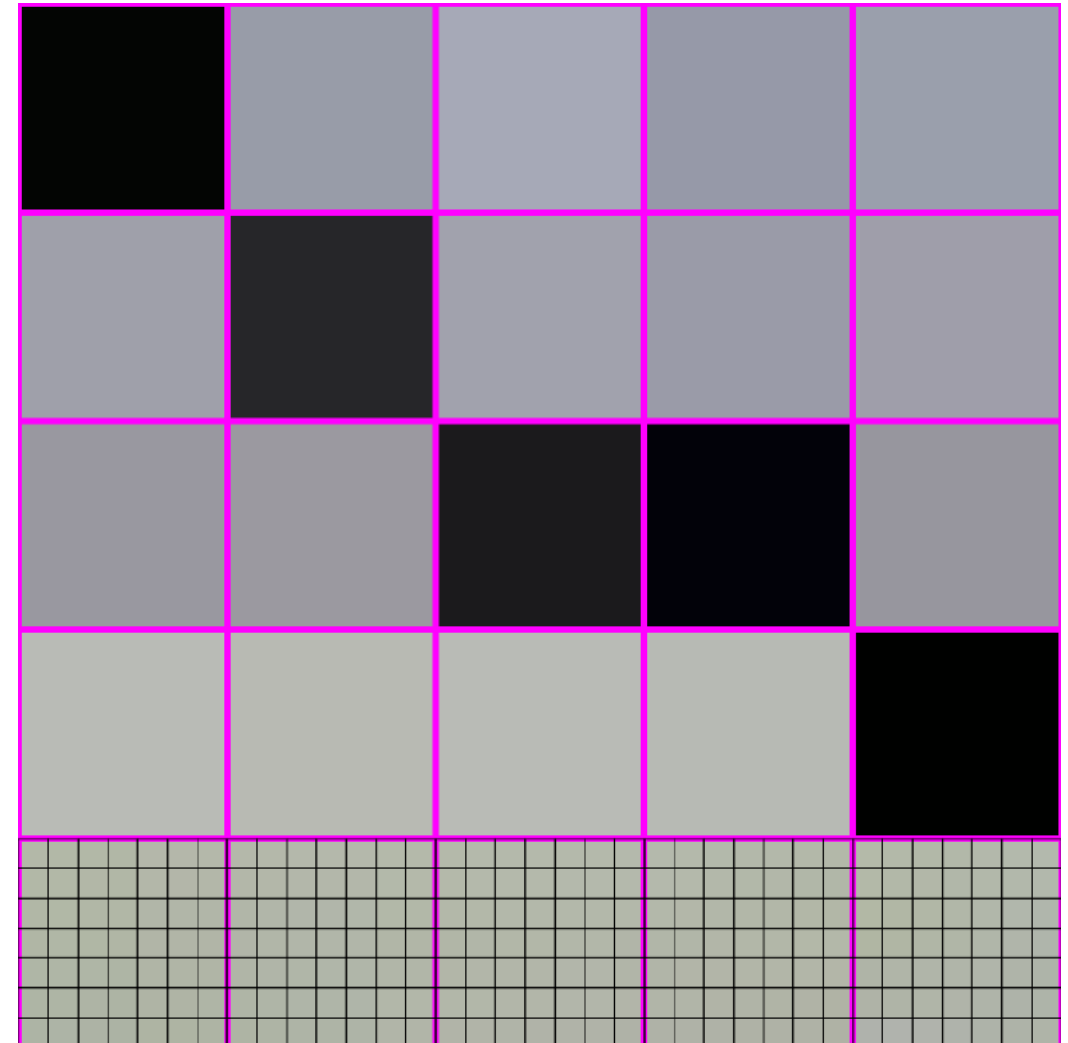


## 448x448 -> 64x64



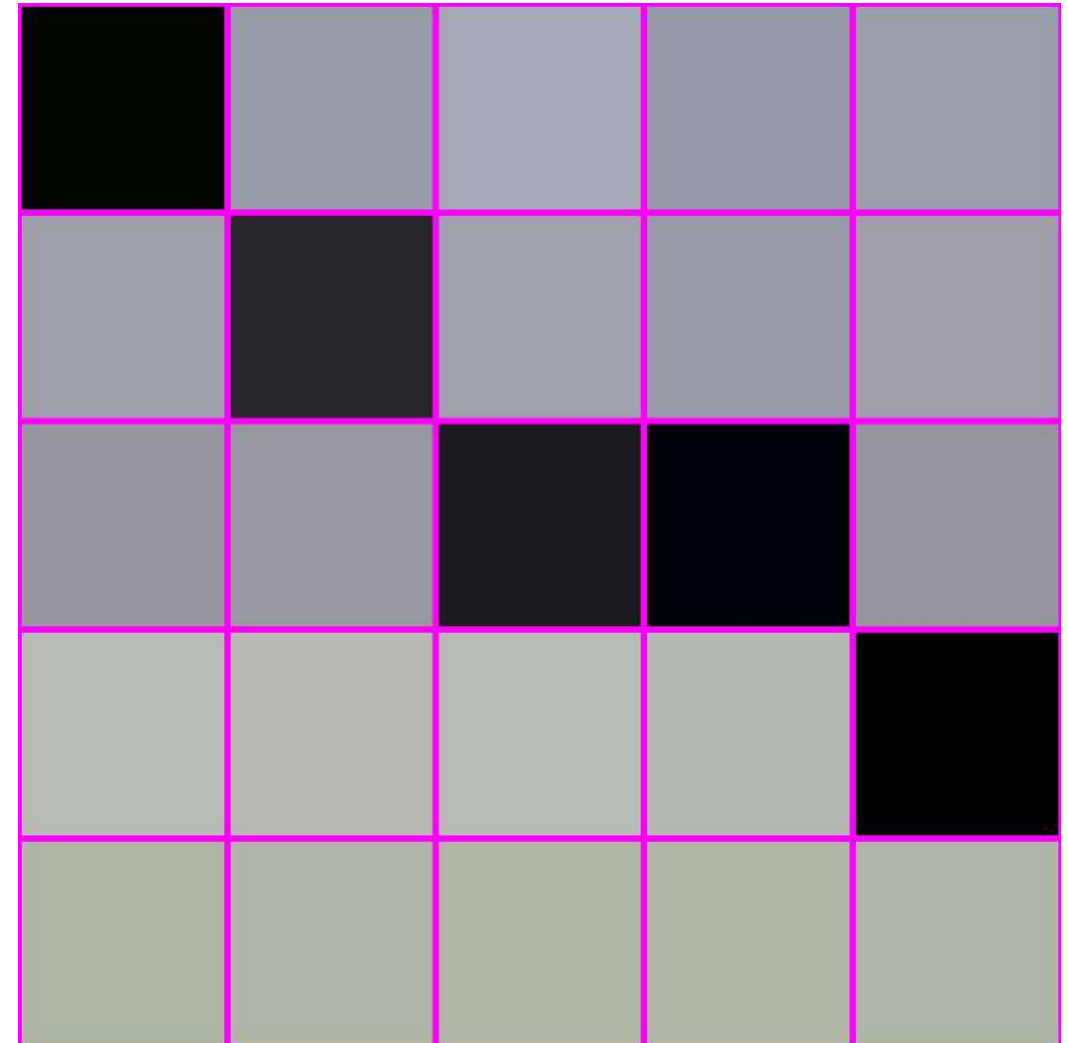


## 448x448 -> 64x64



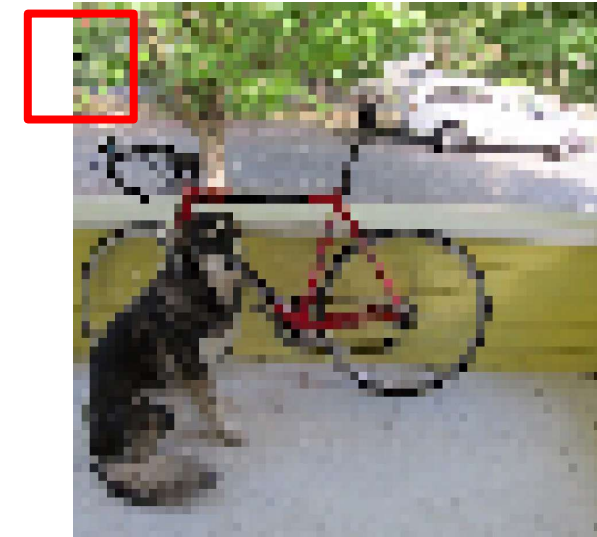
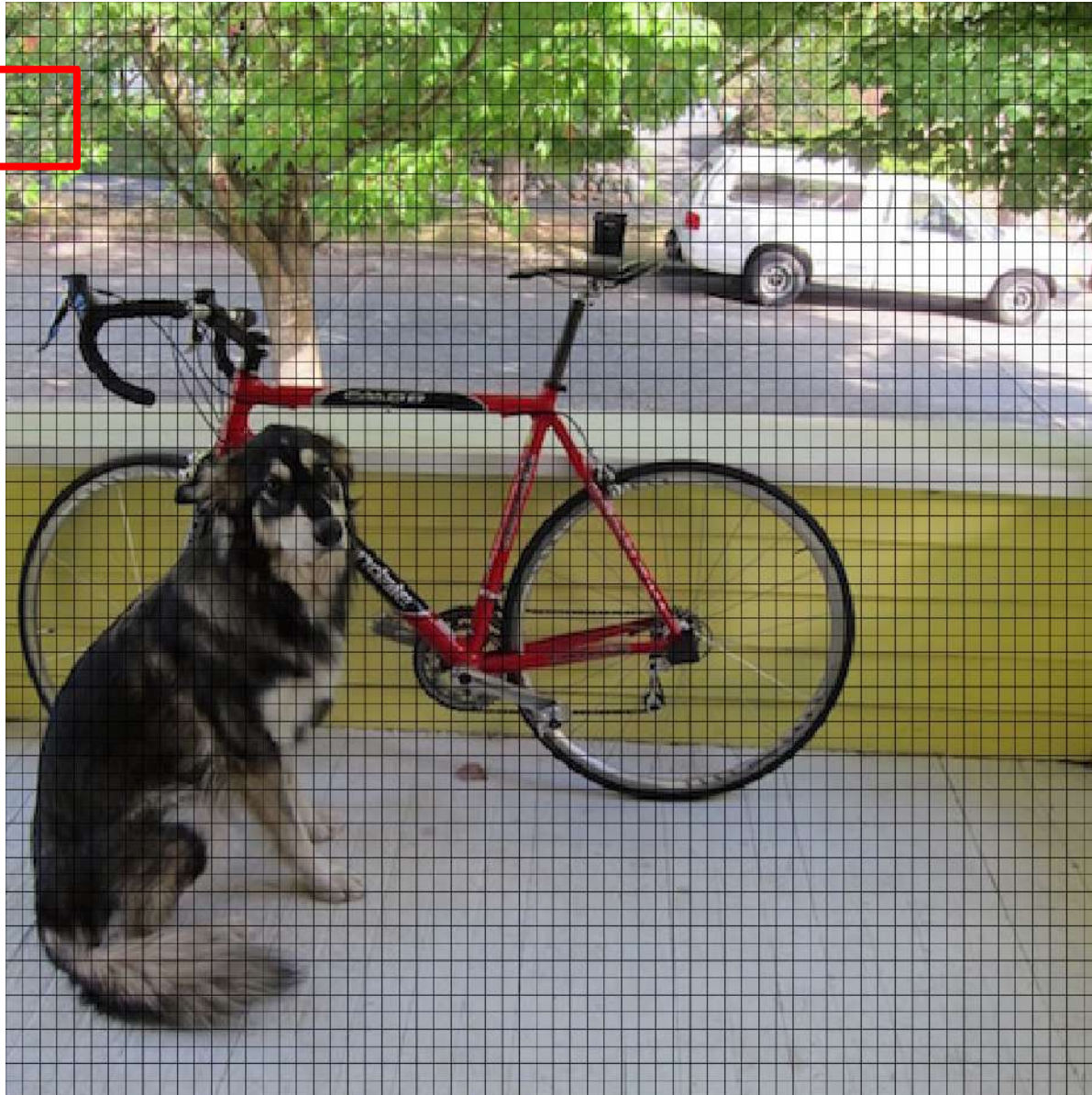


## 448x448 -> 64x64



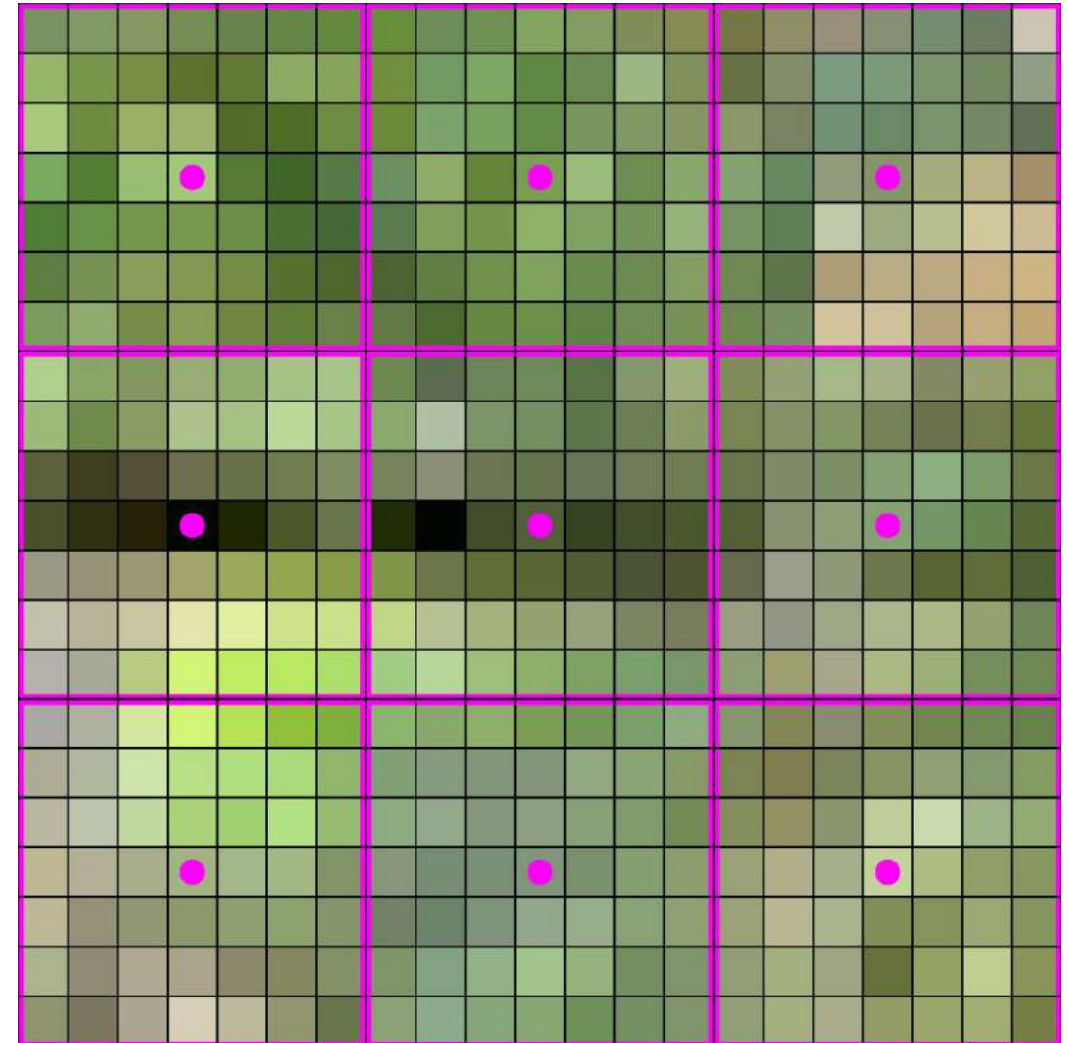
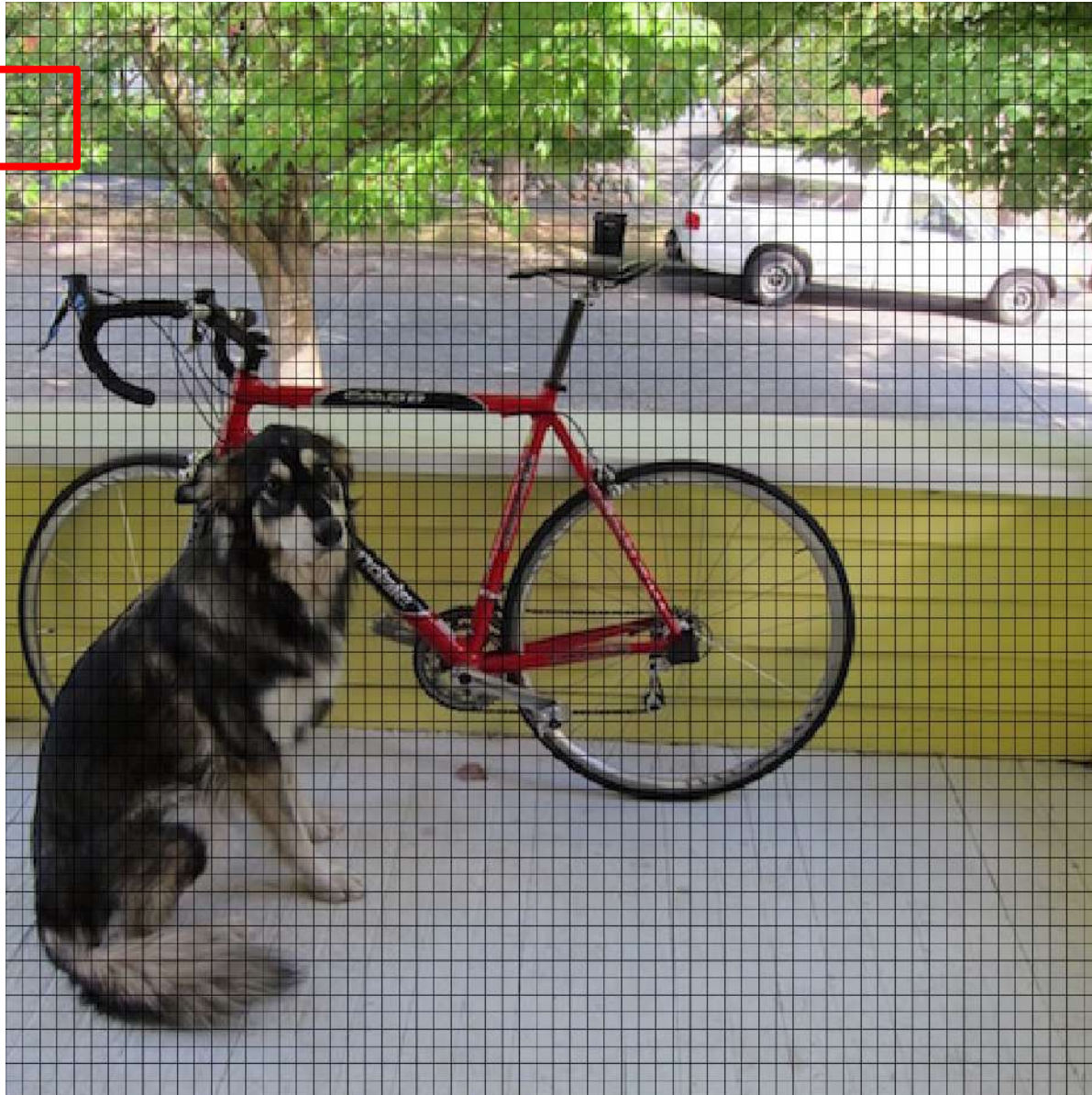


448x448 -> 64x64



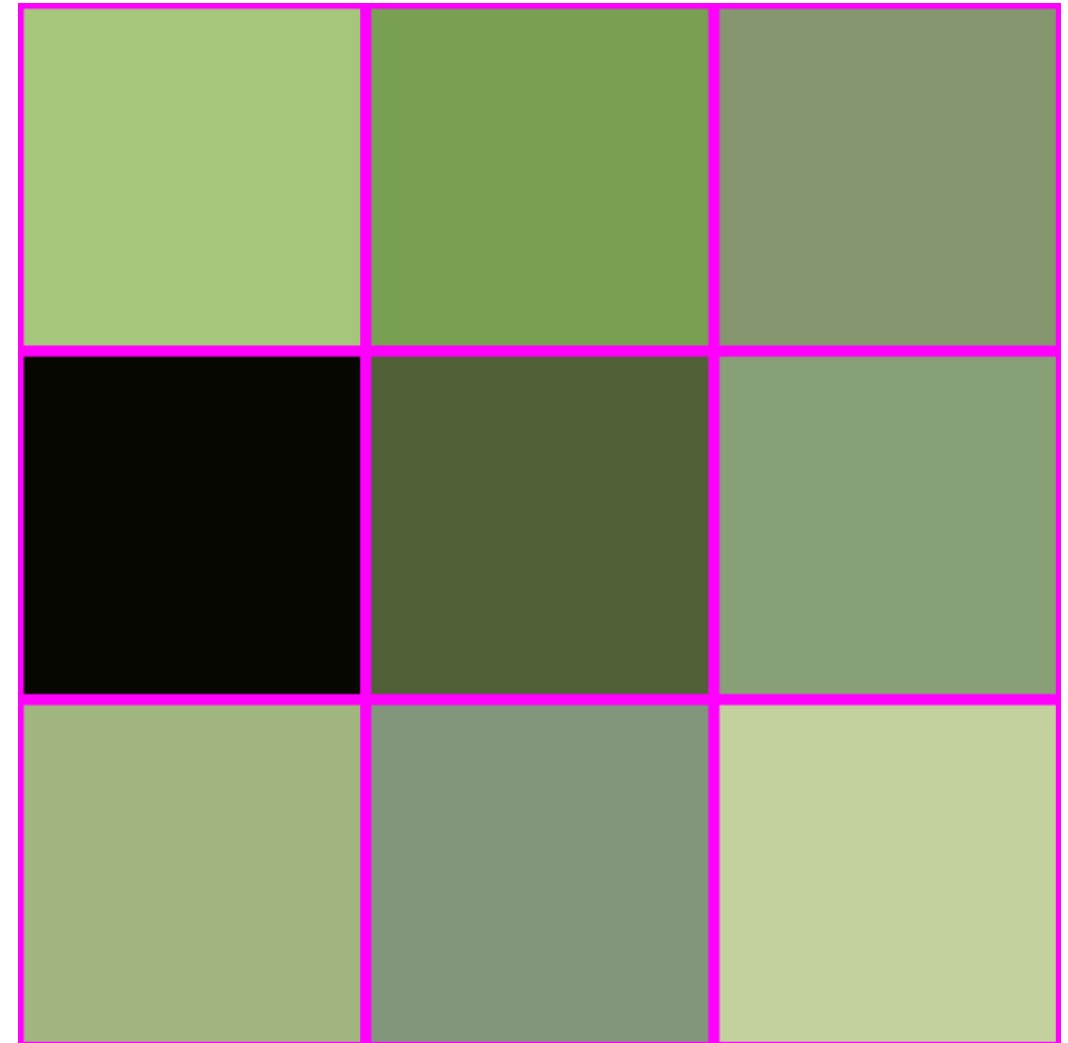
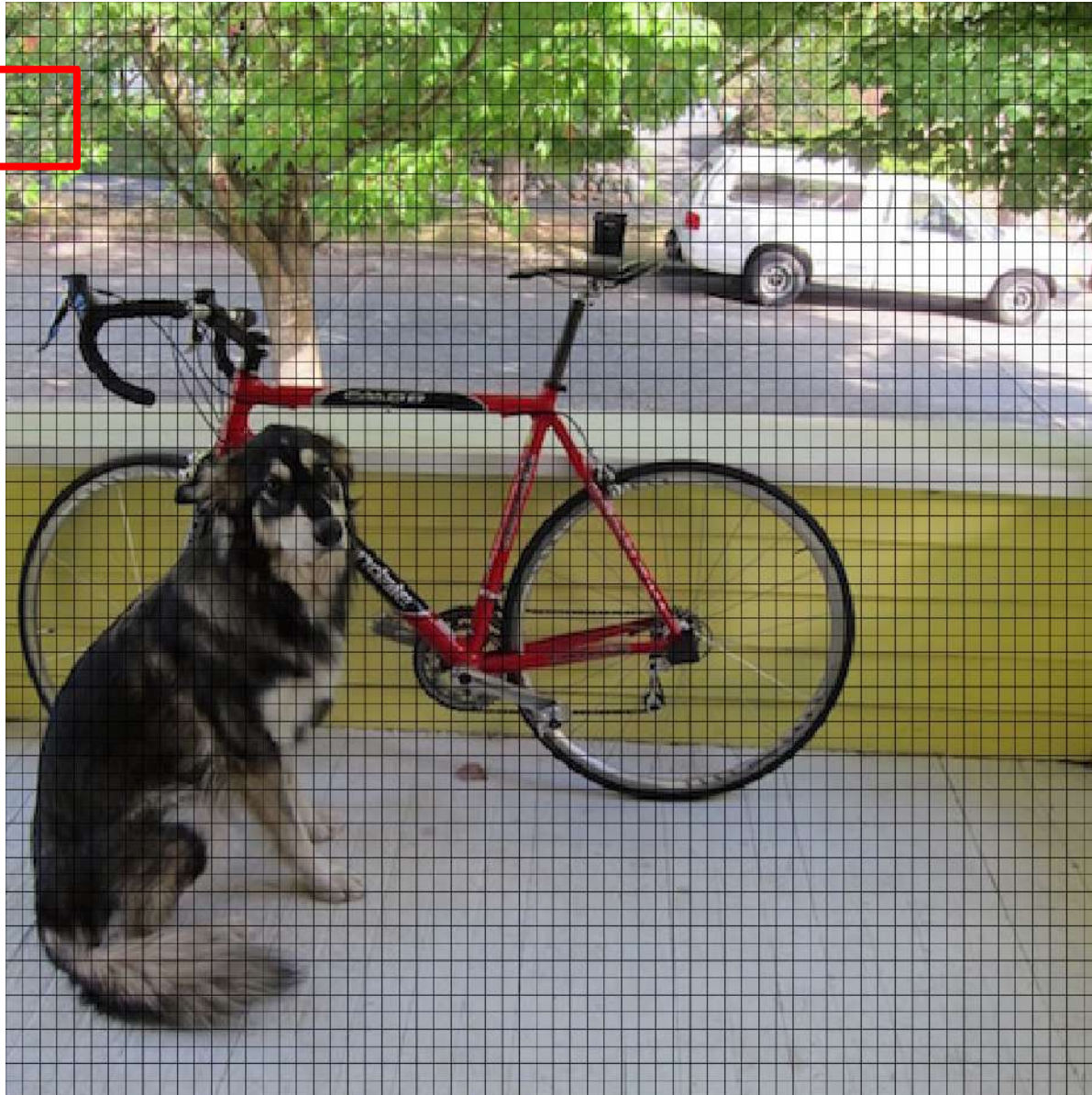


## 448x448 -> 64x64





## 448x448 -> 64x64





## Lots of issues

- NN and Bilinear only look at small area
- Lots of artifacting
- Staircase pattern on diagonal lines
- We'll fix this with filters!



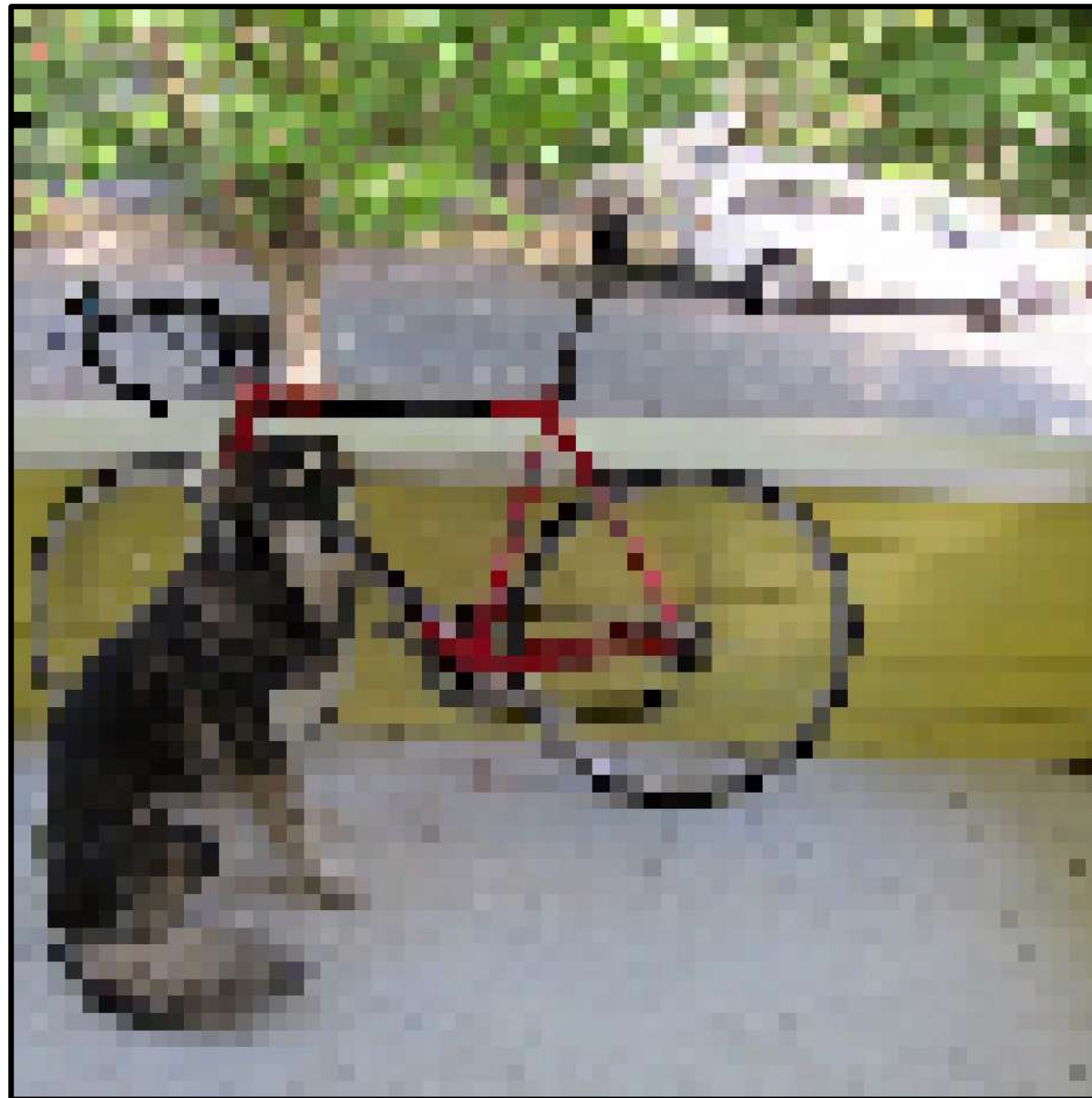
NN



Bilinear

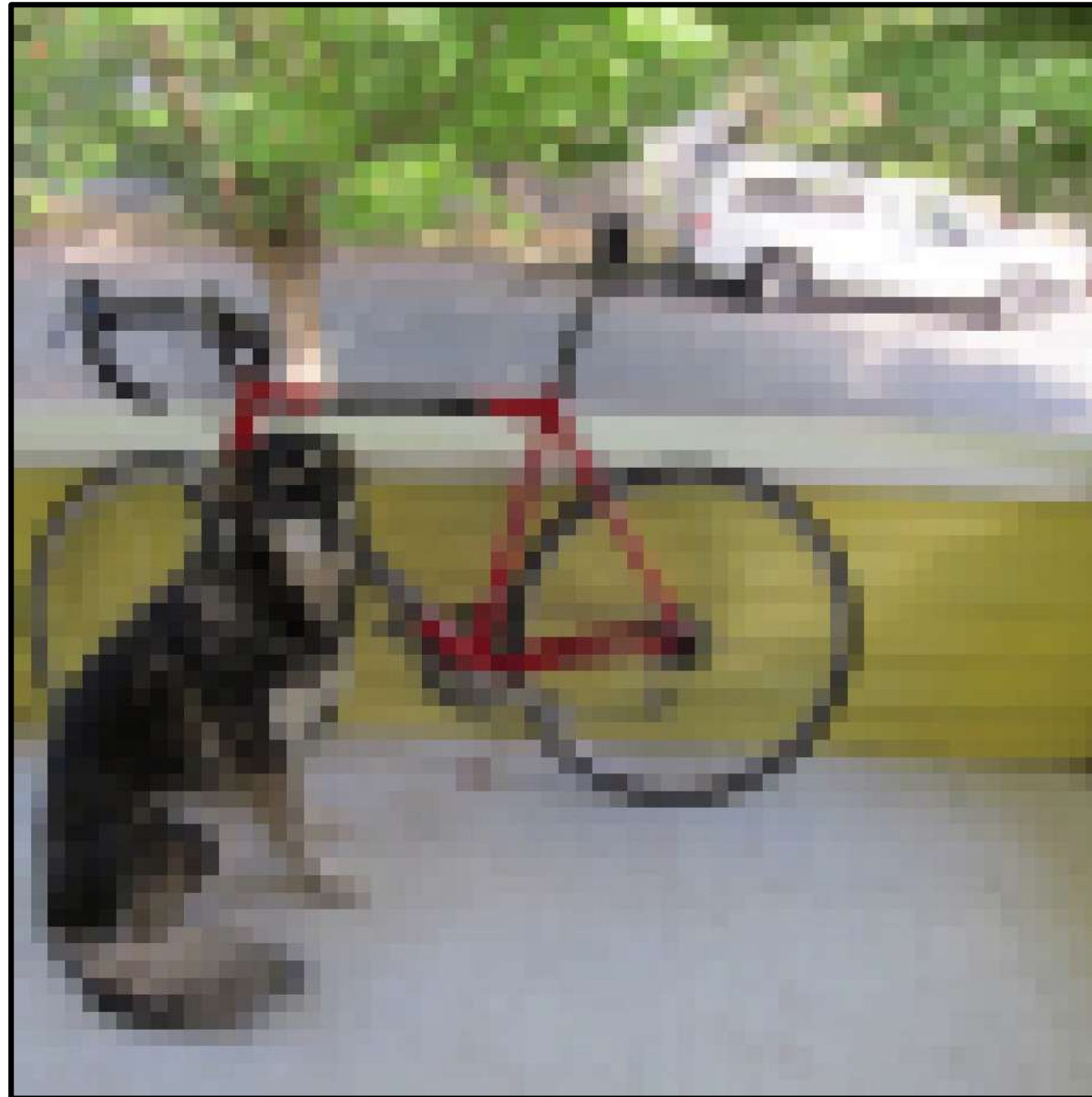


# IS THIS ALL THERE IS??





# THERE IS A BETTER WAY!



**MAI4CAREU**

Master programmes in Artificial  
Intelligence 4 Careers in Europe



# Thank you.







University of Cyprus – MSc Artificial Intelligence

# MAI644 – COMPUTER VISION

## Lecture 5: Filters – Convolution

**Melinos Averkiou**

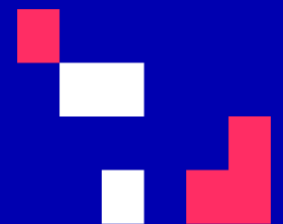
CYENS Centre of Excellence

University of Cyprus - Department of Computer Science

[m.averkiou@cyens.org.cy](mailto:m.averkiou@cyens.org.cy)



**CYENS**  
CENTRE OF EXCELLENCE



# Last time

- Image basics
  - What is an image – addressing pixels
  - Image as a function – image coordinates
- Image interpolation
  - Nearest neighbor
  - Bilinear
  - Bicubic
- Image resizing
  - Enlarge
  - Shrink

# Today's Agenda

- Averaging vs Interpolation
- Systems - filters
- Convolution
  - Box Filter
  - Gaussian
  - Cross correlation vs Convolution
- Examples of filters

**[material based on Joseph Redmon's course]**

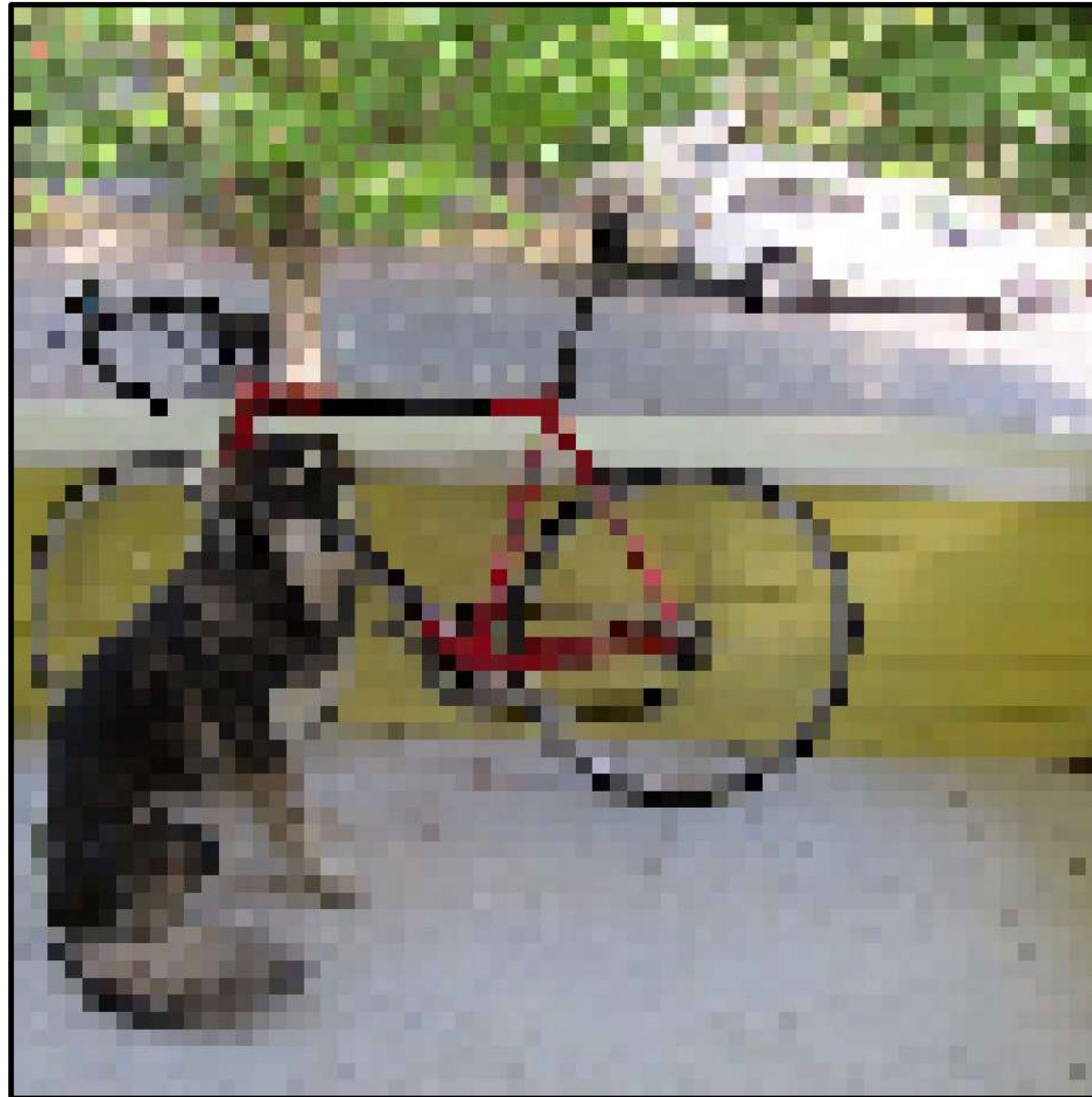
# Today's Agenda

- Averaging vs Interpolation
- Systems - filters
- Convolution
  - Box Filter
  - Gaussian
  - Cross correlation vs Convolution
- Examples of filters





## Is this all there is ??



## Lots of issues

- NN and Bilinear only look at small area
- Lots of artifacting
- Staircase pattern on diagonal lines
- We'll fix this with filters!

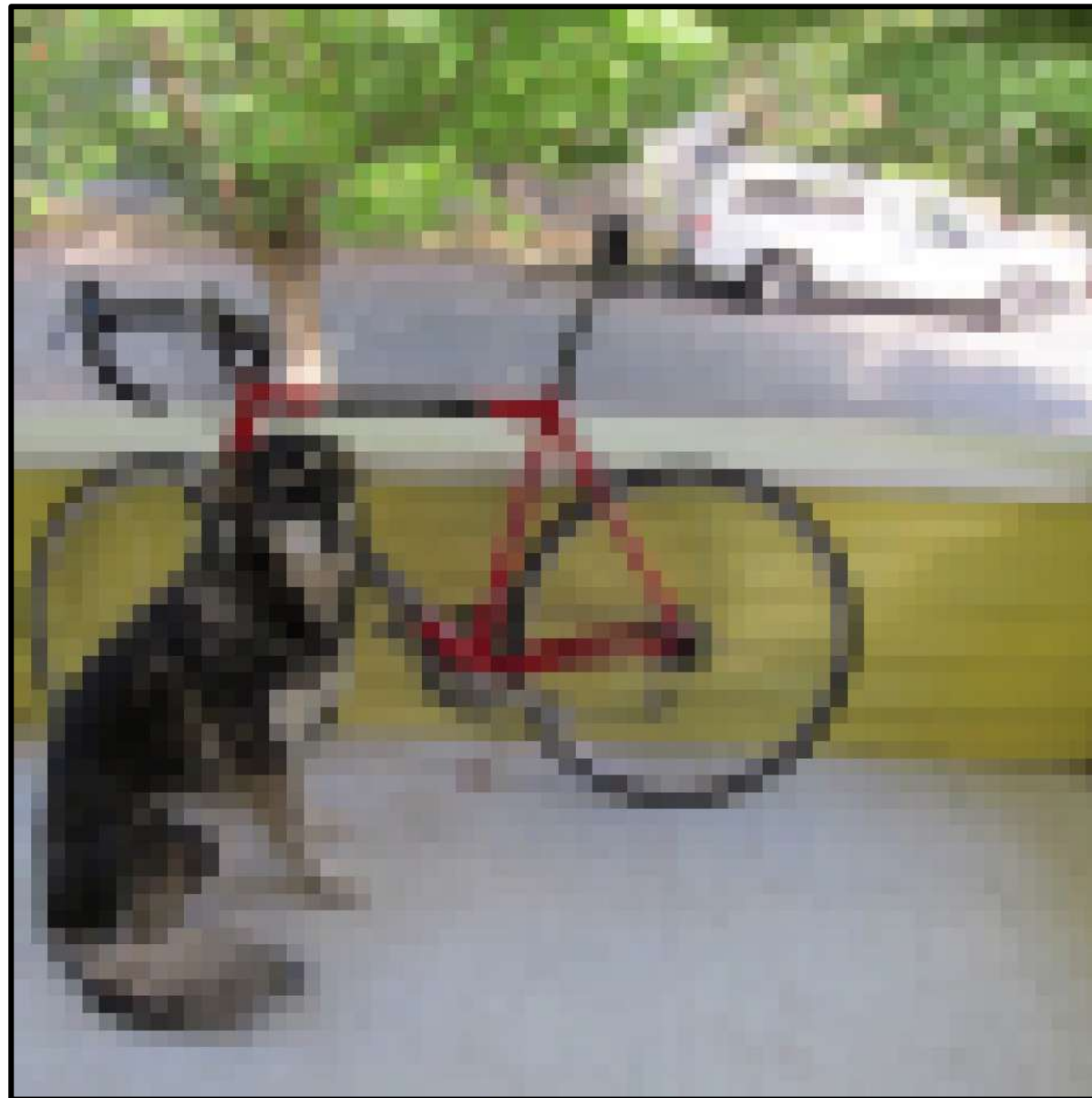


NN



Bilinear

# There is a better way!

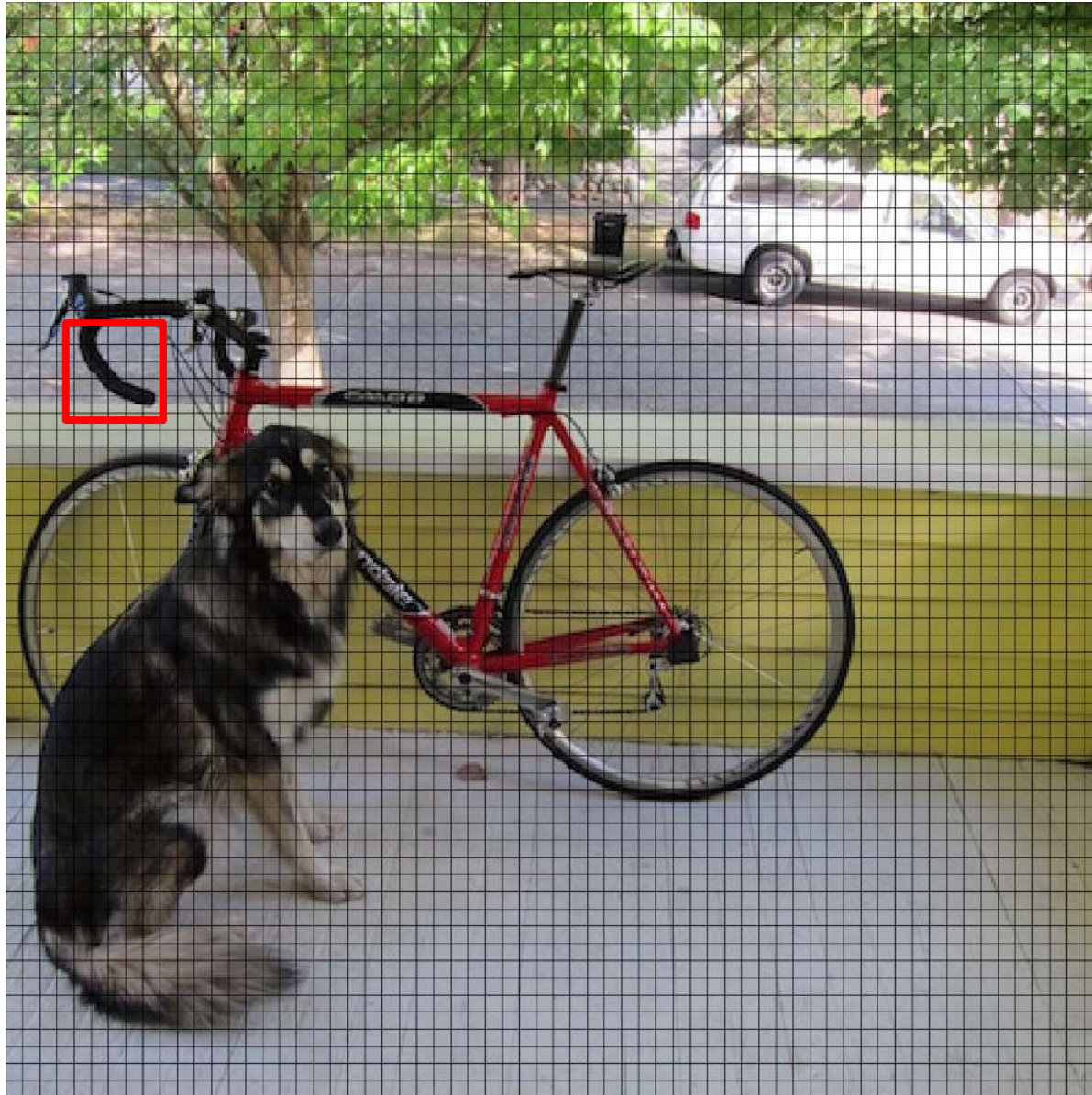


## Look at how much better



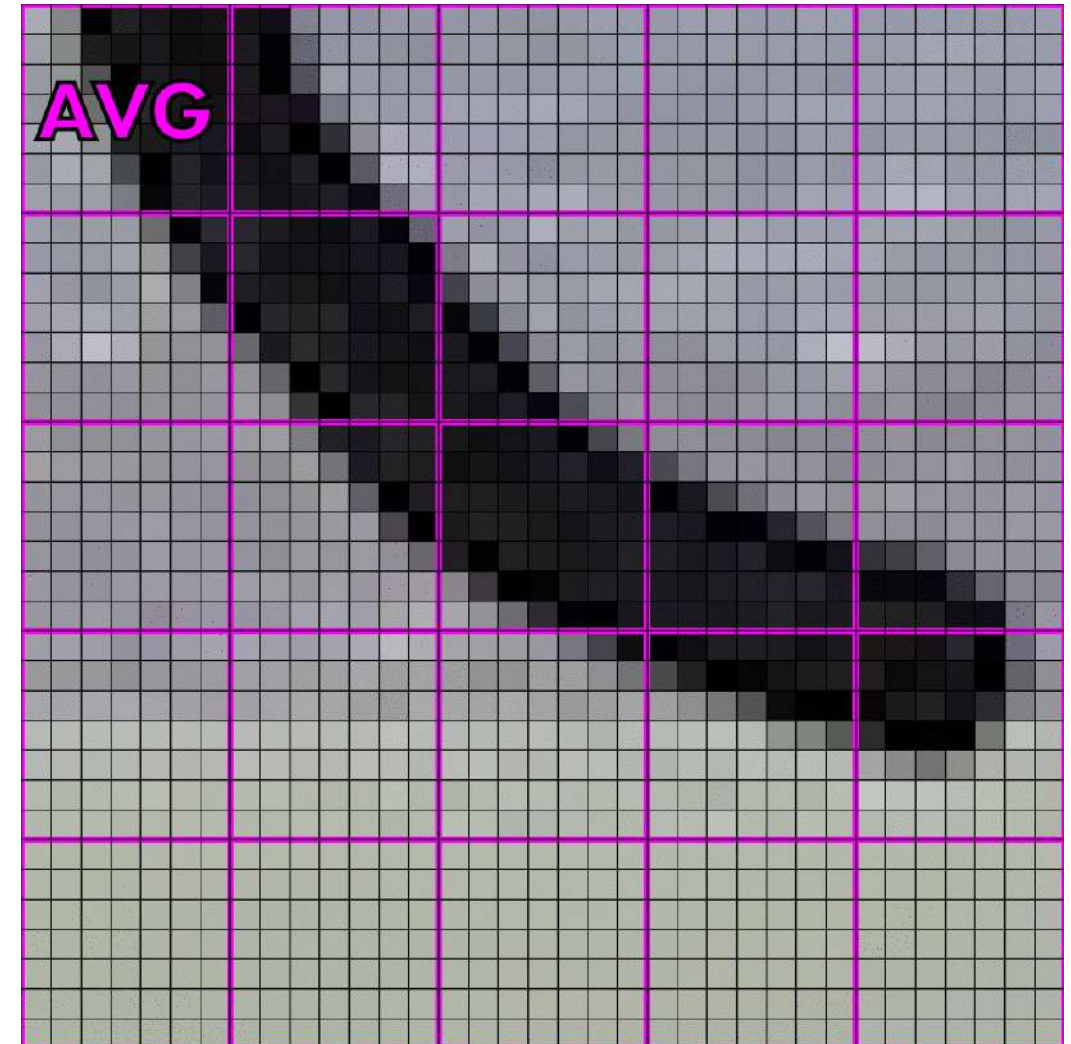
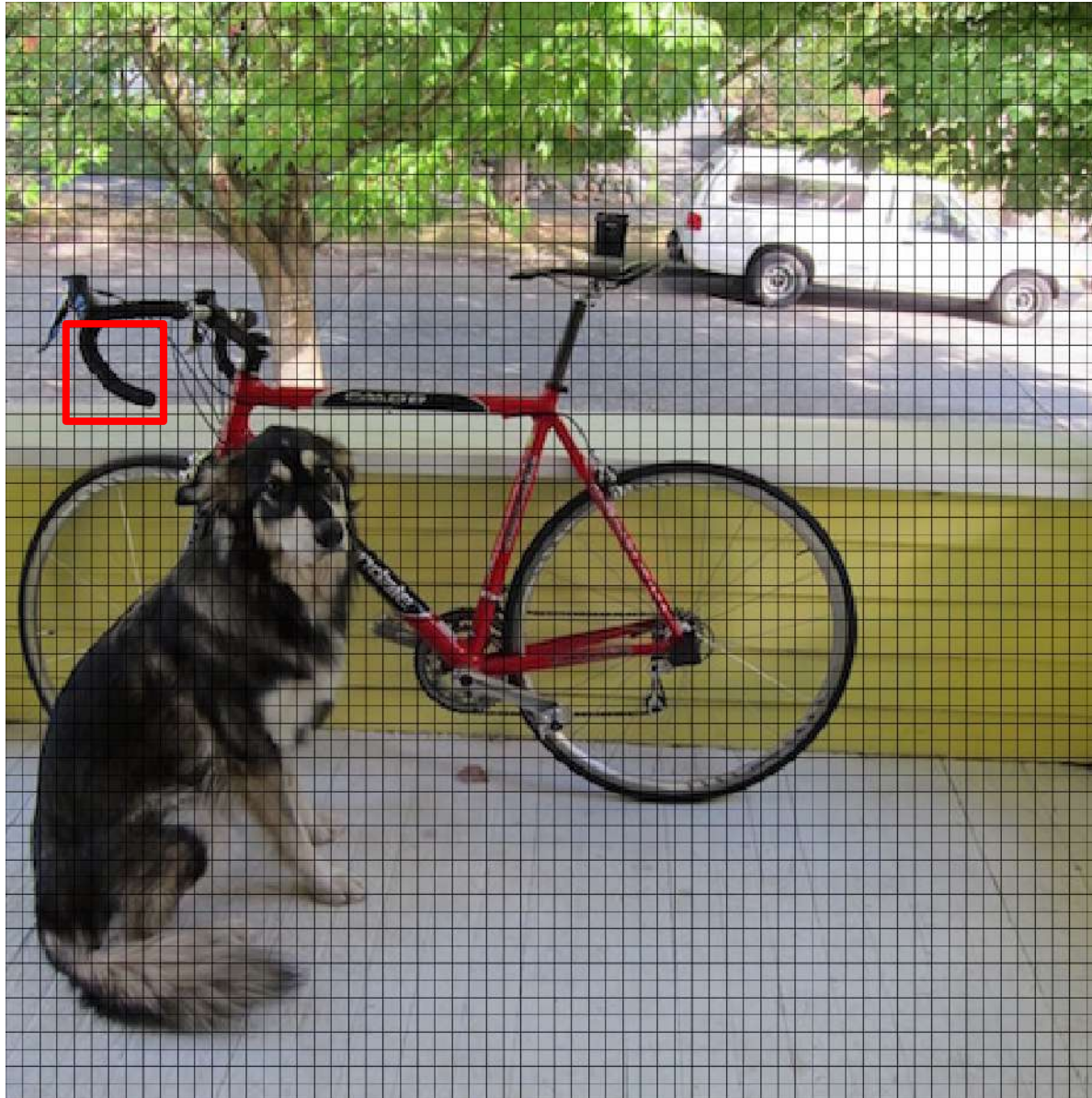


## How?

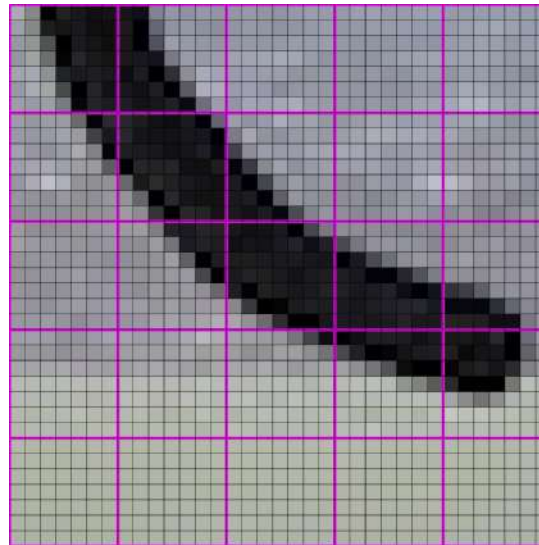




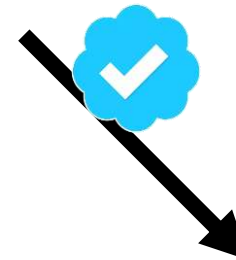
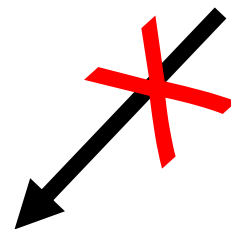
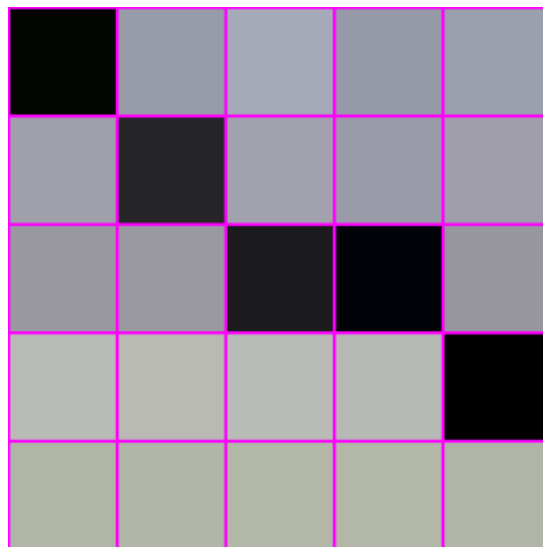
## How? Averaging!



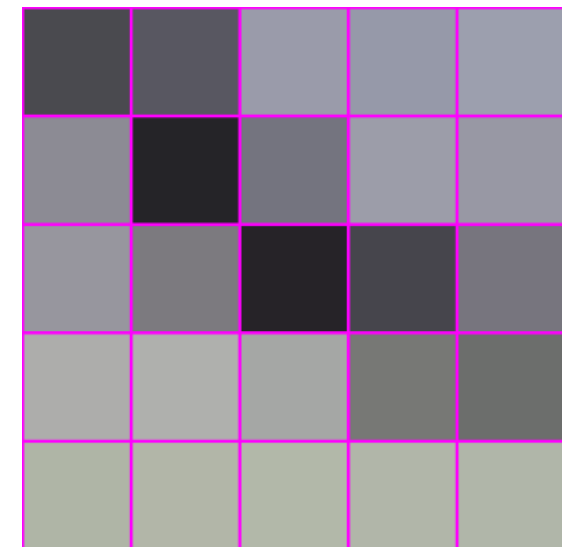
# How? Averaging!



“interpolation”

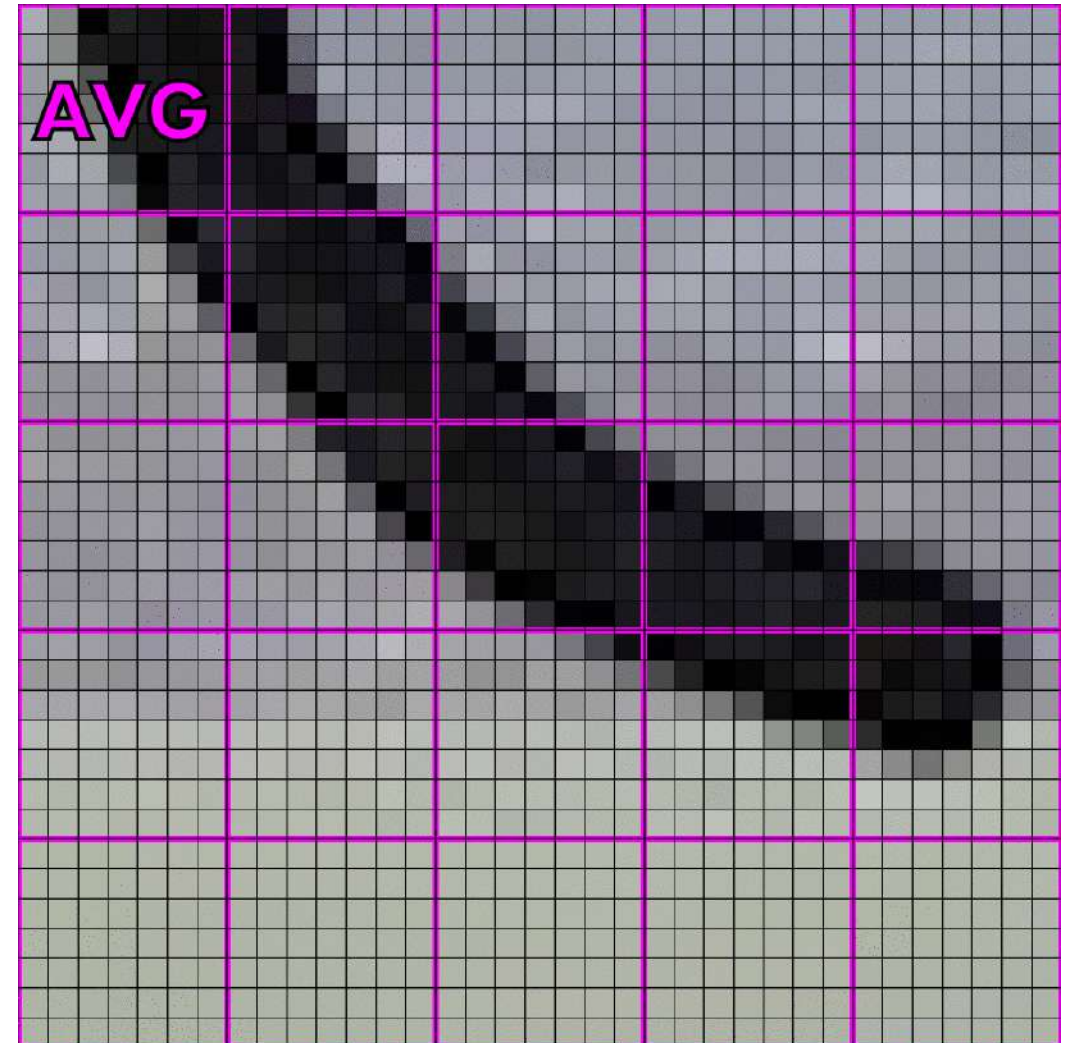


averaging



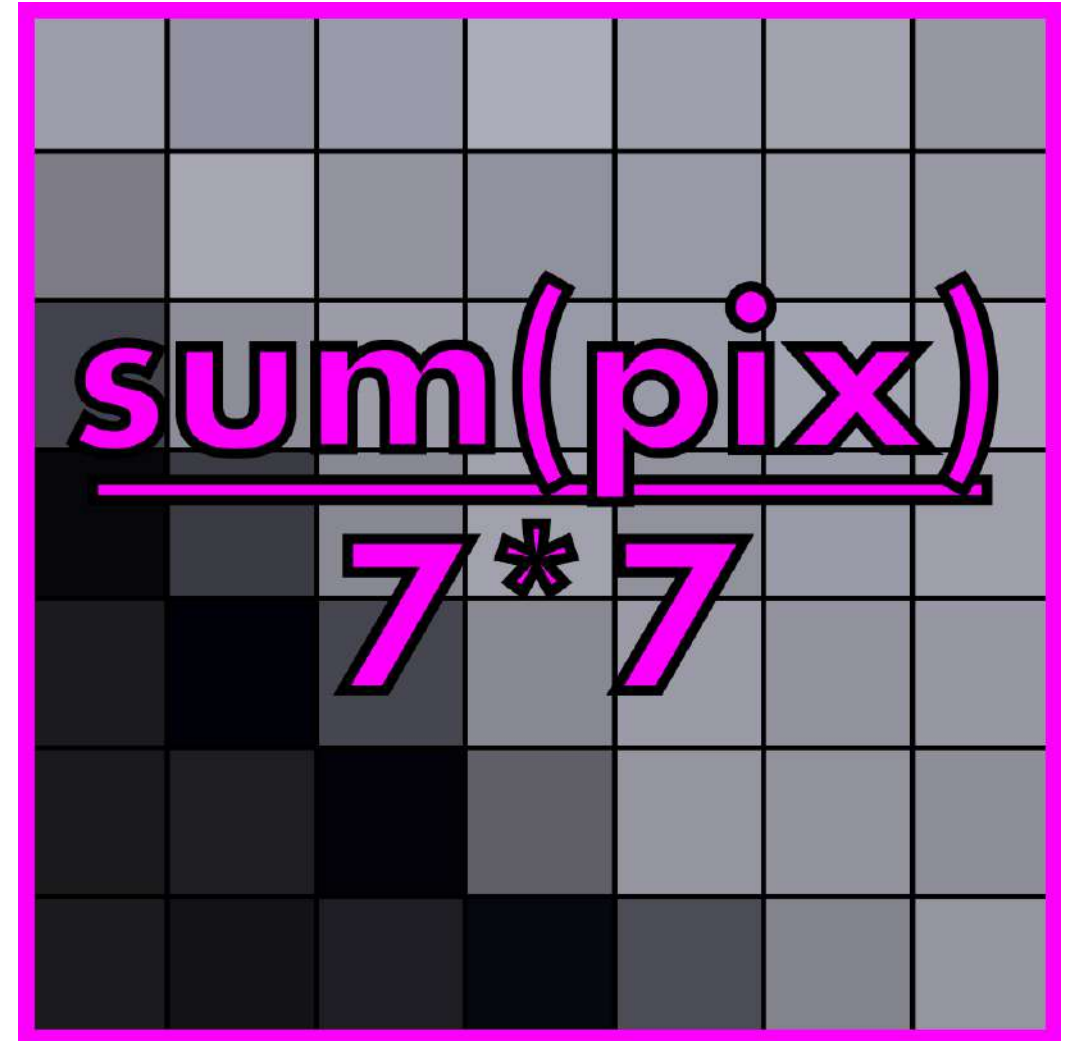


## What is averaging?

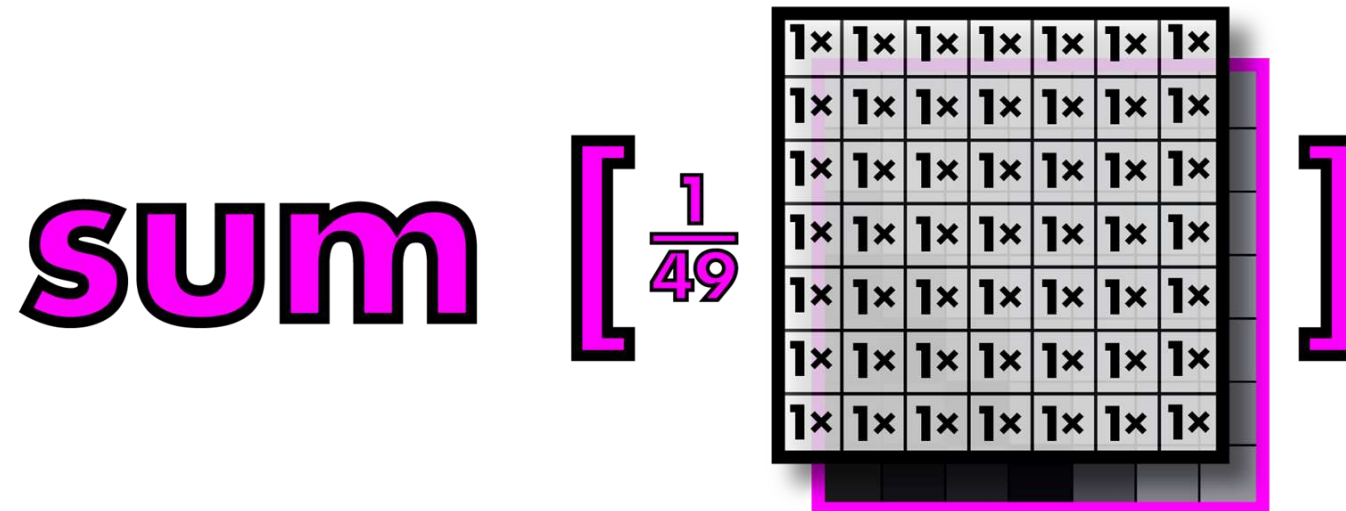




# What is averaging? A weighted sum



# What is averaging? A weighted sum



What are the weights here ?

## Today's Agenda

- Averaging vs Interpolation
- Systems - filters
- Convolution
  - Box Filter
  - Gaussian
  - Cross correlation vs Convolution
- Examples of filters

# Moving average is a filter

*Filter or kernel*

$$\frac{1}{49}$$

1×	1×	1×	1×	1×	1×	1×
1×	1×	1×	1×	1×	1×	1×
1×	1×	1×	1×	1×	1×	1×
1×	1×	1×	1×	1×	1×	1×
1×	1×	1×	1×	1×	1×	1×
1×	1×	1×	1×	1×	1×	1×
1×	1×	1×	1×	1×	1×	1×



# Filtering

- Filtering
  - Forming a new image whose pixel values are transformed from original pixel values
- Goal is to extract useful information from images, or transform images into another domain where we can modify/enhance image properties
  - Features (edges, corners, blobs...)
  - Applications: super-resolution (resizing); in-painting; de-noising;

[Slide by Niebles]

# Applications

### De-noising

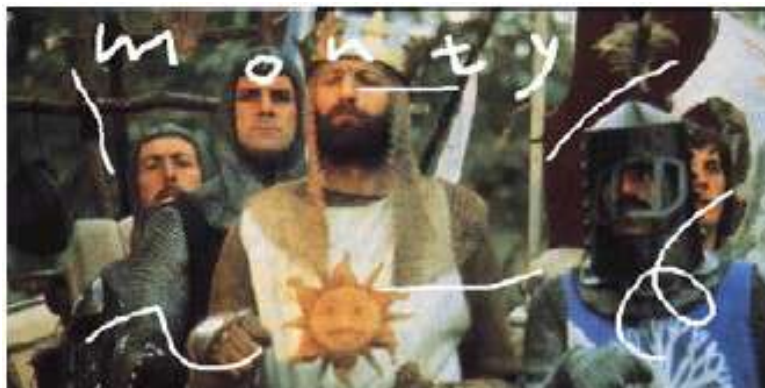


Salt and pepper noise

### Super-resolution



### In-painting



Bertamio et al

[Slide by Niebles]

# Systems

- We define a **system** as a unit that converts an input function  $f[x,y]$  into an output (or response) function  $g[x,y]$ , where  $(x,y)$  are the independent variables.
  - In the case of images,  $(x,y)$  represents the **spatial position in the image**.

[Slide by Niebles]

# Moving average - example

$$F[x, y]$$

0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	0	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	0	0	0	0	0	0	0
0	0	90	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0

$$G[x, y]$$


[Slide by Seitz]





# Moving average - example

$$F[x, y]$$

0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	0	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	0	0	0	0	0	0	0
0	0	90	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0

$$G[x, y]$$

	0	10							

[Slide by Seitz]



# Moving average - example

$$F[x, y]$$

0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	0	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	0	0	0	0	0	0	0
0	0	90	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0

$$G[x, y]$$

	0	10	20						

[Slide by Seitz]

# Moving average - example

$$F[x, y]$$

0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	0	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	0	0	0	0	0	0	0
0	0	90	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0

$$G[x, y]$$

	0	10	20	30					

[Slide by Seitz]

# Moving average - example

$$F[x, y]$$

0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	0	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	0	0	0	0	0	0	0
0	0	90	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0

$$G[x, y]$$

	0	10	20	30	30				

[Slide by Seitz]



# Moving average - example

$$F[x, y]$$

0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	0	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	0	0	0	0	0	0	0
0	0	90	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0

$$G[x, y]$$

	0	10	20	30	30	30	20	10	
	0	20	40	60	60	60	40	20	
	0	30	60	90	90	90	60	30	
	0	30	50	80	80	90	60	30	
	0	30	50	80	80	90	60	30	
	0	20	30	50	50	60	40	20	
	10	20	30	30	30	30	20	10	
	10	10	10	0	0	0	0	0	

[Slide by Seitz]

# Properties of systems

- Amplitude properties

- Additivity  $S[f_i[n, m] + f_j[n, m]] = S[f_i[n, m]] + S[f_j[n, m]]$

- Homogeneity  $S[\alpha f_i[n, m]] = \alpha S[f_i[n, m]]$

- Superposition  $S[\alpha f_i[n, m] + \beta f_j[n, m]] = \alpha S[f_i[n, m]] + \beta S[f_j[n, m]]$

- Stability  $|f[n, m]| \leq k \implies |g[n, m]| \leq ck$

- Invertibility  $S^{-1}[S[f_i[n, m]]] = f_i[n, m]$

[Slide by Niebles]

# Properties of systems

- Spatial properties

- Causality

for  $n < n_0, m < m_0$ , if  $f[n, m] = 0 \implies g[n, m] = 0$

- Shift invariance

$f[n - n_0, m - m_0] \xrightarrow{\mathcal{S}} g[n - n_0, m - m_0]$

[Slide by Niebles]

# Linear Systems - filters

- Linear filtering
  - Form a new image whose pixels are a weighted sum of original pixel values
  - Use the same set of weights at each point
- $S$  is a linear system (function) iff  $S$  satisfies

$$S[\alpha f_i[n, m] + \beta f_j[h, m]] = \alpha S[f_i[n, m]] + \beta S[f_j[h, m]]$$

superposition property

[Slide by Niebles]



# Linear Shift Invariant Systems

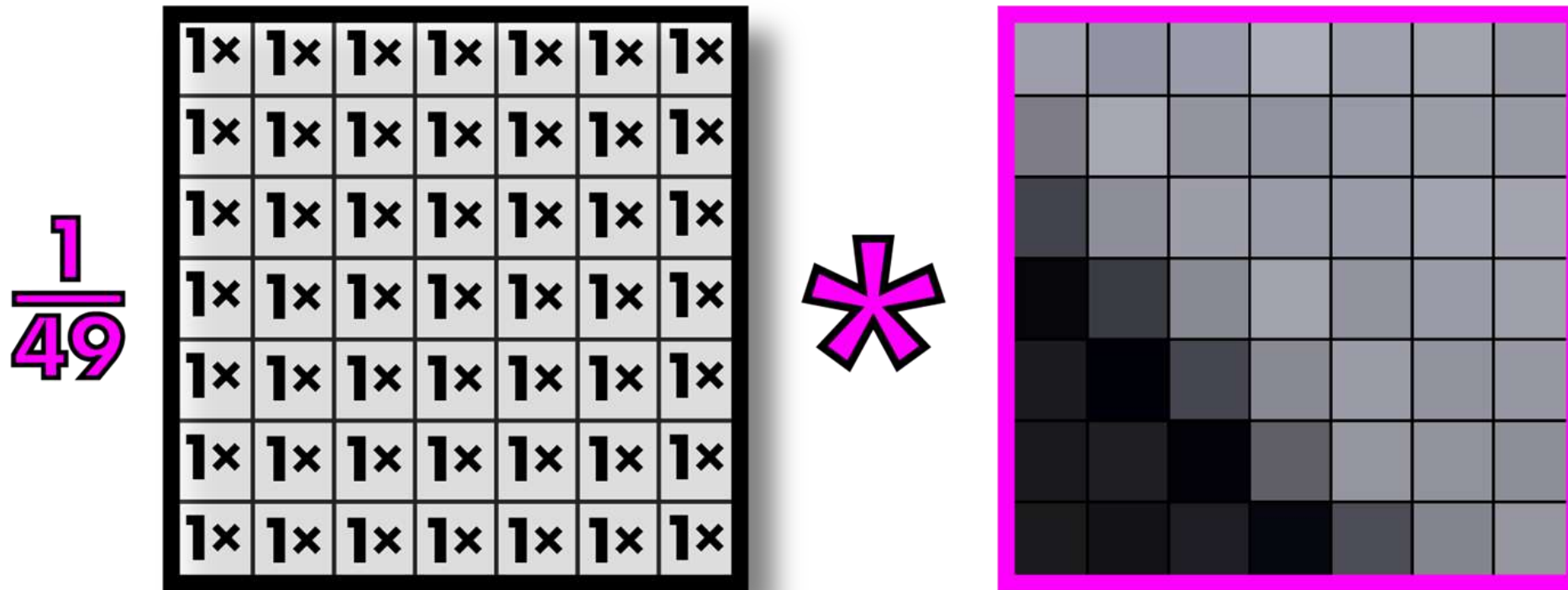
- We call systems which satisfy the superposition and shift-invariant property *Linear Shift Invariant Systems* (LSI)
- Not all filters are LSI
  - Is thresholding linear? 
$$g[n, m] = \begin{cases} 1, & f[n, m] > 100 \\ 0, & \text{otherwise.} \end{cases}$$
  - Consider:
    - $f_1[n, m] + f_2[n, m] > T$
    - $f_1[n, m] < T$
    - $f_2[n, m] < T$
- LSI systems can be described by the **convolution** operation

# Today's Agenda

- Averaging vs Interpolation
- Systems - filters
- Convolution
  - Box Filter
  - Gaussian
  - Cross correlation vs Convolution
- Examples of filters

Call this operation “*convolution*”

*Filter or kernel*



Note: multiplying an image section by a filter is actually called “correlation” and convolution involves inverting the filter first

# Convolutions on larger images

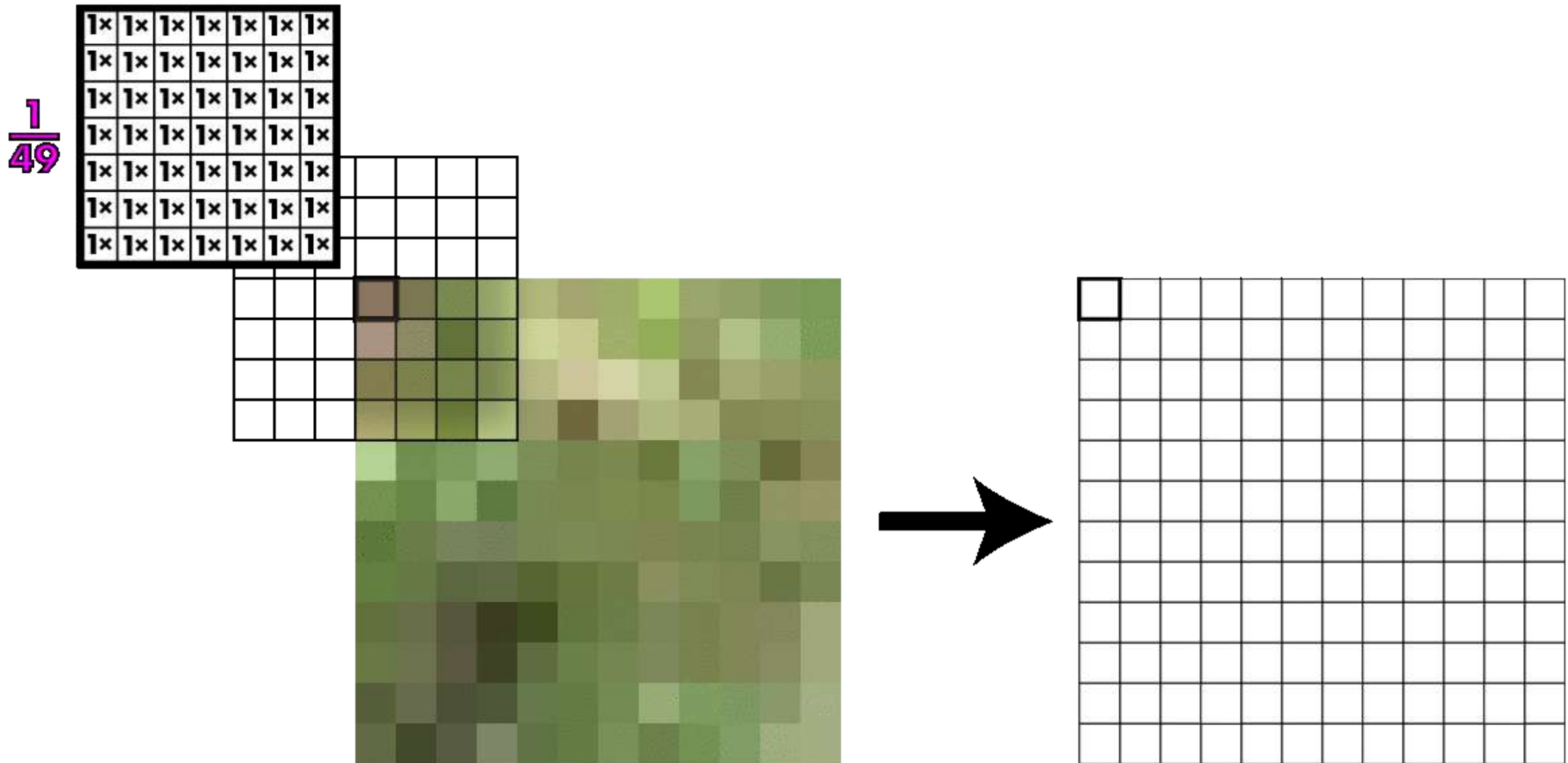
$\frac{1}{49}$

1x	1x	1x	1x	1x	1x	1x
1x	1x	1x	1x	1x	1x	1x
1x	1x	1x	1x	1x	1x	1x
1x	1x	1x	1x	1x	1x	1x
1x	1x	1x	1x	1x	1x	1x
1x	1x	1x	1x	1x	1x	1x
1x	1x	1x	1x	1x	1x	1x

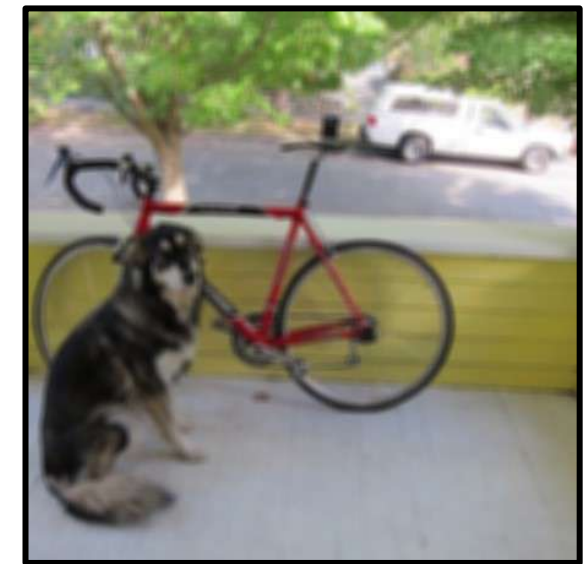
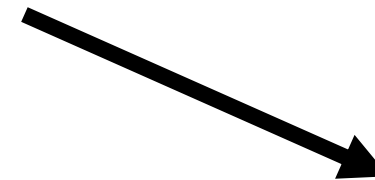
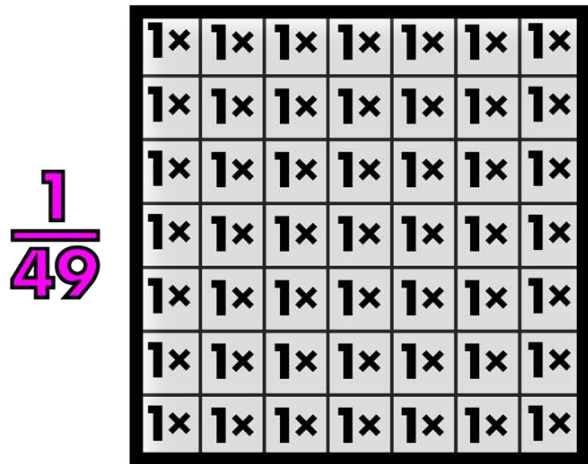




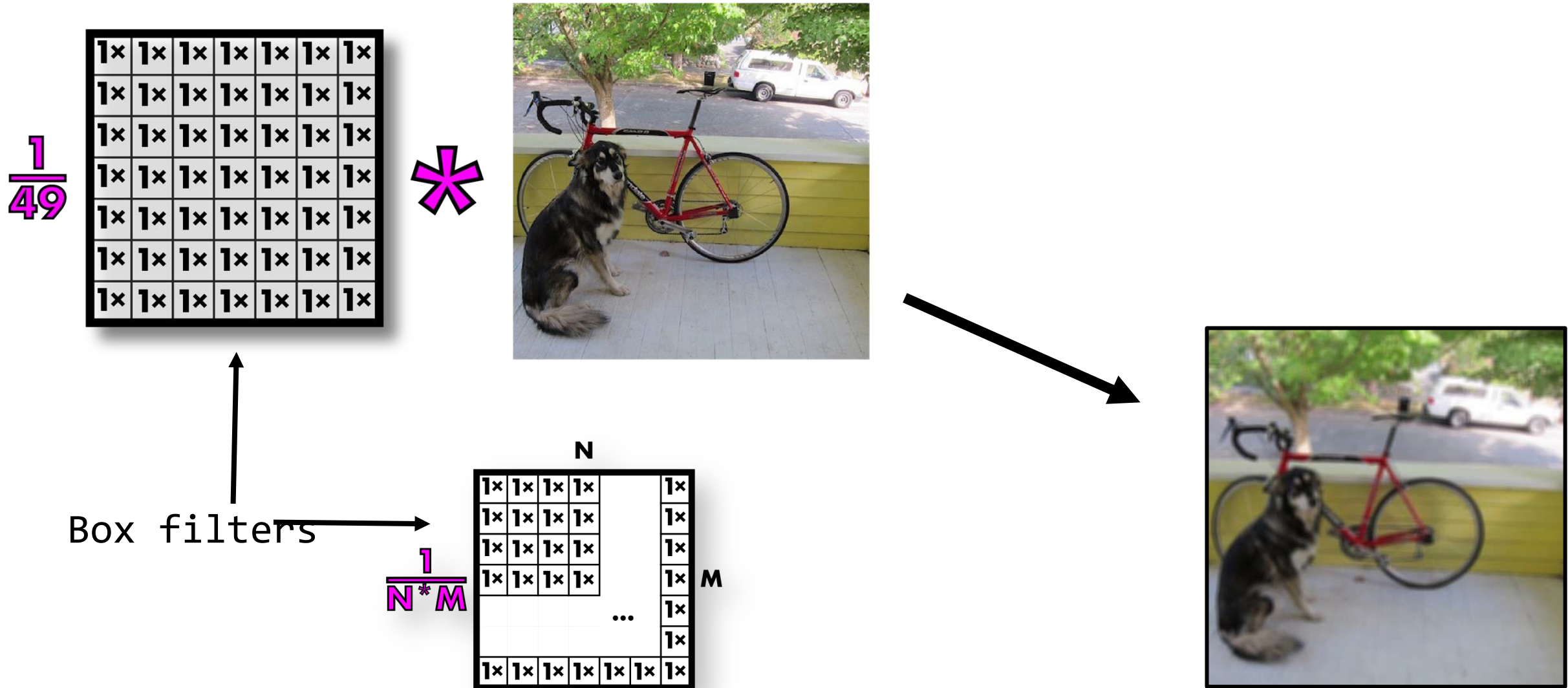
## Kernel slides across image



# Convolutions on larger images

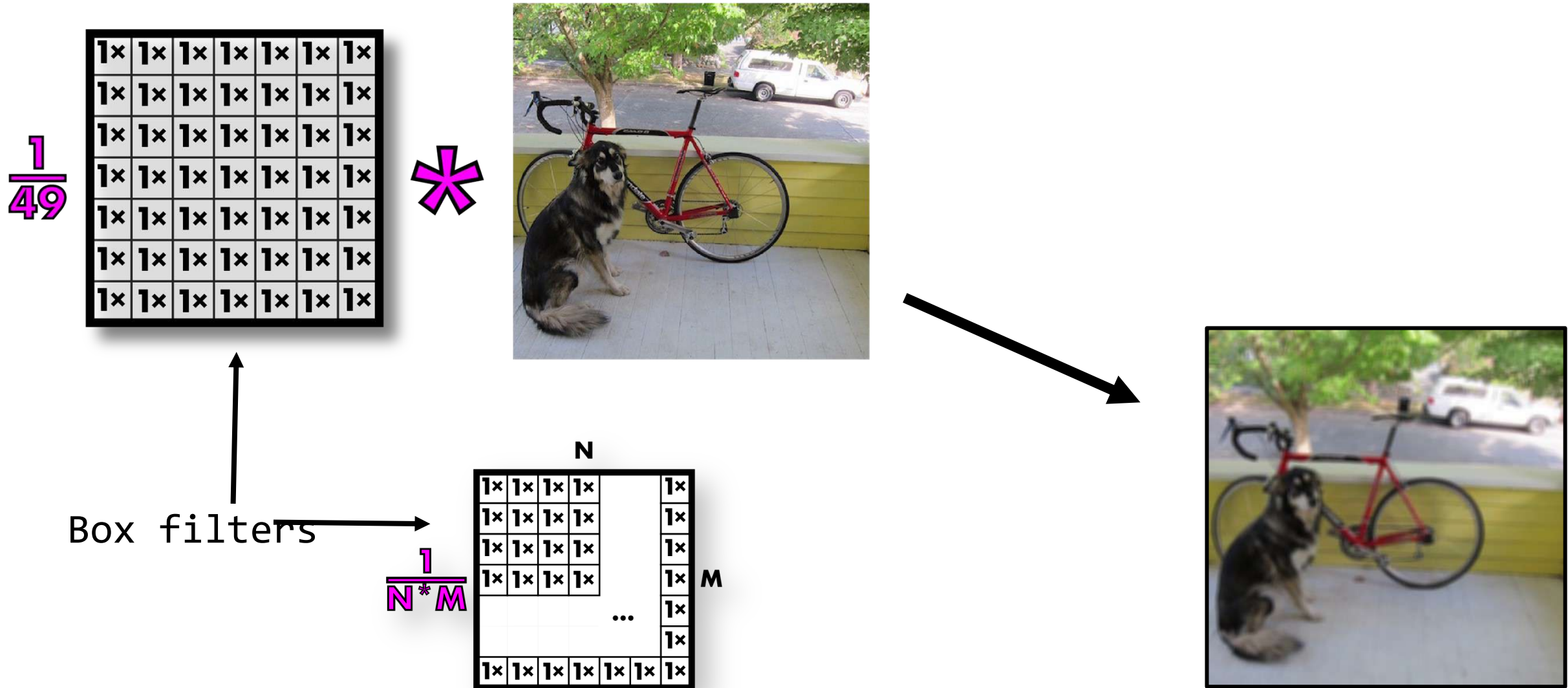


# This is called box filter



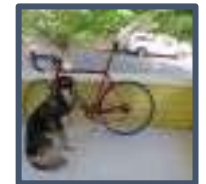


# Box filters smooth image

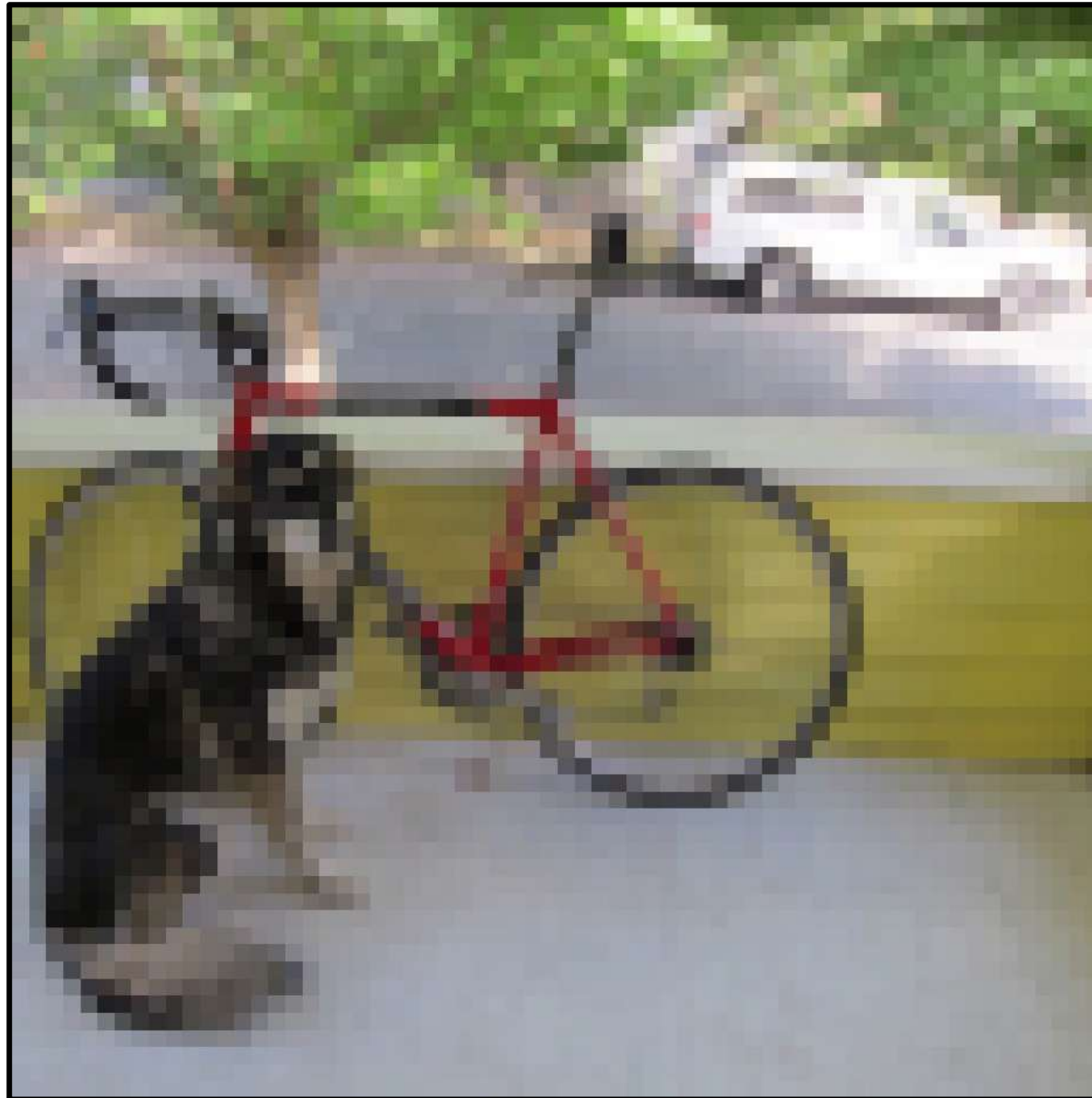




## Now we resize our smoothed image

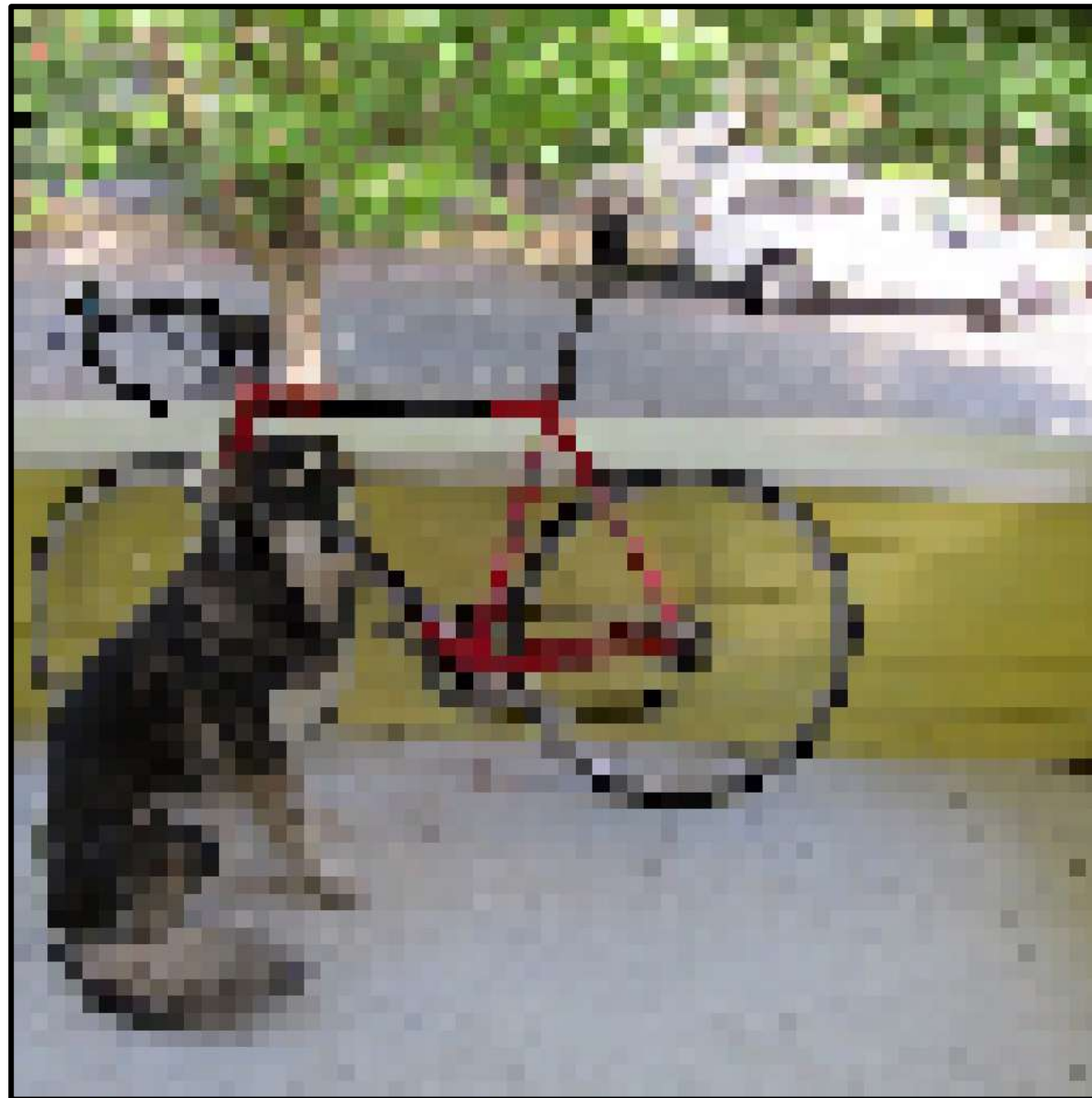


## So much better!



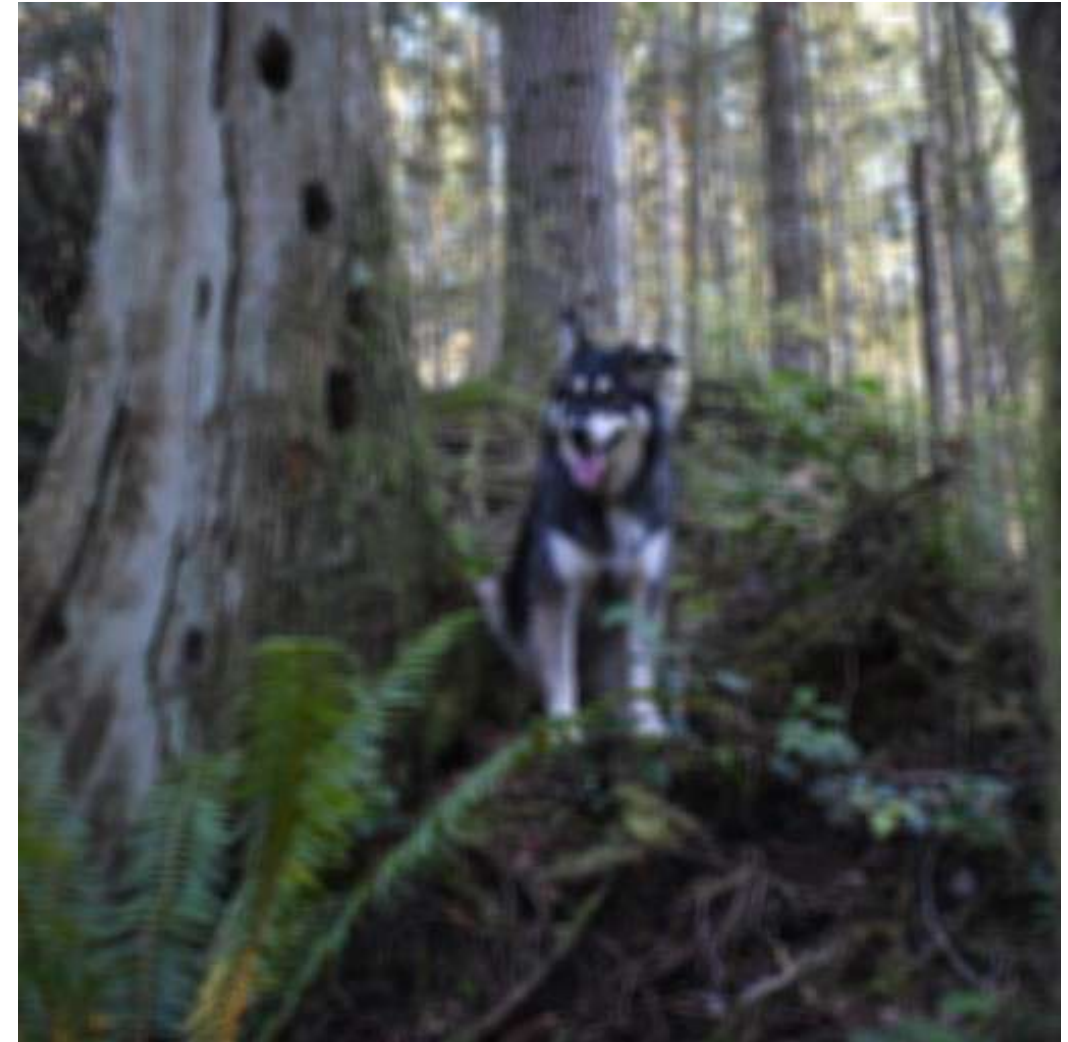


## Compare to interpolation



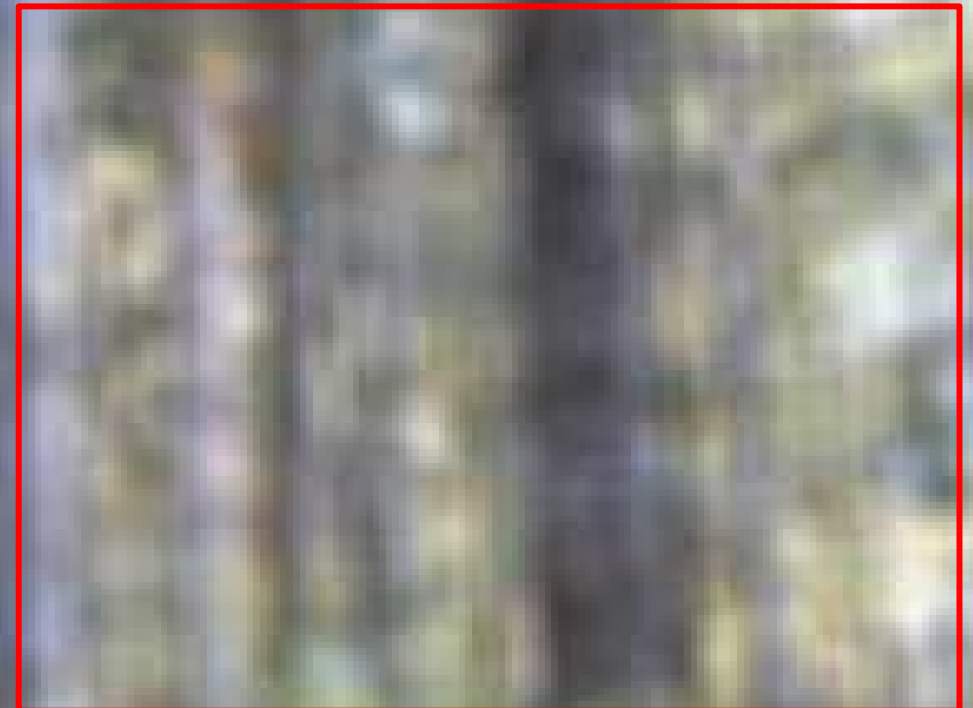


## Box filters have artifacts



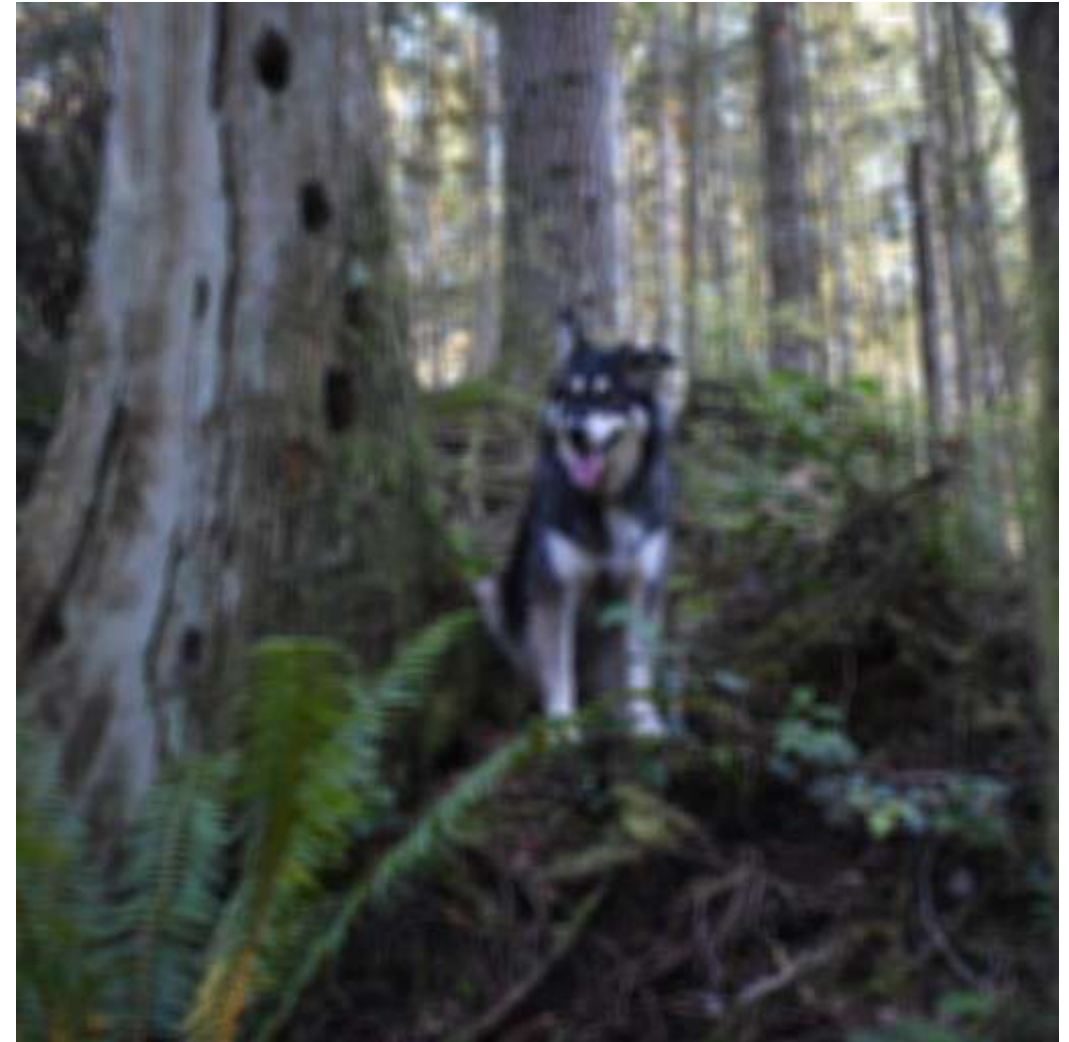
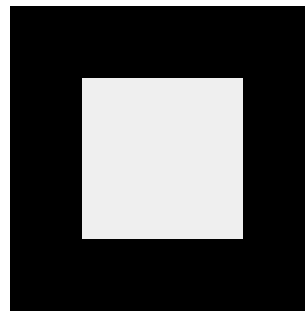


## Box filters have artifacts

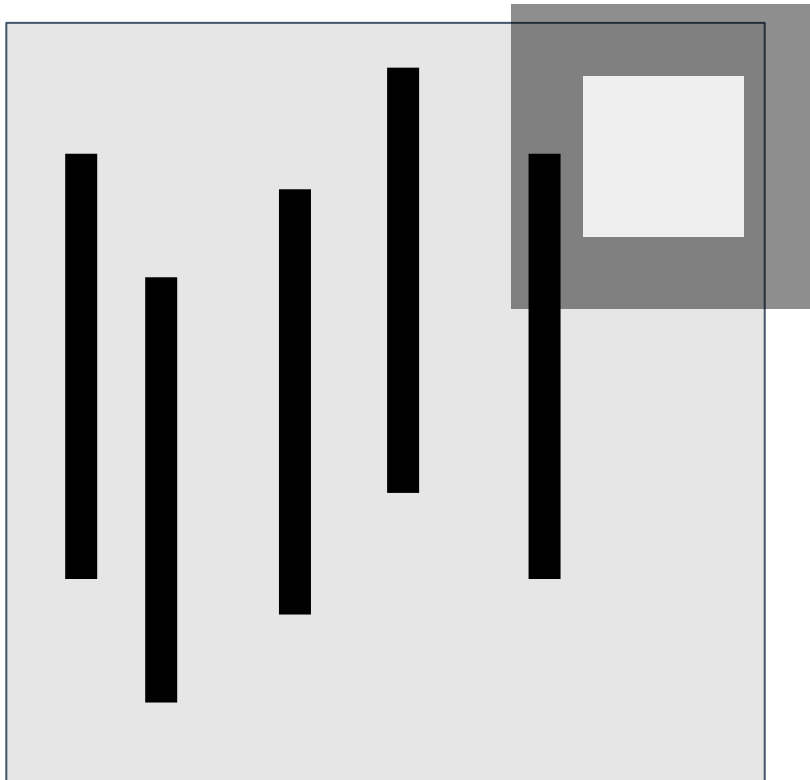




## Box filters: vertical + horizontal streaking

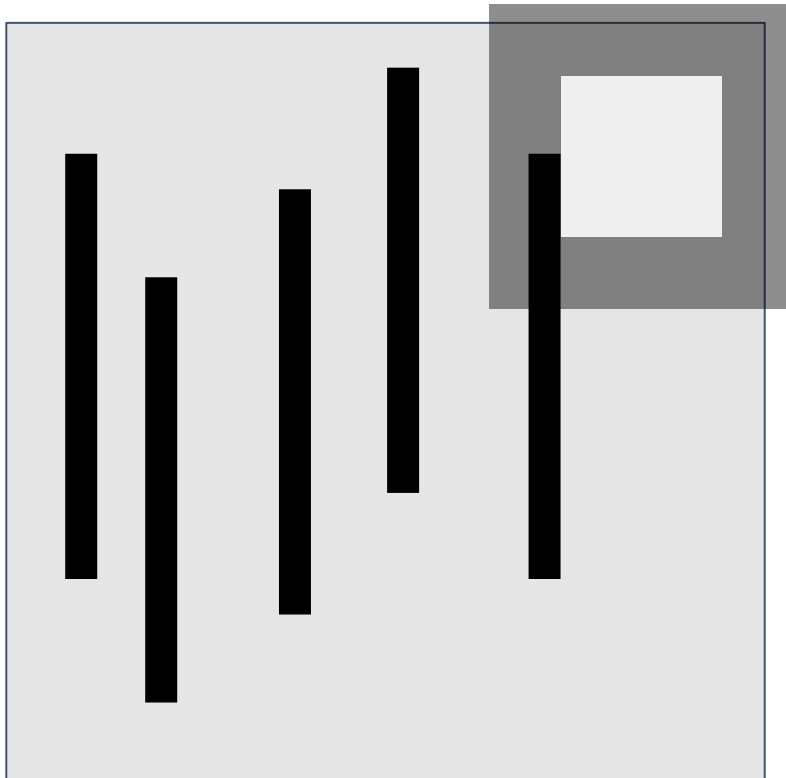


## Box filters: vertical + horizontal streaking



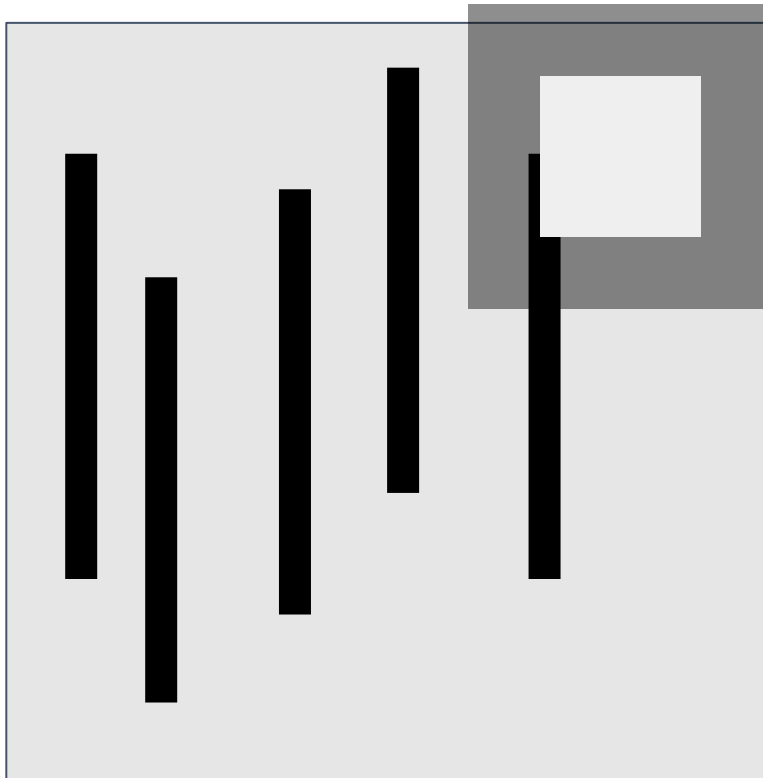


## Box filters: vertical + horizontal streaking

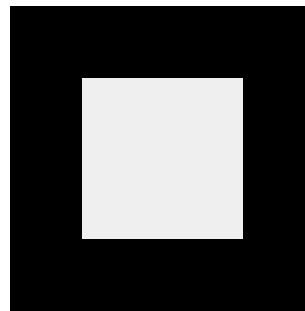




## Box filters: vertical + horizontal streaking



## We want a smoothly weighted kernel

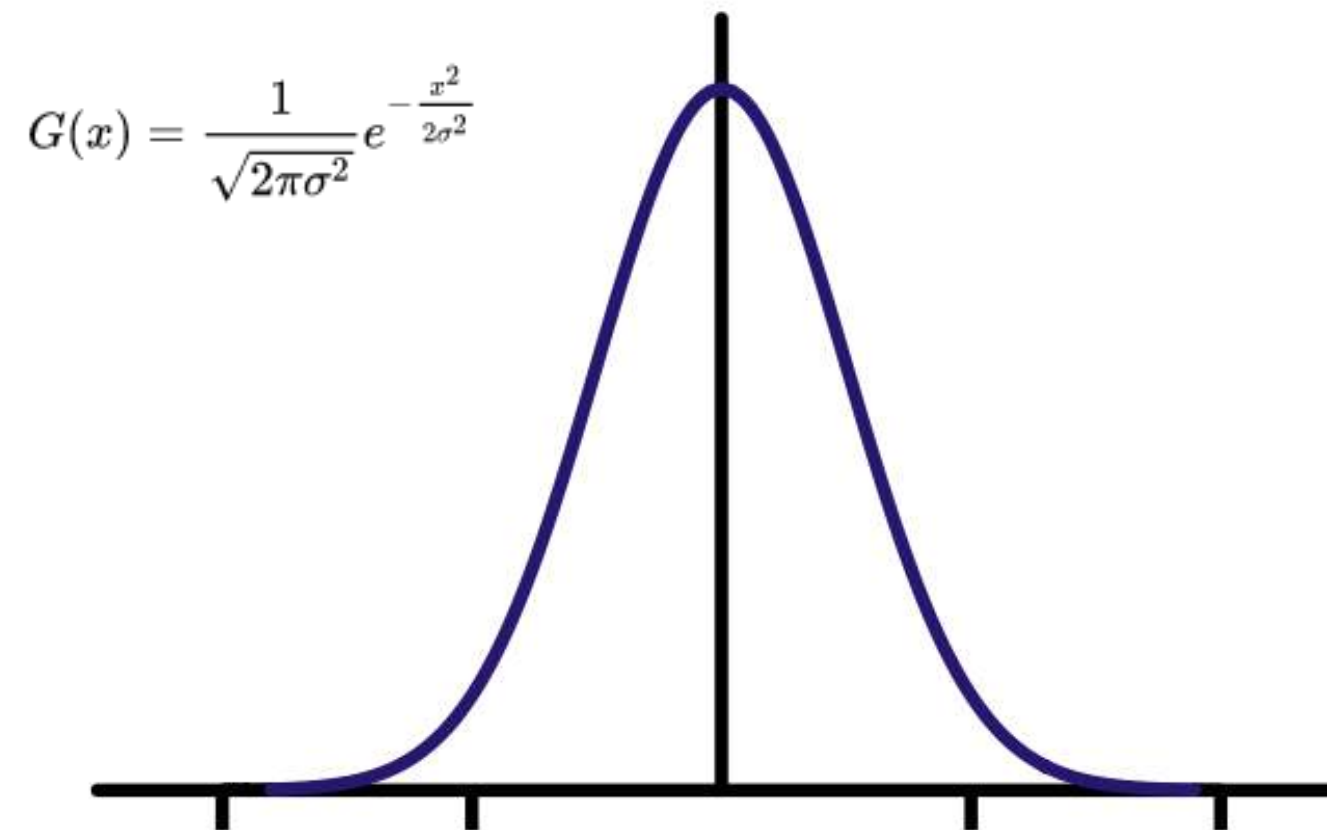


## Today's Agenda

- Averaging vs Interpolation
- Systems - filters
- Convolution
  - Box Filter
  - Gaussian
  - Cross correlation vs Convolution
- Examples of filters

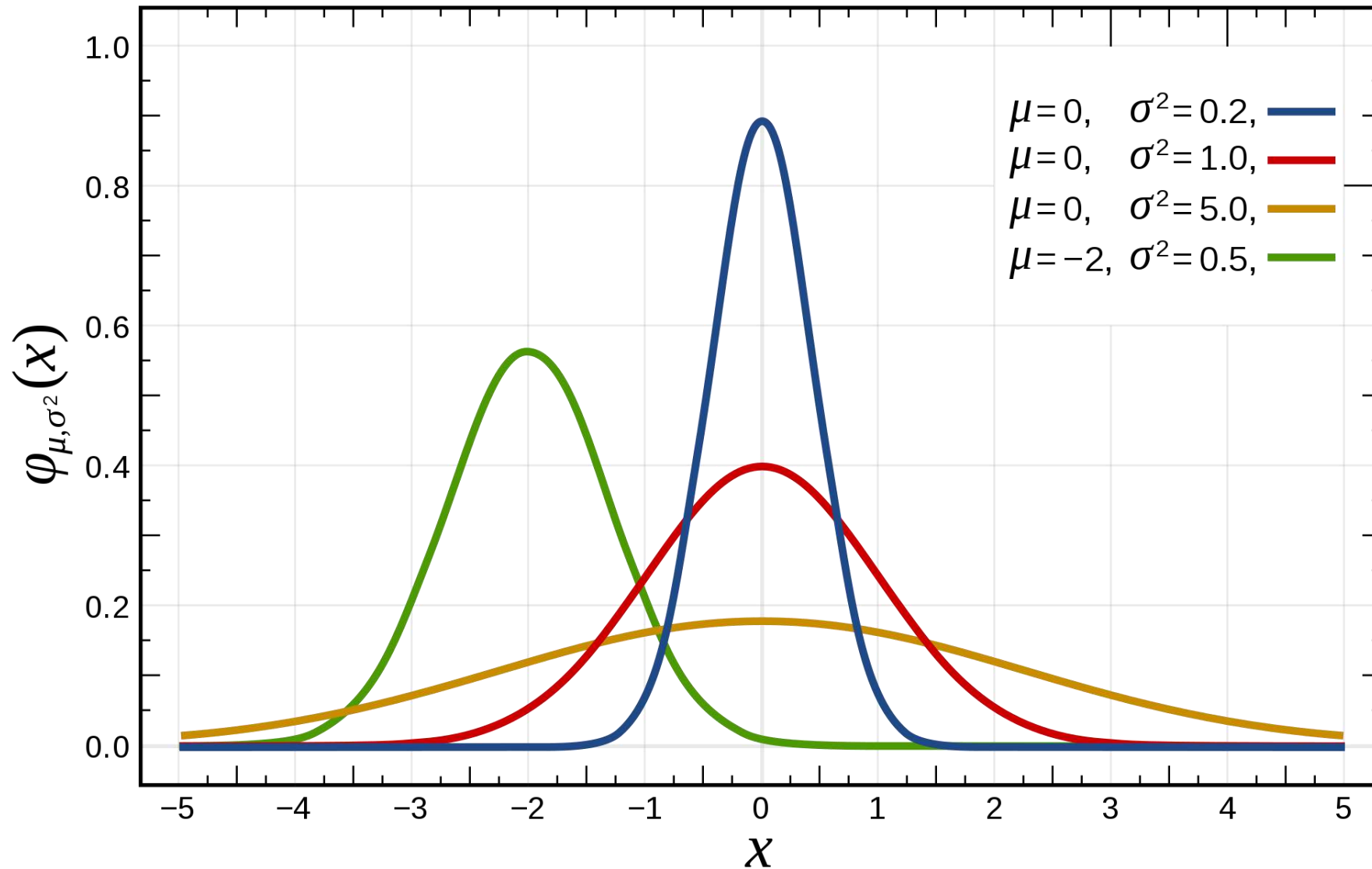


# Gaussians



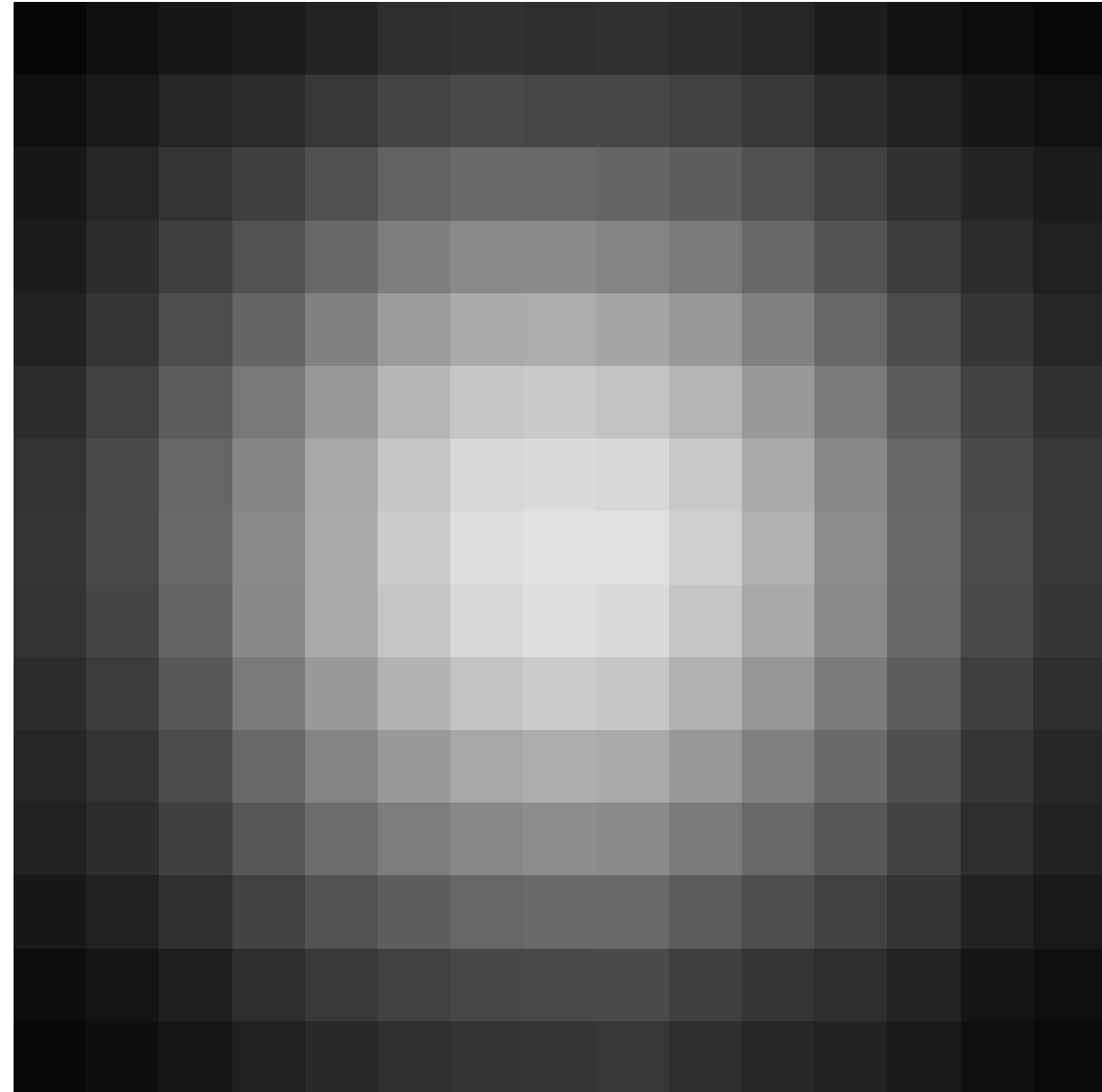
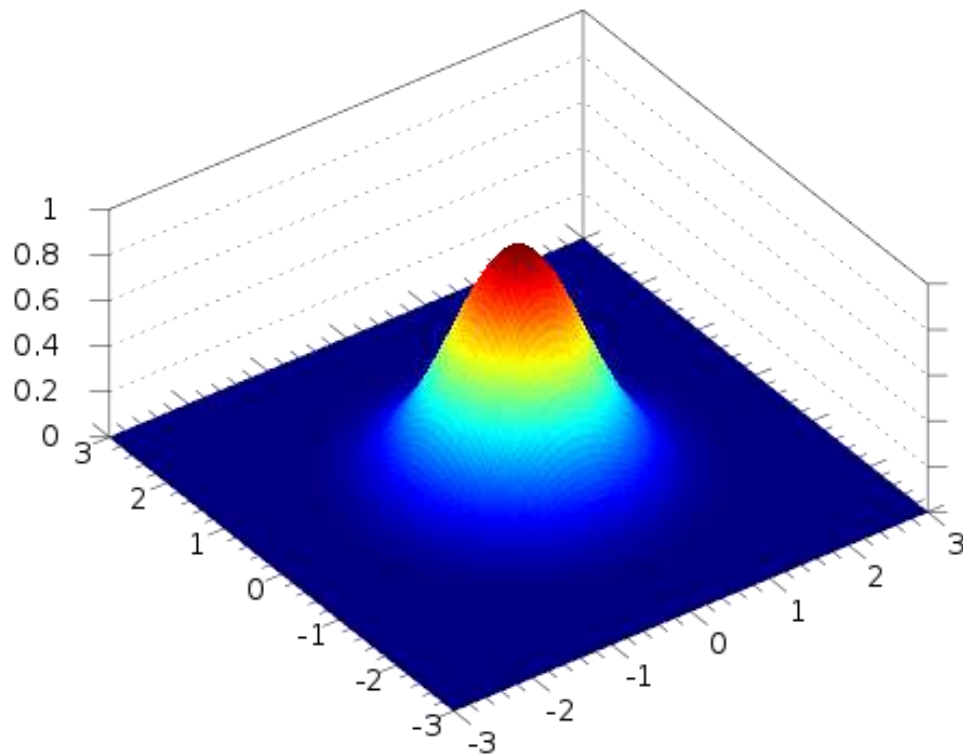


# Gaussians – how $\sigma$ affects the shape



## 2D Gaussian

$$g(x, y) = \frac{1}{2\pi\sigma^2} \cdot e^{-\frac{x^2+y^2}{2\sigma^2}}$$



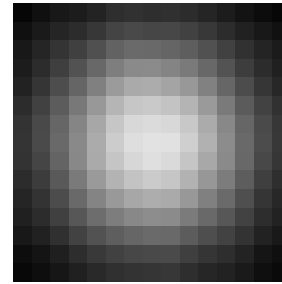
# Example 7x7 Gaussian

		0.000	0.000	0.001	0.001	0.001	0.000	0.000	
		0.000	0.002	0.012	0.020	0.012	0.002	0.000	
		0.001	0.012	0.068	0.109	0.068	0.012	0.001	
		0.001	0.020	0.109	0.172	0.109	0.020	0.001	
		0.001	0.012	0.068	0.109	0.068	0.012	0.001	
		0.000	0.002	0.012	0.020	0.012	0.002	0.000	
		0.000	0.000	0.001	0.001	0.001	0.000	0.000	

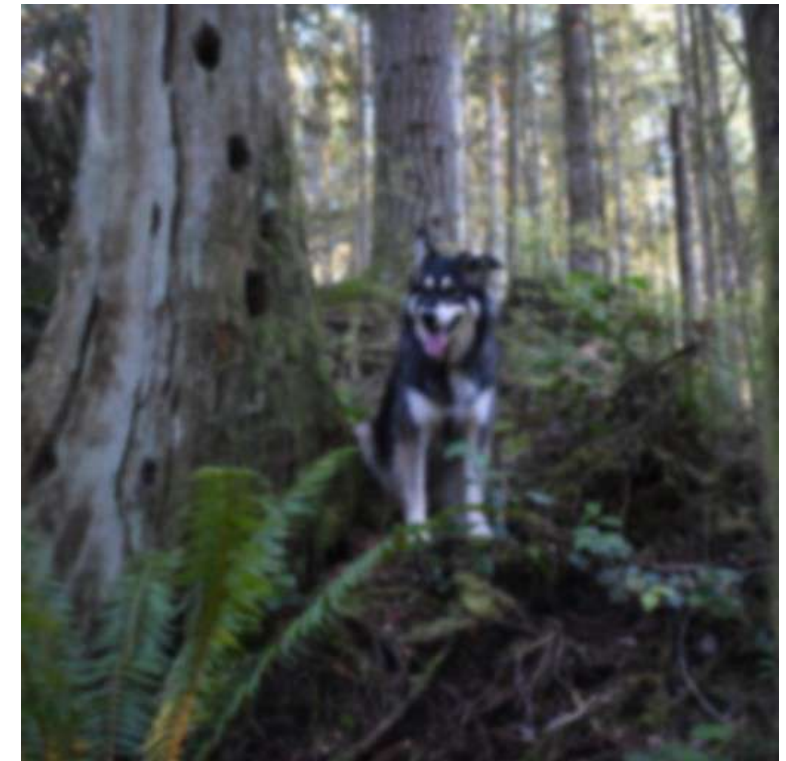
# Better smoothing with Gaussians



\*

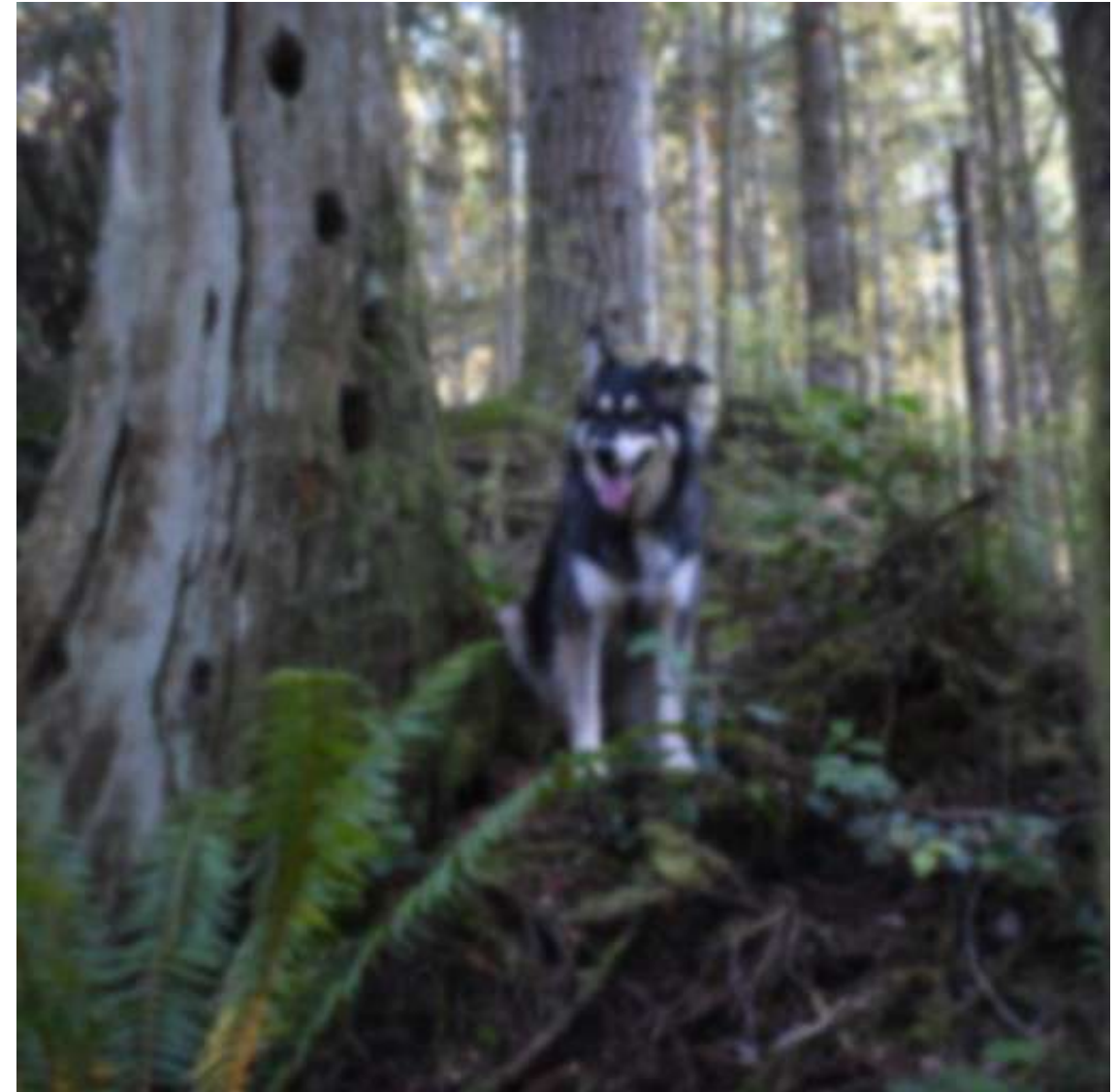
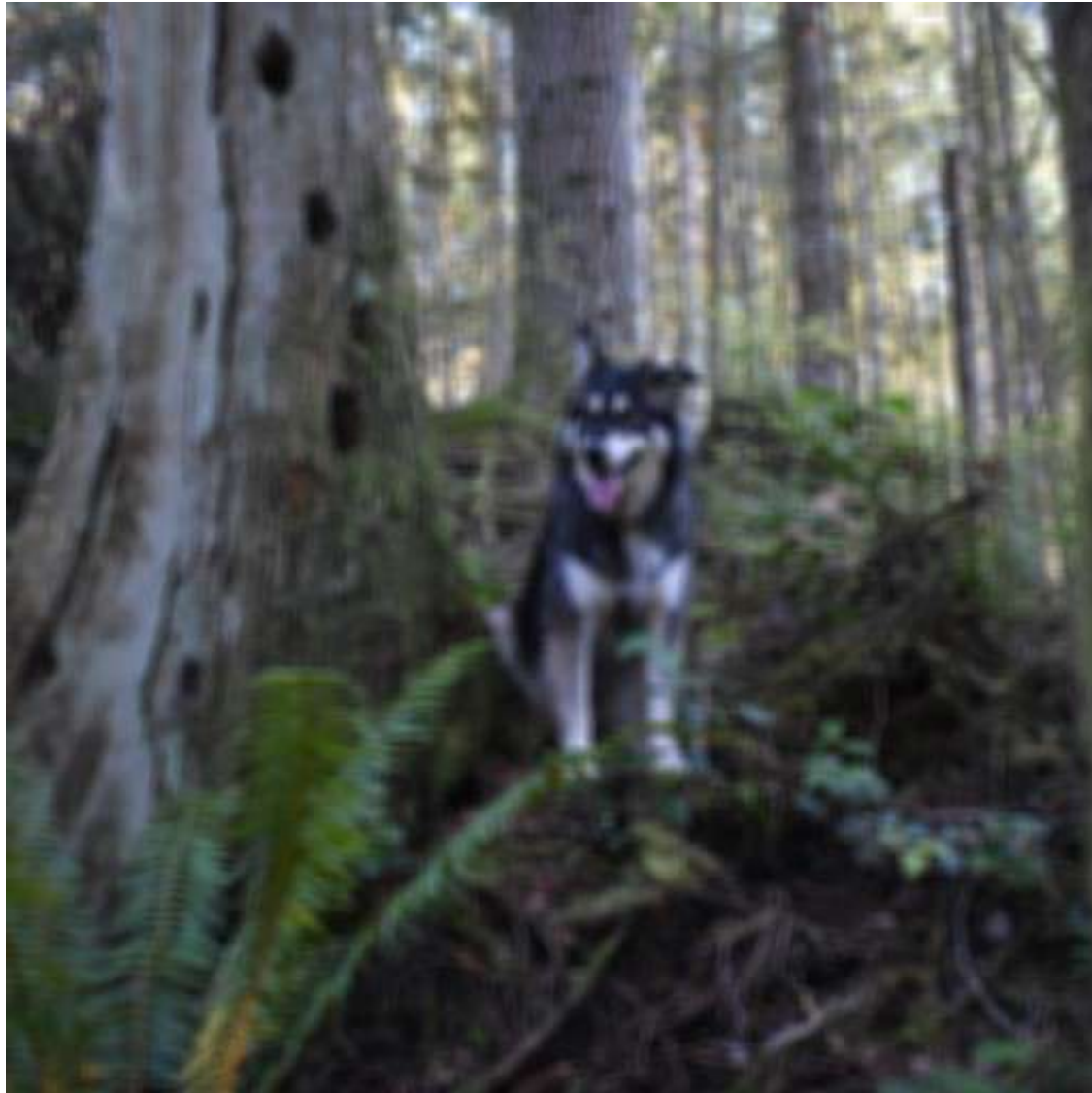


=

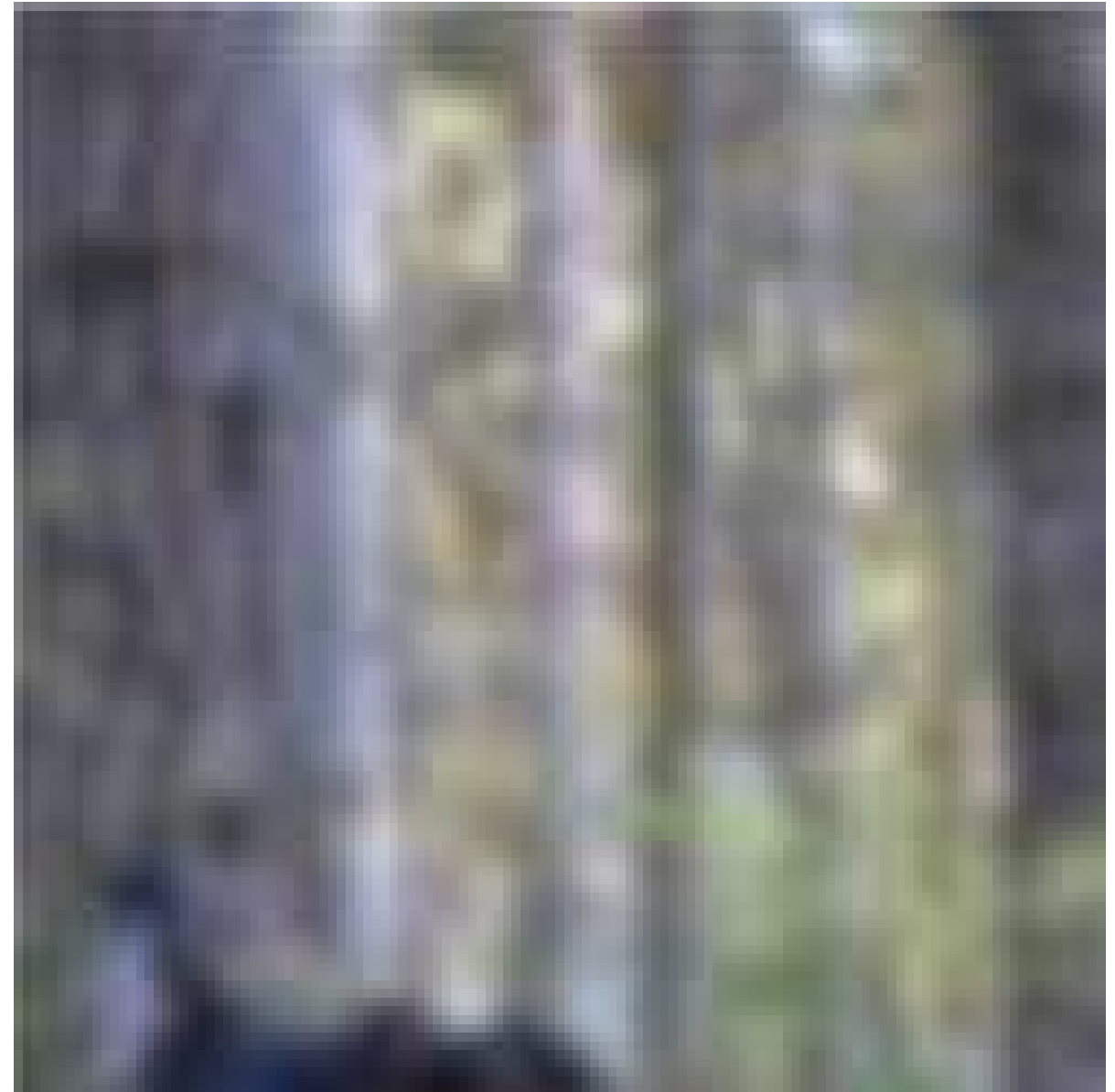
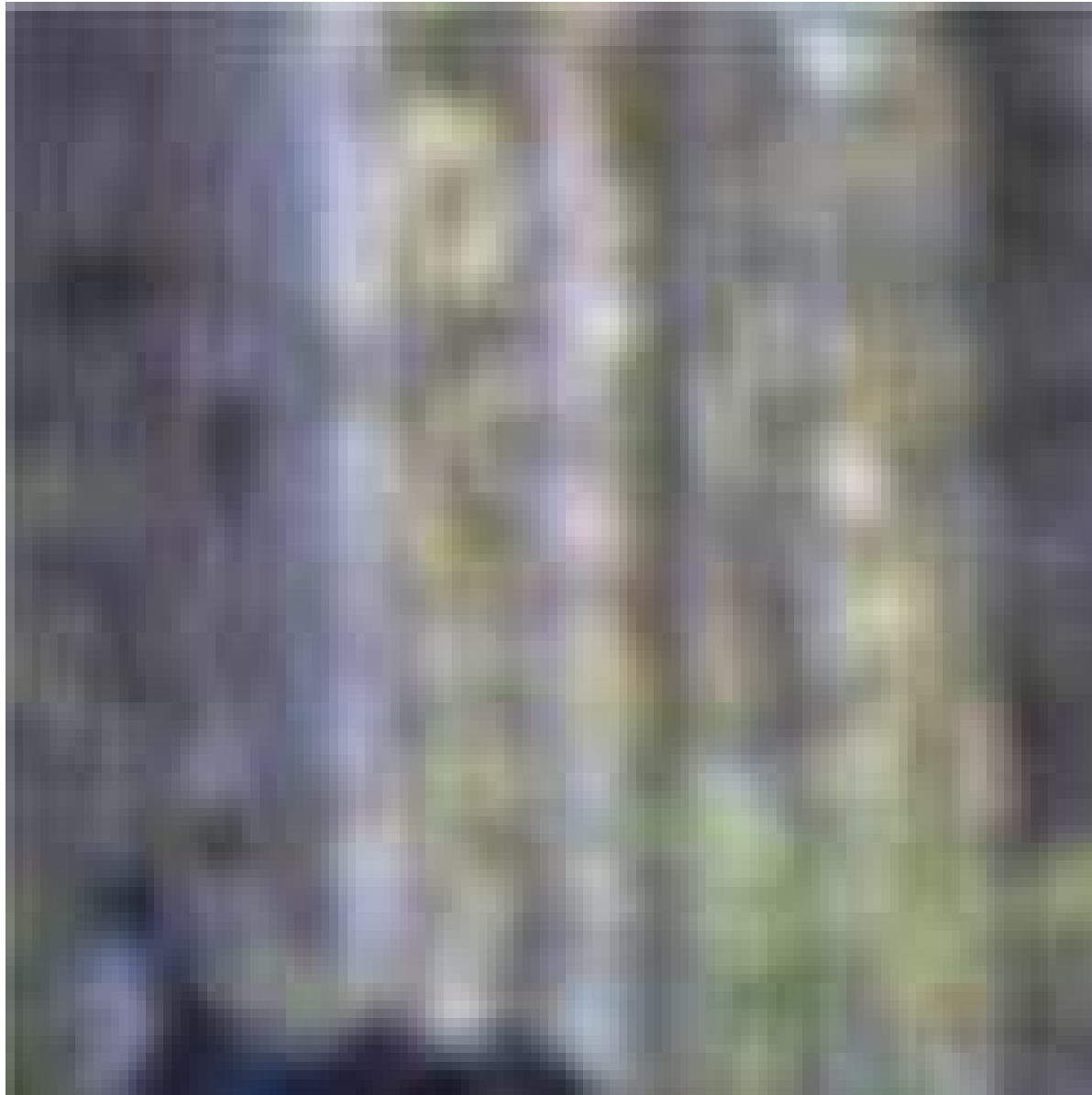




# Better smoothing with Gaussians



## Better smoothing with Gaussians

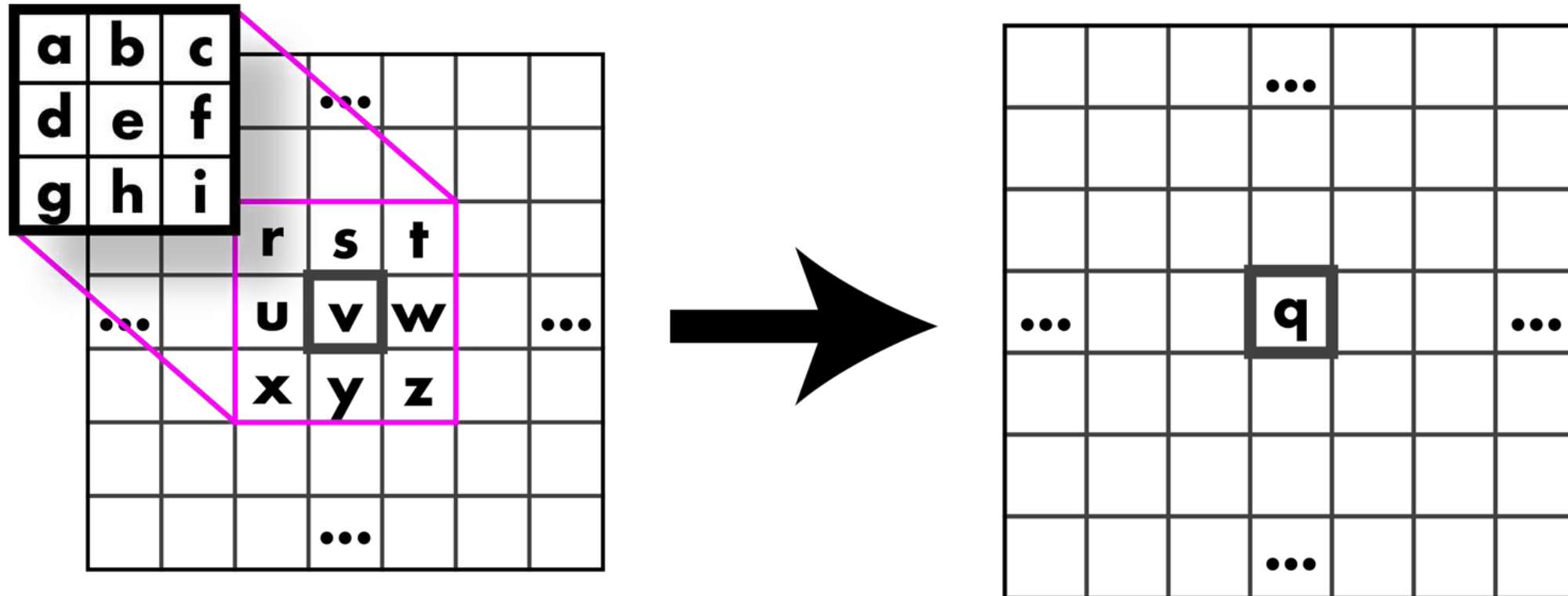


## Today's Agenda

- Averaging vs Interpolation
- Systems - filters
- Convolution
  - Box Filter
  - Gaussian
  - Cross correlation vs Convolution
- Examples of filters



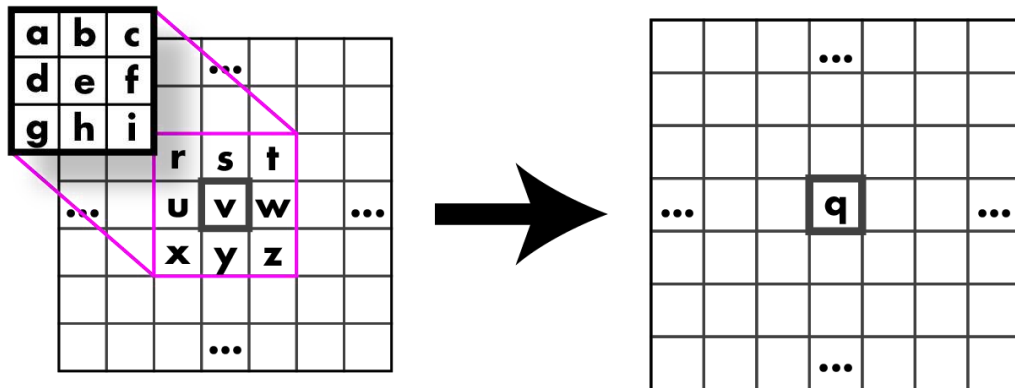
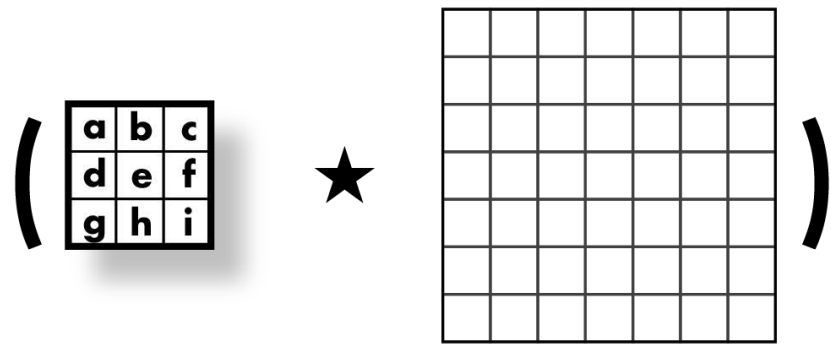
## So what is convolution??





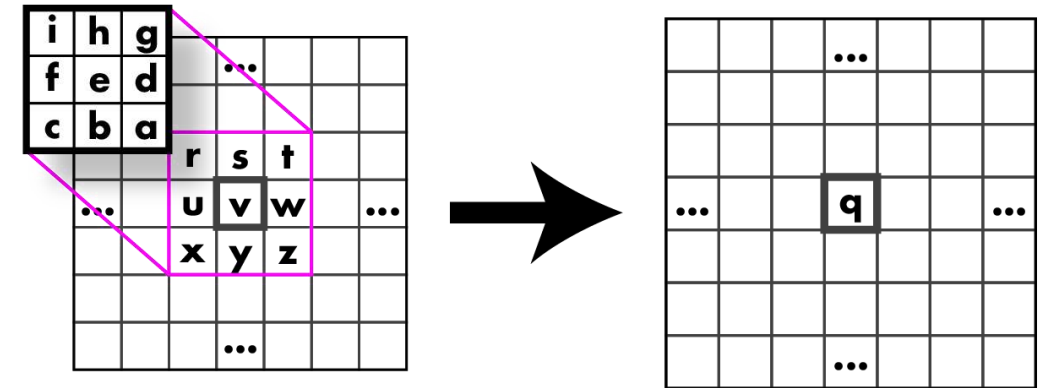
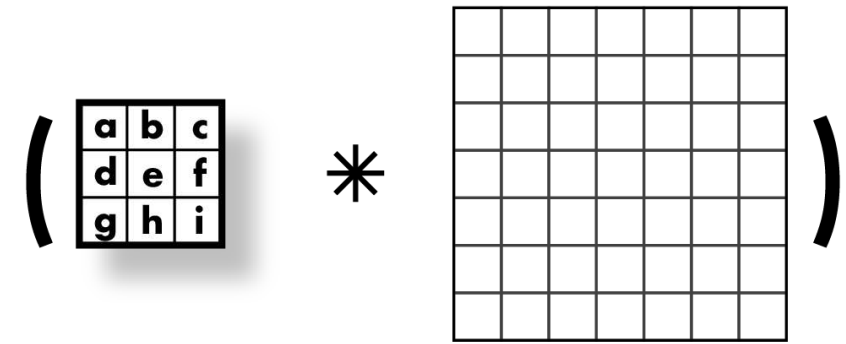
# Cross-Correlation vs Convolution

## Cross-Correlation



$$q = a \times r + b \times s + c \times t + d \times u + e \times v + f \times w + g \times x + h \times y + i \times z$$

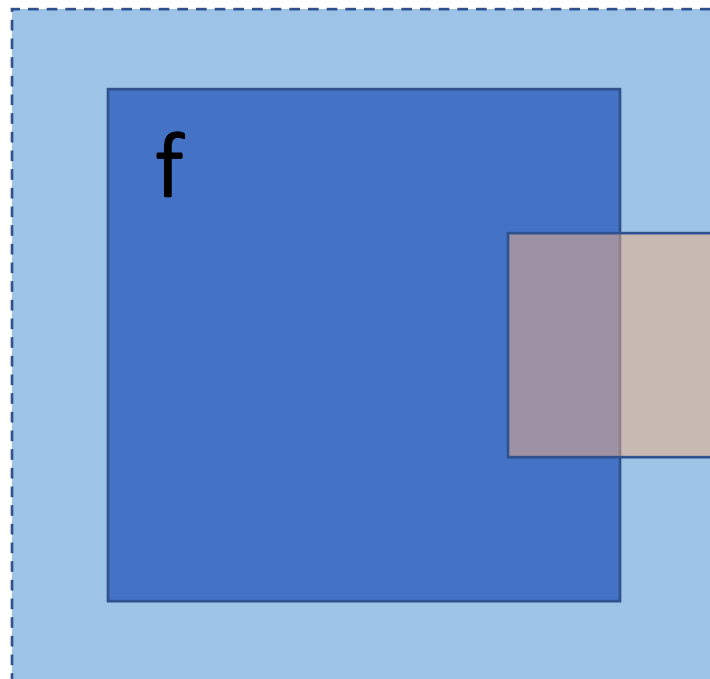
## Convolution



$$q = i \times r + h \times s + g \times t + f \times u + e \times v + d \times w + c \times x + b \times y + a \times z$$

# Image support and edge effect

- A computer will only convolve **finite support signals**
- What happens at the edge?

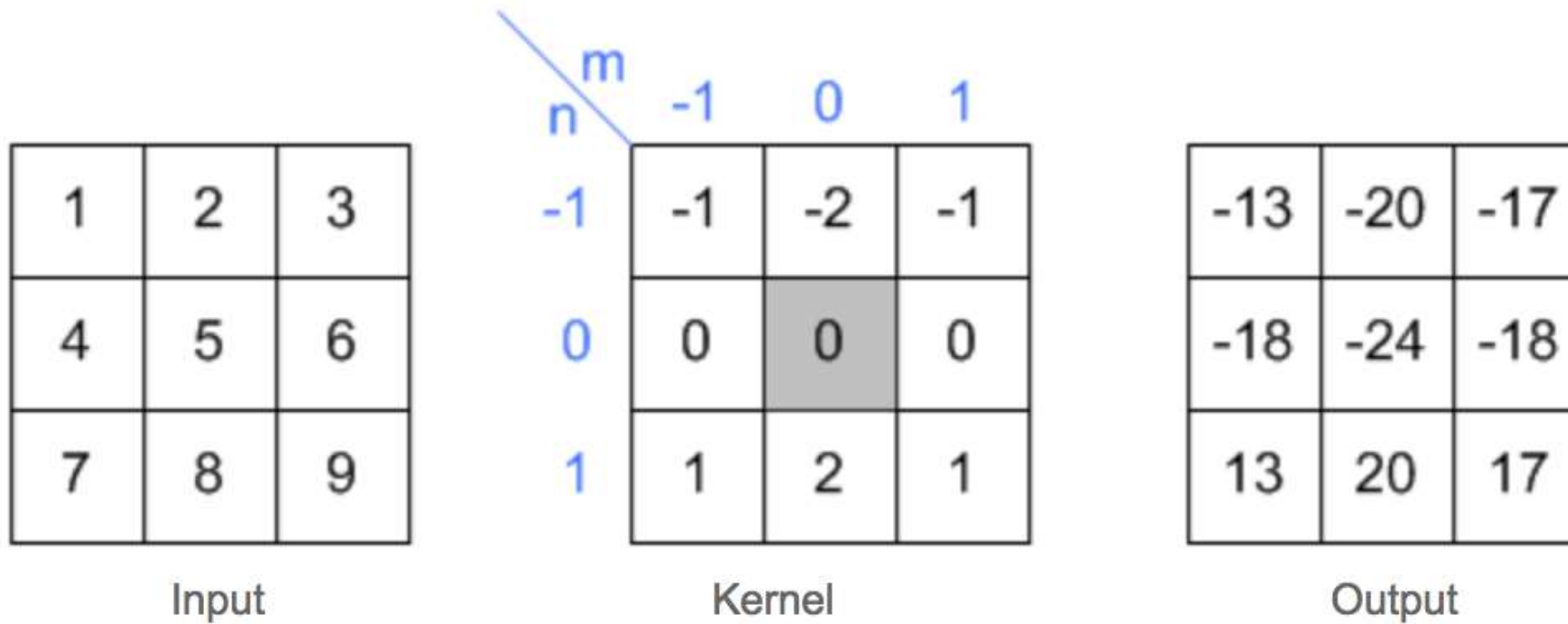


h

- zero “padding”
- edge value replication
- mirror extension
- ...

[Slide by Niebles]

## 2D convolution example



[Slide by Song Ho Ahn]

## 2D convolution example

1	2	1	
0	0	0	3
-1	-2	-1	6
	7	8	9

$$\begin{aligned}
 &= x[-1,-1] \cdot h[1,1] + x[0,-1] \cdot h[0,1] + x[1,-1] \cdot h[-1,1] \\
 &\quad + x[-1,0] \cdot h[1,0] + x[0,0] \cdot h[0,0] + x[1,0] \cdot h[-1,0] \\
 &\quad + x[-1,1] \cdot h[1,-1] + x[0,1] \cdot h[0,-1] + x[1,1] \cdot h[-1,-1] \\
 &= 0 \cdot 1 + 0 \cdot 2 + 0 \cdot 1 + 0 \cdot 0 + 1 \cdot 0 + 2 \cdot 0 + 0 \cdot (-1) + 4 \cdot (-2) + 5 \cdot (-1) = -13
 \end{aligned}$$

-13	-20	-17
-18	-24	-18
13	20	17

Output

[Slide by Song Ho Ahn]



## 2D convolution example

	1	2	1
0	1	2	3
-1	4	5	6
	7	8	9

$$\begin{aligned}
 &= x[0,-1] \cdot h[1,1] + x[1,-1] \cdot h[0,1] + x[2,-1] \cdot h[-1,1] \\
 &\quad + x[0,0] \cdot h[1,0] + x[1,0] \cdot h[0,0] + x[2,0] \cdot h[-1,0] \\
 &\quad + x[0,1] \cdot h[1,-1] + x[1,1] \cdot h[0,-1] + x[2,1] \cdot h[-1,-1] \\
 &= 0 \cdot 1 + 0 \cdot 2 + 0 \cdot 1 + 1 \cdot 0 + 2 \cdot 0 + 3 \cdot 0 + 4 \cdot (-1) + 5 \cdot (-2) + 6 \cdot (-1) = -20
 \end{aligned}$$

-13	-20	-17
-18	-24	-18
13	20	17

Output

[Slide by Song Ho Ahn]

## 2D convolution example

		1	2	1
1	0	2	0	3
4	-1	5	-2	6
7	8	9		

$$\begin{aligned}
 &= x[1,-1] \cdot h[1,1] + x[2,-1] \cdot h[0,1] + x[3,-1] \cdot h[-1,1] \\
 &+ x[1,0] \cdot h[1,0] + x[2,0] \cdot h[0,0] + x[3,0] \cdot h[-1,0] \\
 &+ x[1,1] \cdot h[1,-1] + x[2,1] \cdot h[0,-1] + x[3,1] \cdot h[-1,-1] \\
 &= 0 \cdot 1 + 0 \cdot 2 + 0 \cdot 1 + 2 \cdot 0 + 3 \cdot 0 + 0 \cdot 0 + 5 \cdot (-1) + 6 \cdot (-2) + 0 \cdot (-1) = -17
 \end{aligned}$$

-13	-20	-17
-18	-24	-18
13	20	17

Output

[Slide by Song Ho Ahn]

## 2D convolution example

1	2	1	3
0	0	0	6
-1	-2	-1	9

$$\begin{aligned}
 &= x[-1,0] \cdot h[1,1] + x[0,0] \cdot h[0,1] + x[1,0] \cdot h[-1,1] \\
 &\quad + x[-1,1] \cdot h[1,0] + x[0,1] \cdot h[0,0] + x[1,1] \cdot h[-1,0] \\
 &\quad + x[-1,2] \cdot h[1,-1] + x[0,2] \cdot h[0,-1] + x[1,2] \cdot h[-1,-1] \\
 &= 0 \cdot 1 + 1 \cdot 2 + 2 \cdot 1 + 0 \cdot 0 + 4 \cdot 0 + 5 \cdot 0 + 0 \cdot (-1) + 7 \cdot (-2) + 8 \cdot (-1) = -18
 \end{aligned}$$

-13	-20	-17
-18	-24	-18
13	20	17

Output

[Slide by Song Ho Ahn]

## 2D convolution example

1	2	1
1	2	3
0	0	0
4	5	6
-1	-2	-1
7	8	9

$$\begin{aligned}
 &= x[0,0] \cdot k[1,1] + x[1,0] \cdot k[0,1] + x[2,0] \cdot k[-1,1] \\
 &\quad + x[0,1] \cdot k[1,0] + x[1,1] \cdot k[0,0] + x[2,1] \cdot k[-1,0] \\
 &\quad + x[0,2] \cdot k[1,-1] + x[1,2] \cdot k[0,-1] + x[2,2] \cdot k[-1,-1] \\
 &= 1 \cdot 1 + 2 \cdot 2 + 3 \cdot 1 + 4 \cdot 0 + 5 \cdot 0 + 6 \cdot 0 + 7 \cdot (-1) + 8 \cdot (-2) + 9 \cdot (-1) = -24
 \end{aligned}$$

-13	-20	-17
-18	-24	-18
13	20	17

Output

[Slide by Song Ho Ahn]





## 2D convolution example

1	1	2	3	1
4	0	5	0	0
7	-1	8	-2	-1

$$\begin{aligned}
 &= x[1,0] \cdot h[1,1] + x[2,0] \cdot h[0,1] + x[3,0] \cdot h[-1,1] \\
 &\quad + x[1,1] \cdot h[1,0] + x[2,1] \cdot h[0,0] + x[3,1] \cdot h[-1,0] \\
 &\quad + x[1,2] \cdot h[1,-1] + x[2,2] \cdot h[0,-1] + x[3,2] \cdot h[-1,-1] \\
 &= 2 \cdot 1 + 3 \cdot 2 + 0 \cdot 1 + 5 \cdot 0 + 6 \cdot 0 + 0 \cdot 0 + 8 \cdot (-1) + 9 \cdot (-2) + 0 \cdot (-1) = -18
 \end{aligned}$$

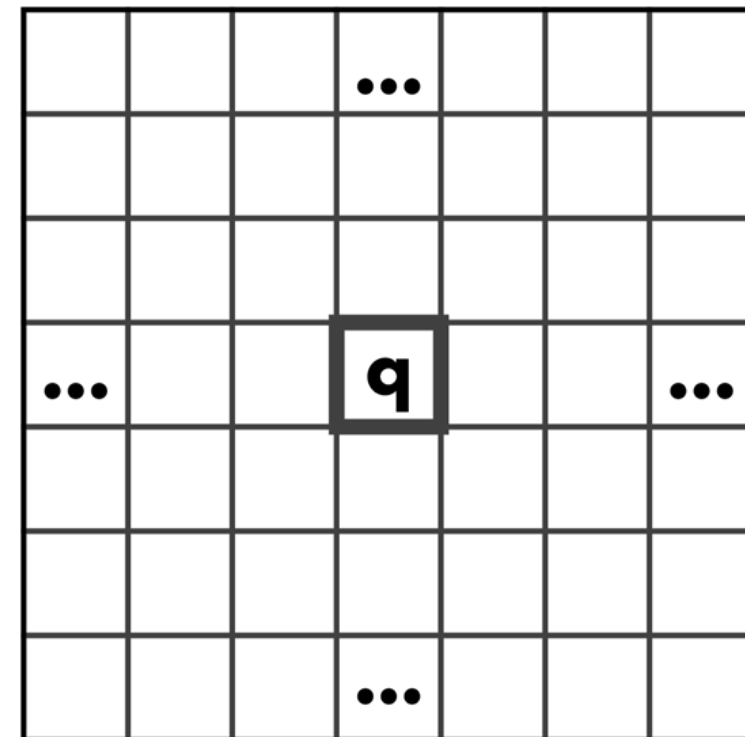
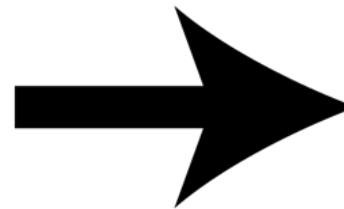
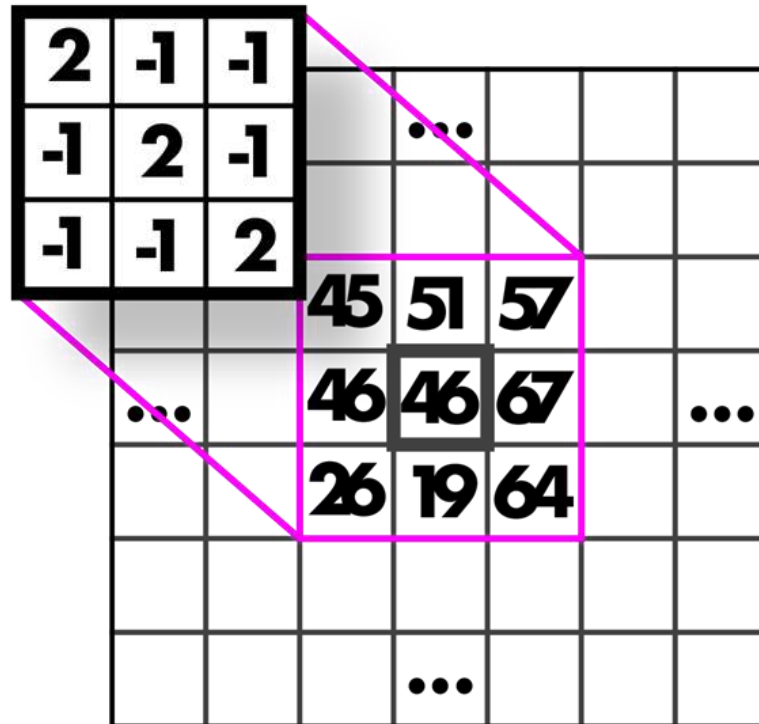
-13	-20	-17
-18	-24	-18
13	20	17

Output

[Slide by Song Ho Ahn]

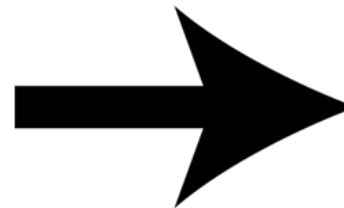


## Calculate it!



## Calculate it!

-1	-1	-1				
-1	8	-1	...			
-1	-1	-1				
			16	105	153	
...			15	104	113	...
			18	111	136	
			...			



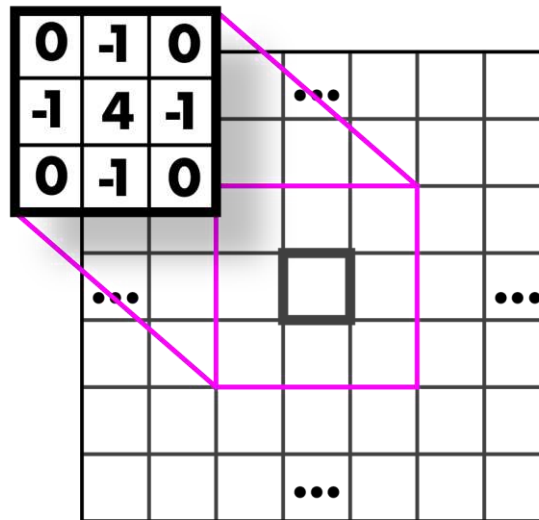
			...			
...			9			...
			...			

## Today's Agenda

- Averaging vs Interpolation
- Systems - filters
- Convolution
  - Box Filter
  - Gaussian
  - Cross correlation vs Convolution
- Examples of filters

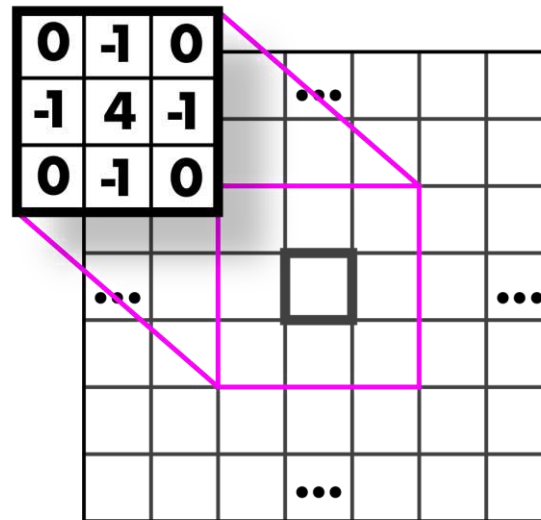


## Guess that kernel!



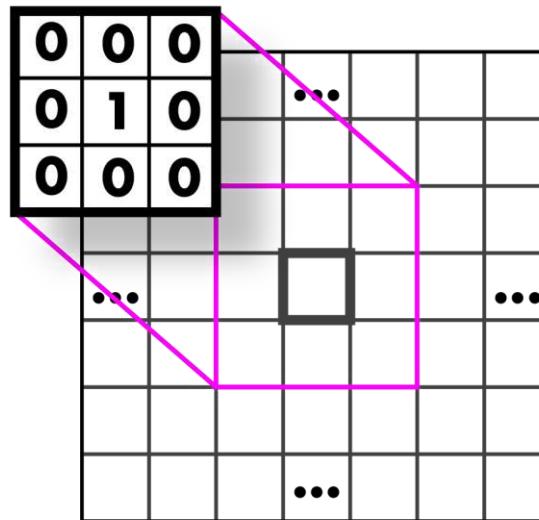
## Highpass Kernel: finds edges

applied to grayscale

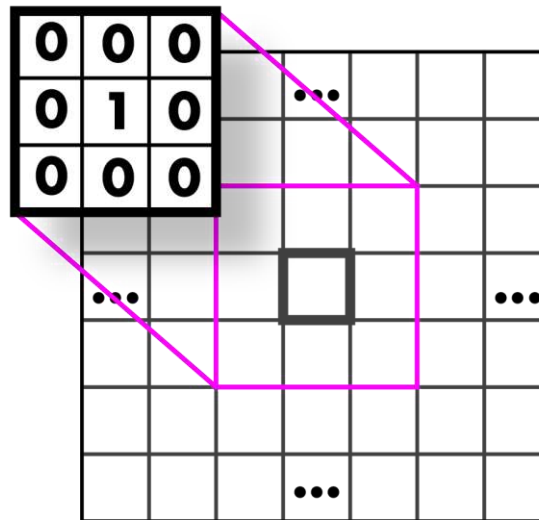




## Guess that kernel!

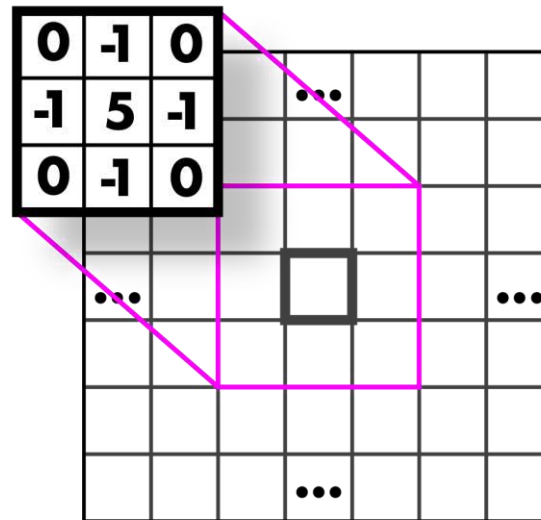


## Identity Kernel: Does nothing!



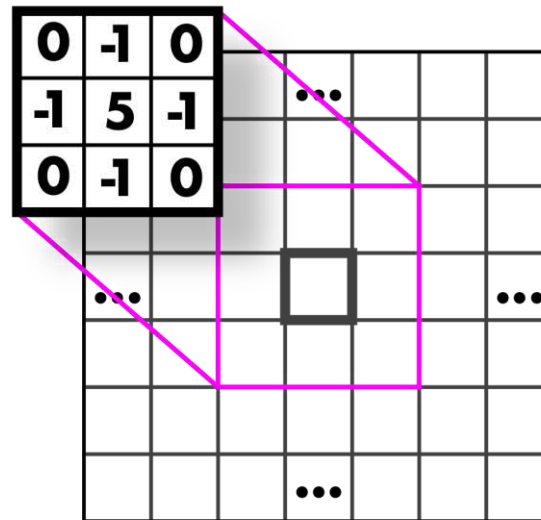


## Guess that kernel!



# Sharpen Kernel: sharpens!

applied to all three channels



Note: sharpen = highpass + identity! Why ?



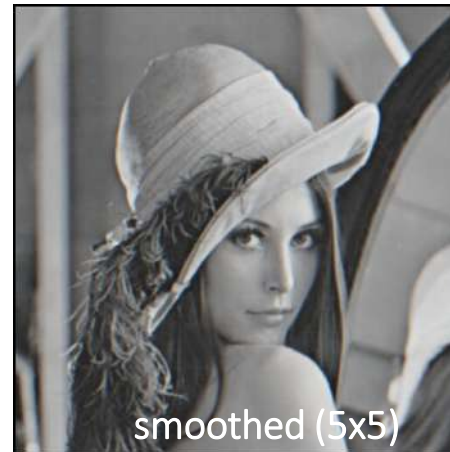
# Sharpen Kernel: sharpens!

What does blurring take away?

Highpass



-



=



Let's add it back:

Identity + Highpass



+

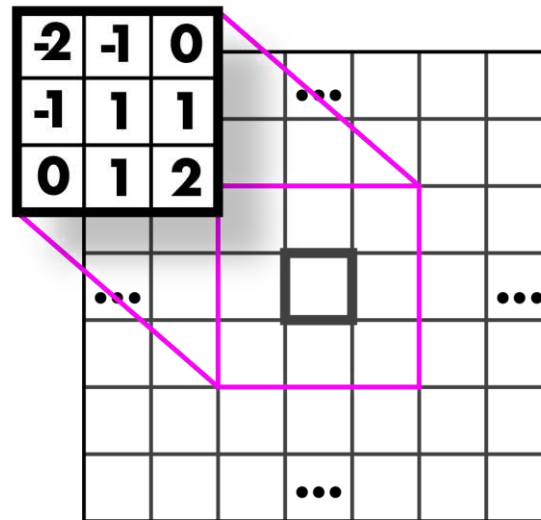


=



[Slide by D. Lowe]

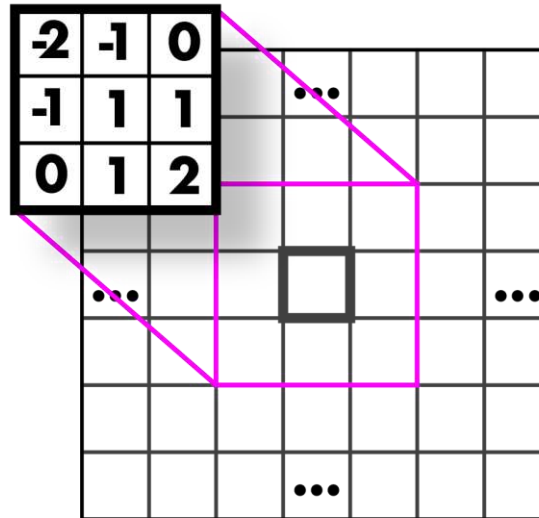
## Guess that kernel!



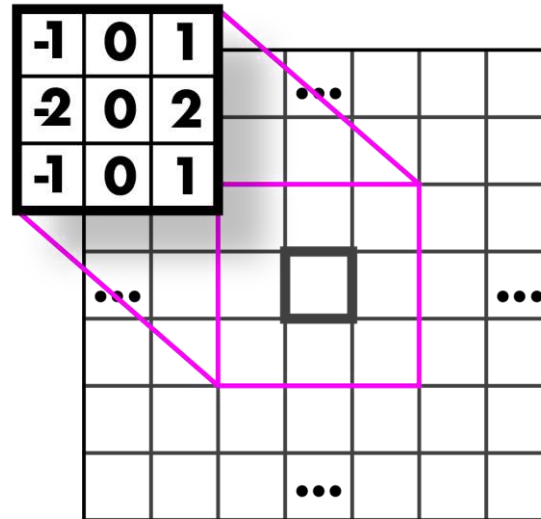
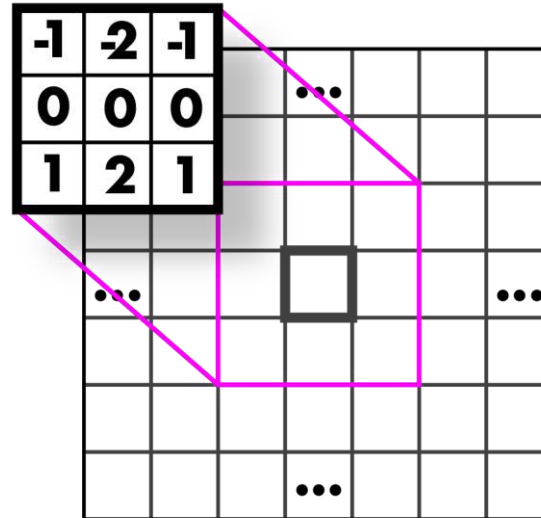


## Emboss Kernel: styling

applied to all three channels



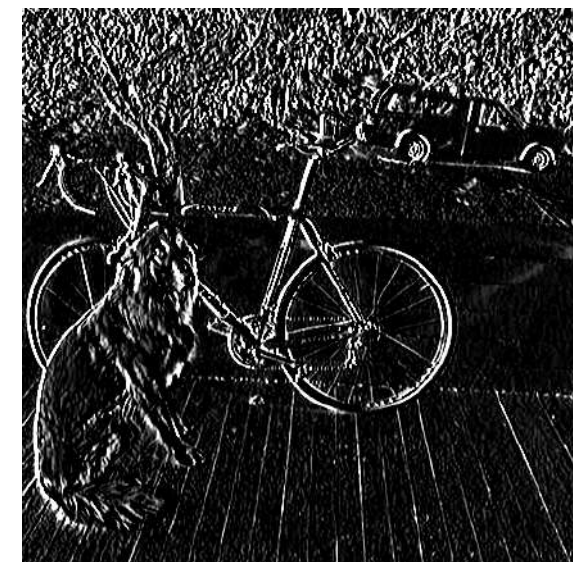
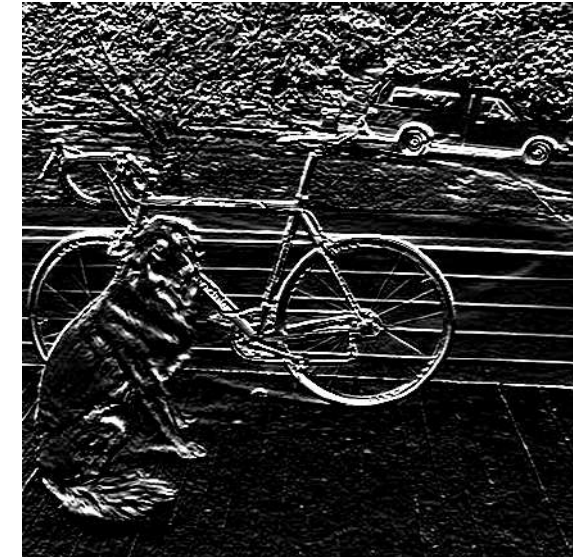
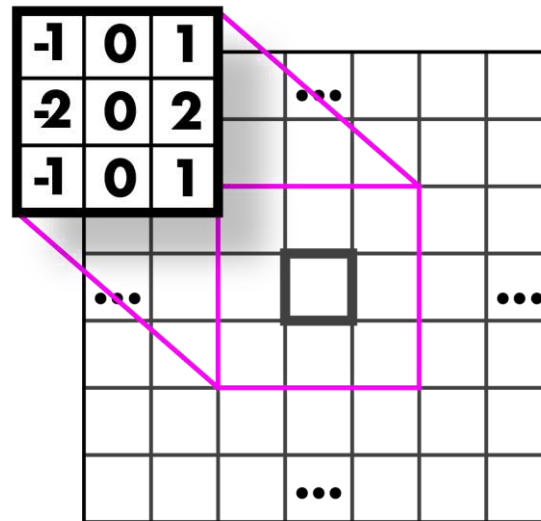
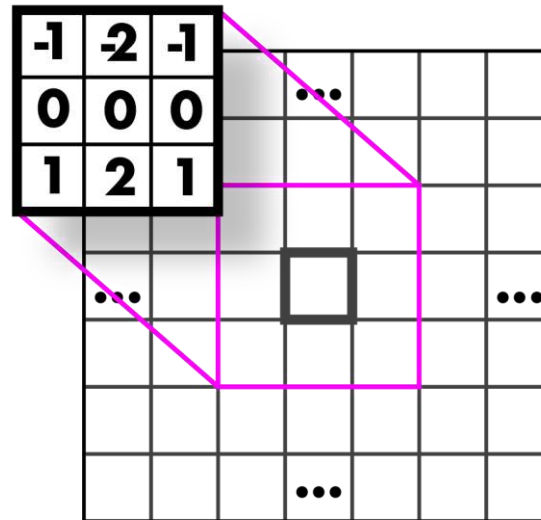
## Guess those kernels!





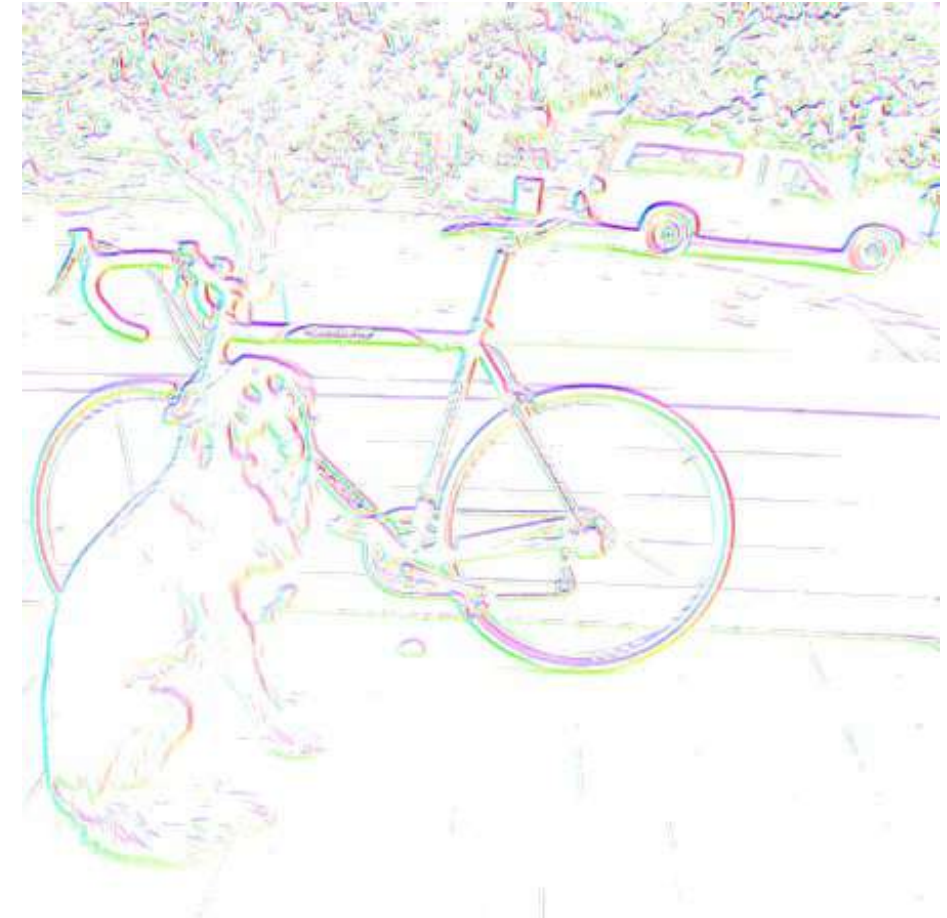
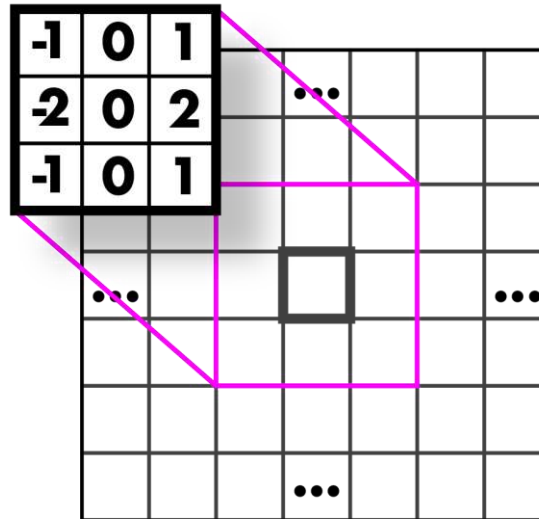
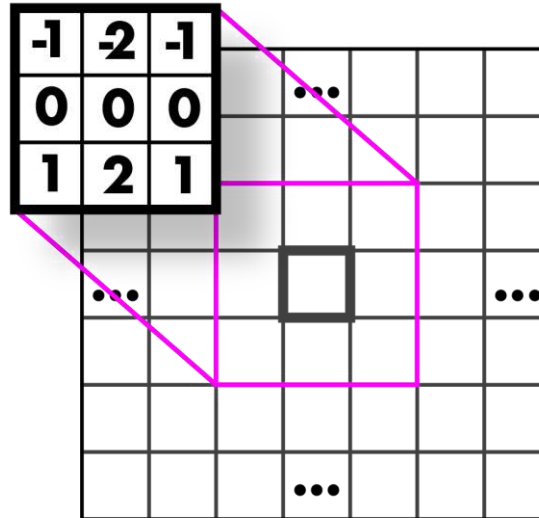
## Sobel Kernels: edges and...

applied to grayscale and thresholded



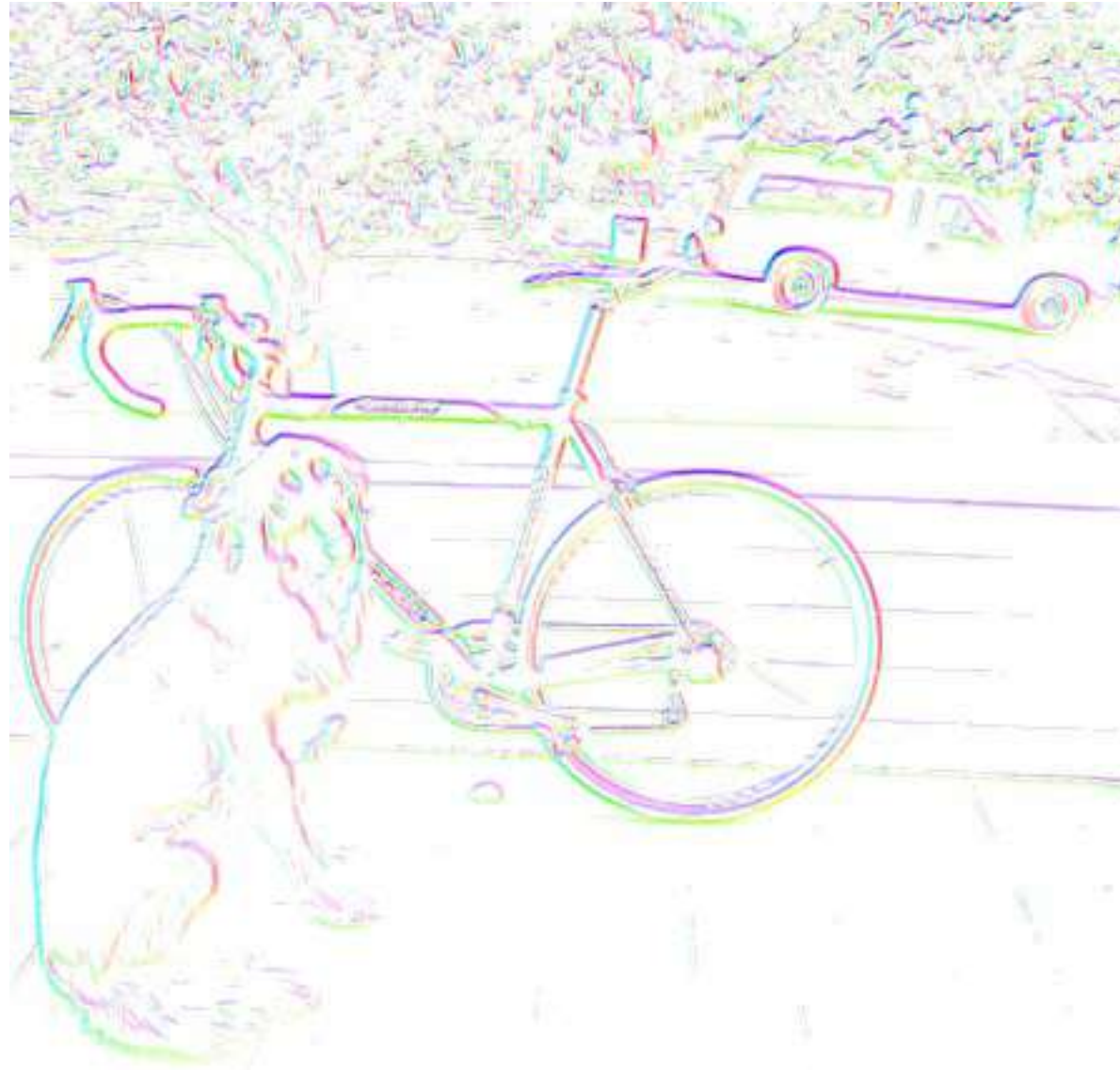


## Sobel Kernels: edges and gradient!





# Sobel Kernels: edges and gradient!



This visualization is showing the magnitude and direction of the gradient

# And so much more!!

- Next time
  - Edges
  - Features

**MAI4CAREU**

Master programmes in Artificial  
Intelligence 4 Careers in Europe



# Thank you.





University of Cyprus – MSc Artificial Intelligence

# MAI644 – COMPUTER VISION

## Lecture 6: Edges

**Melinos Averkiou**

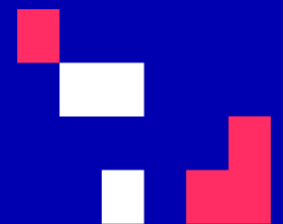
CYENS Centre of Excellence

University of Cyprus - Department of Computer Science

[m.averkiou@cyens.org.cy](mailto:m.averkiou@cyens.org.cy)



**CYENS**  
CENTRE OF EXCELLENCE





## Last time

- Averaging vs Interpolation
- Systems - filters
- Convolution
  - Box Filter
  - Gaussian
  - Cross correlation vs Convolution
- Examples of filters

# Today's Agenda

- What can we do with convolutions
- What is an edge – image derivatives
- Sobel filters
- Laplacian filters
- Difference of Gaussian filters
- Canny edge detection

**[material based on Joseph Redmon's course]**

# Today's Agenda

- What can we do with convolutions
- What is an edge – image derivatives
- Sobel filters
- Laplacian filters
- Difference of Gaussian filters
- Canny edge detection

# What can we do with convolutions

Mathematically: nice linear properties

- Commutative
  - $A*B = B*A$
- Associative
  - $A*(B*C) = (A*B)*C$
- Distributes over addition
  - $A*(B+C) = A*B + A*C$
- Plays well with scalars
  - $x(A*B) = (xA)*B = A*(xB)$



# What can we do with convolutions

This means some convolutions decompose:

- 2d gaussian is just composition of 1d gaussians
  - Faster to run 2 1d convolutions

# What can we do with convolutions

- Blurring
- Sharpening
- Edges
- Features
- Derivatives
- Super-resolution
- Classification
- Detection
- Image captioning
- ...

# What can we do with convolutions

- Blurring
- Sharpening
- Edges
- Features
- Derivatives
- Super-resolution
- Classification
- Detection
- Image captioning
- ...

A big part of computer vision  
is **convolutions**

So what can we do with these convolutions anyway?

- Blurring
- Sharpening
- Edges
- Features
- Derivatives
- Super-resolution
- Classification
- Detection
- Image captioning
- ...

A big part of computer vision  
is **convolutions**

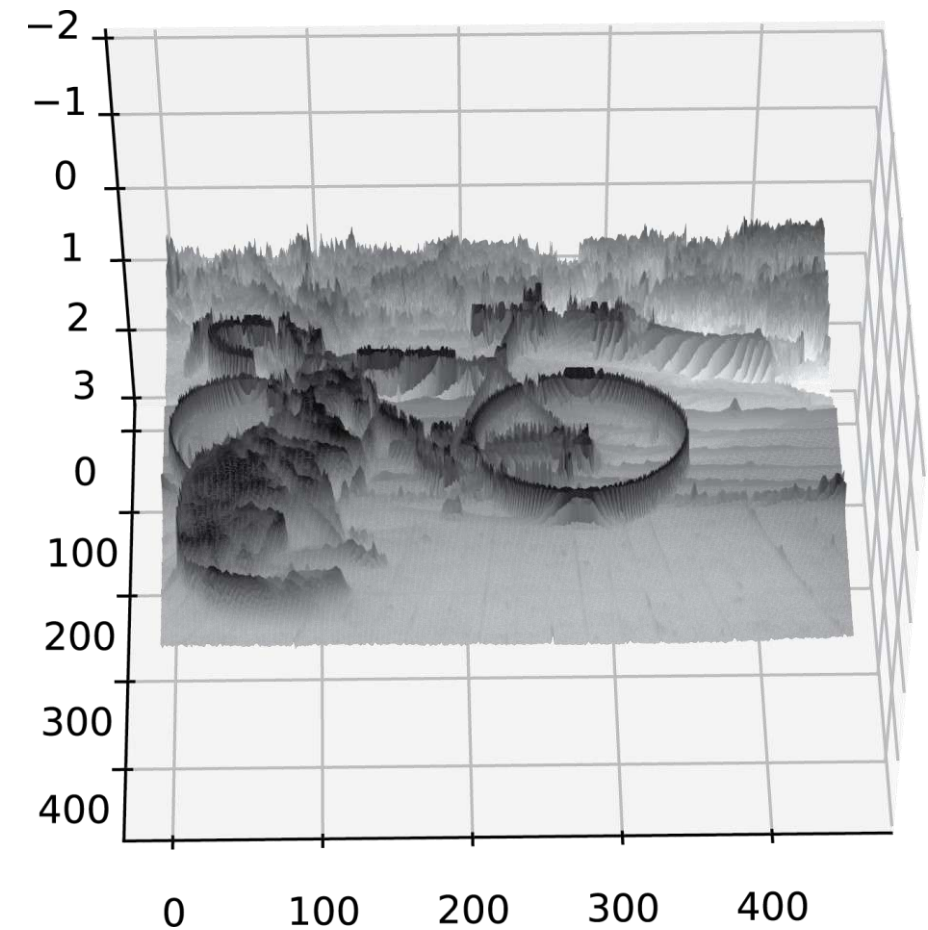
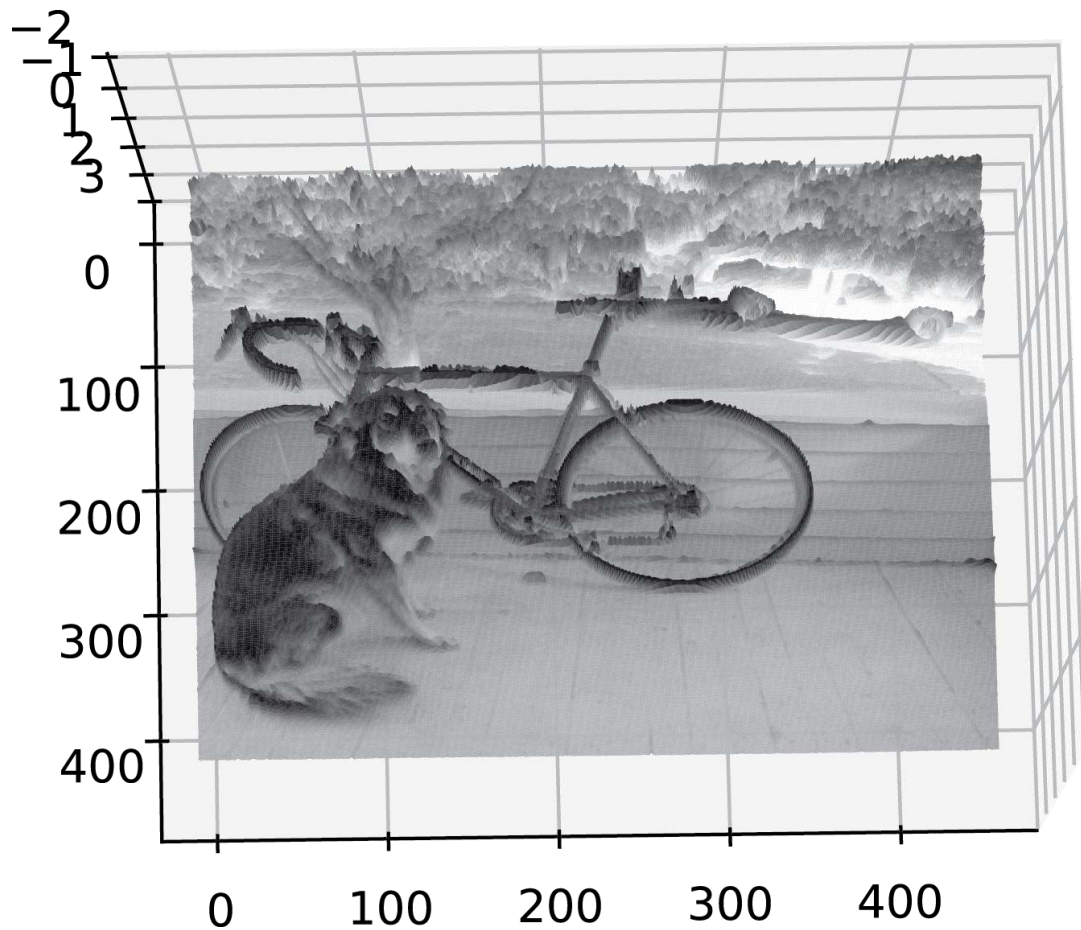


# Today's Agenda

- What can we do with convolutions
- What is an edge – image derivatives
- Sobel filters
- Laplacian filters
- Difference of Gaussian filters
- Canny edge detection

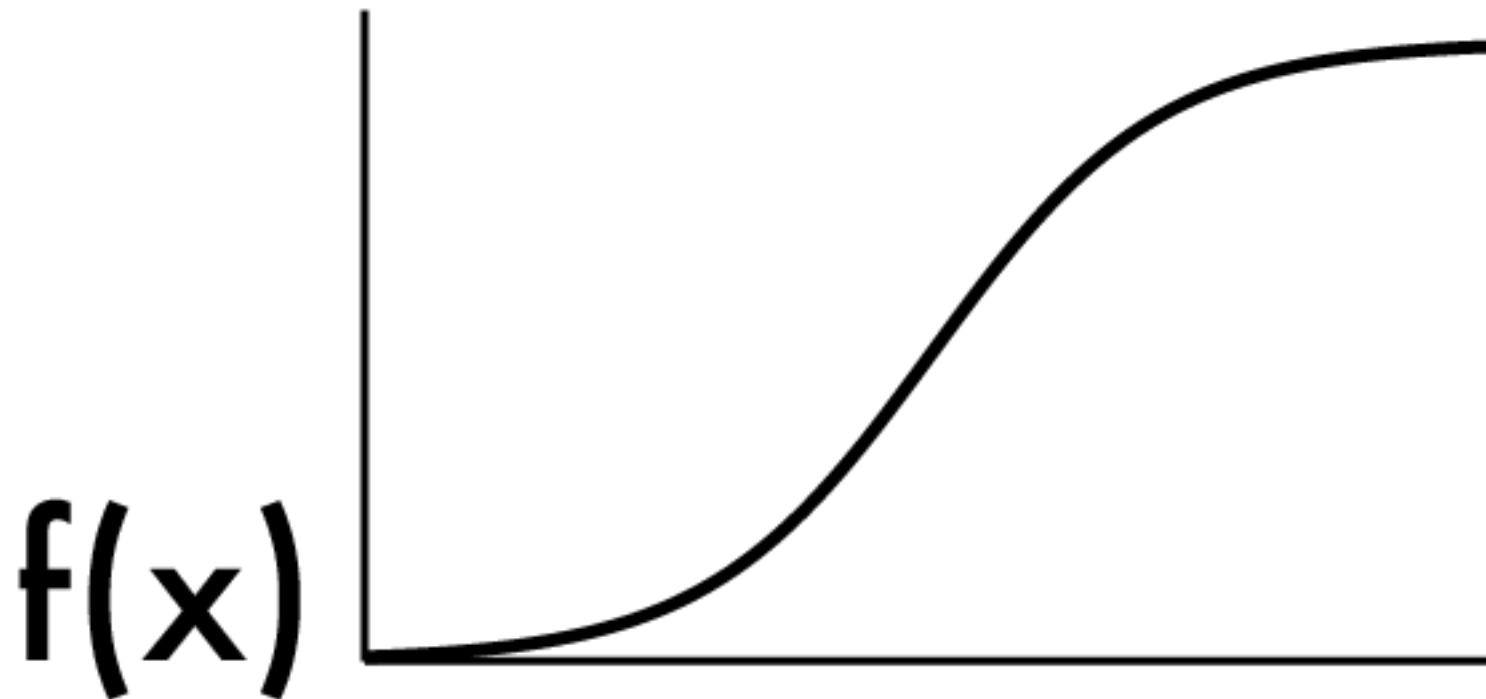
# What's an edge?

- Image is a function
- Edges are rapid changes in this function



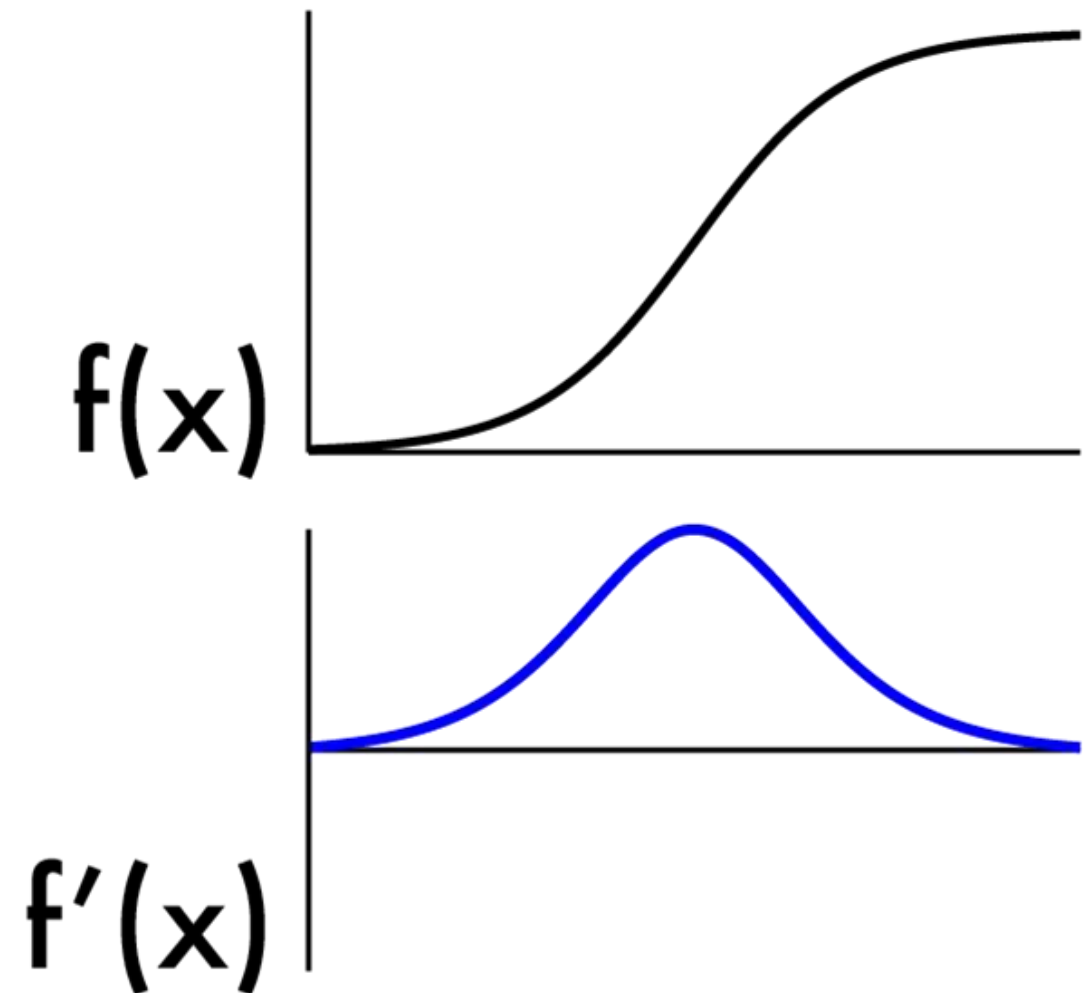
# What's an edge?

- Image is a function
- Edges are rapid changes in this function



# Finding edges

- Could take derivative
- Edges = high response



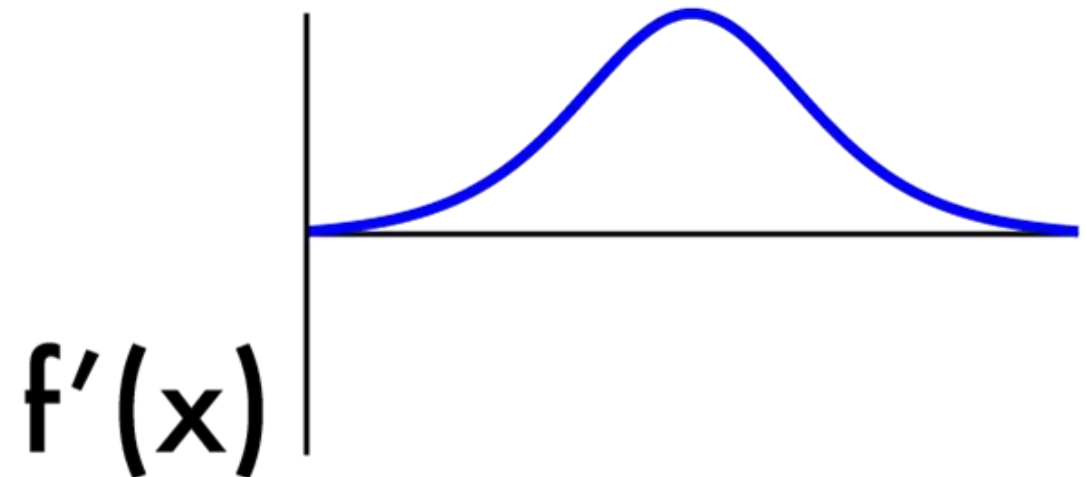
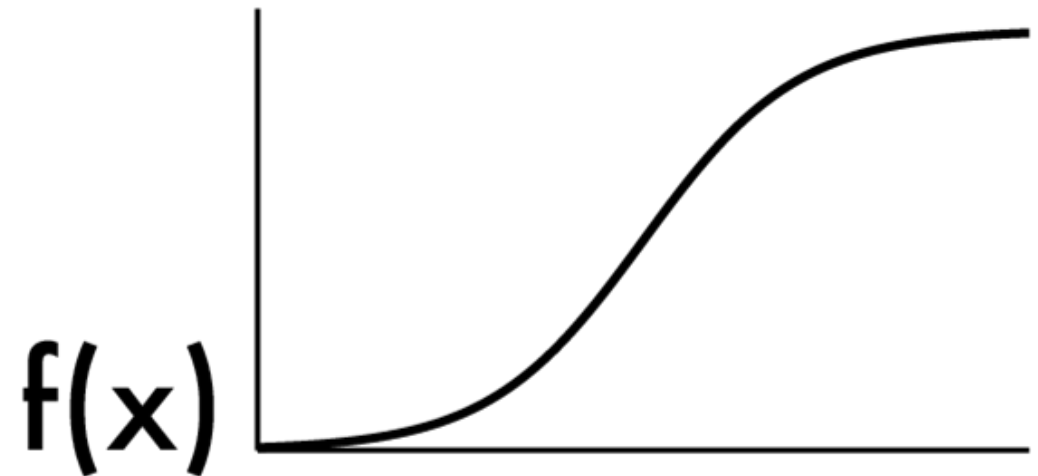


# Image derivatives

- Recall:

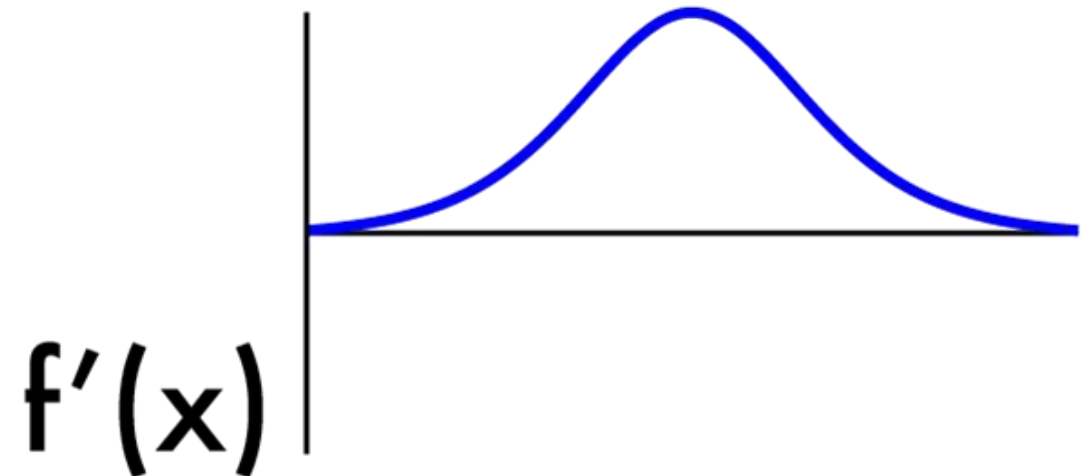
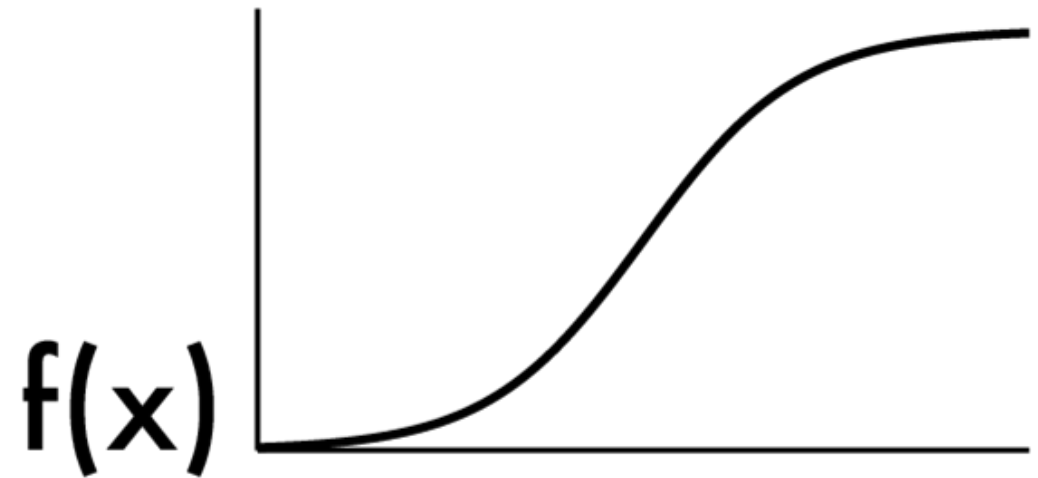
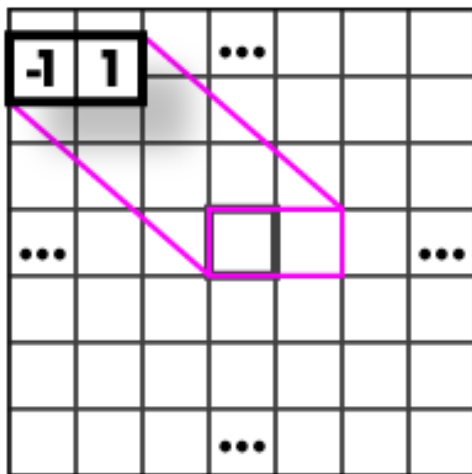
$$f'(a) = \lim_{h \rightarrow 0} \frac{f(a+h) - f(a)}{h}$$

- We don't have an "actual" function, must estimate
- Possibility: set  $h = 1$
- What will that look like?



## Image derivatives

- Recall:
  - $f'(a) = \lim_{h \rightarrow 0} \frac{f(a+h) - f(a)}{h}$ .
- We don't have an "actual" function, must estimate
- Possibility: set  $h = 1$
- What will that look like?

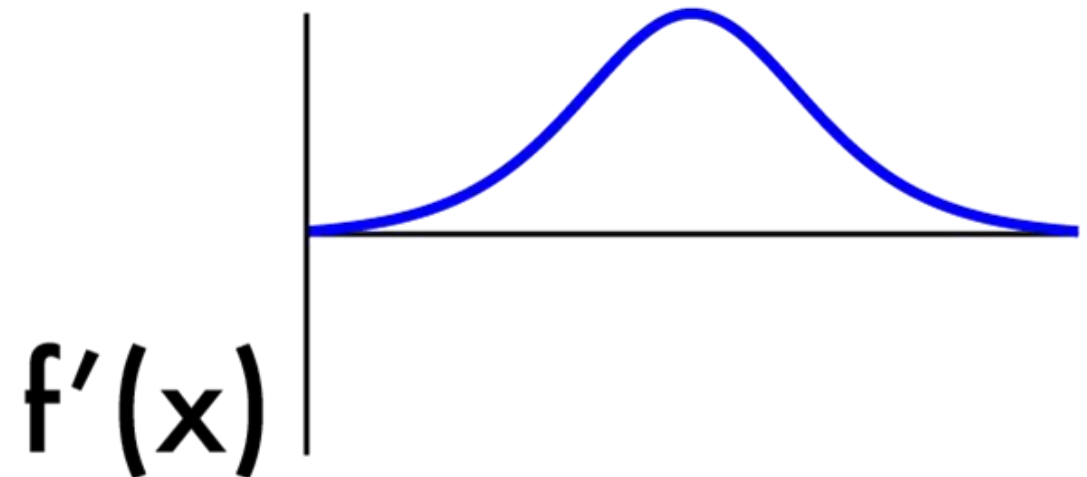
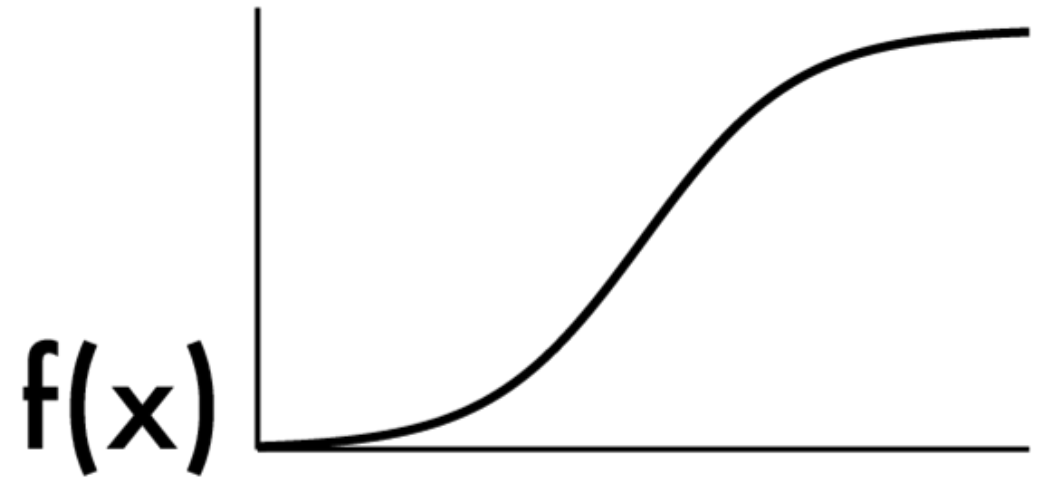


# Image derivatives

- Recall:

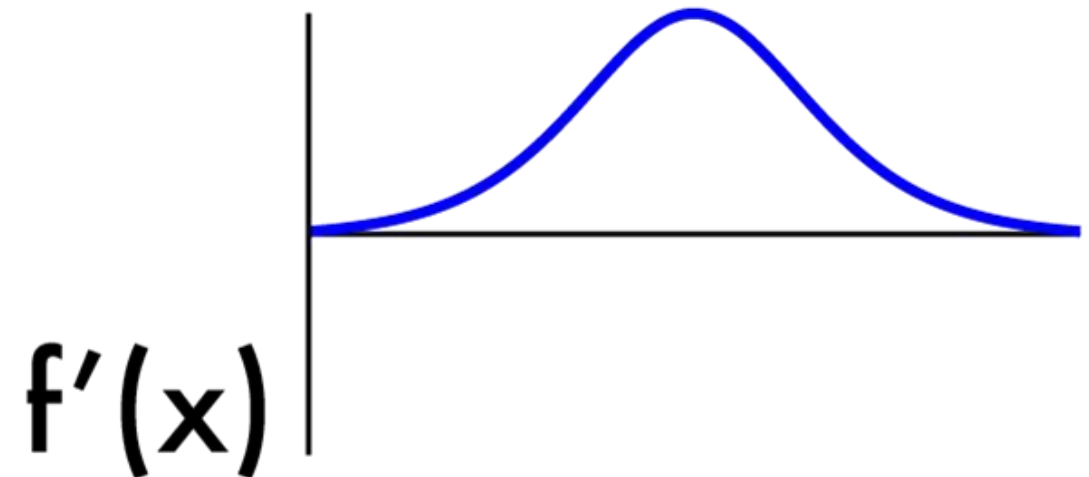
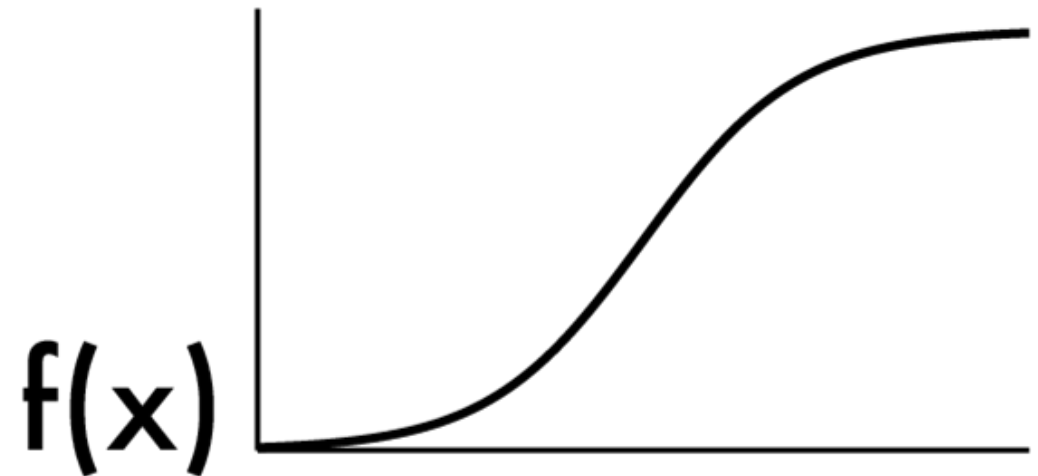
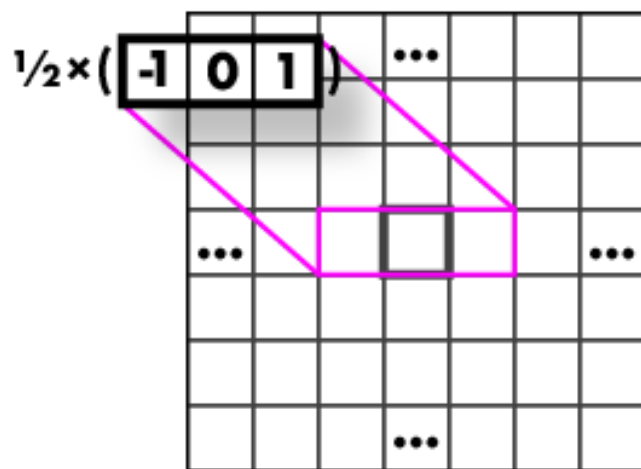
$$f'(a) = \lim_{h \rightarrow 0} \frac{f(a+h) - f(a)}{h}$$

- We don't have an "actual" function, must estimate
- Possibility: set  $h = 2$
- What will that look like?



# Image derivatives

- Recall:
  - $f'(a) = \lim_{h \rightarrow 0} \frac{f(a+h) - f(a)}{h}$ .
- We don't have an "actual" function, must estimate
- Possibility: set  $h = 2$
- What will that look like?

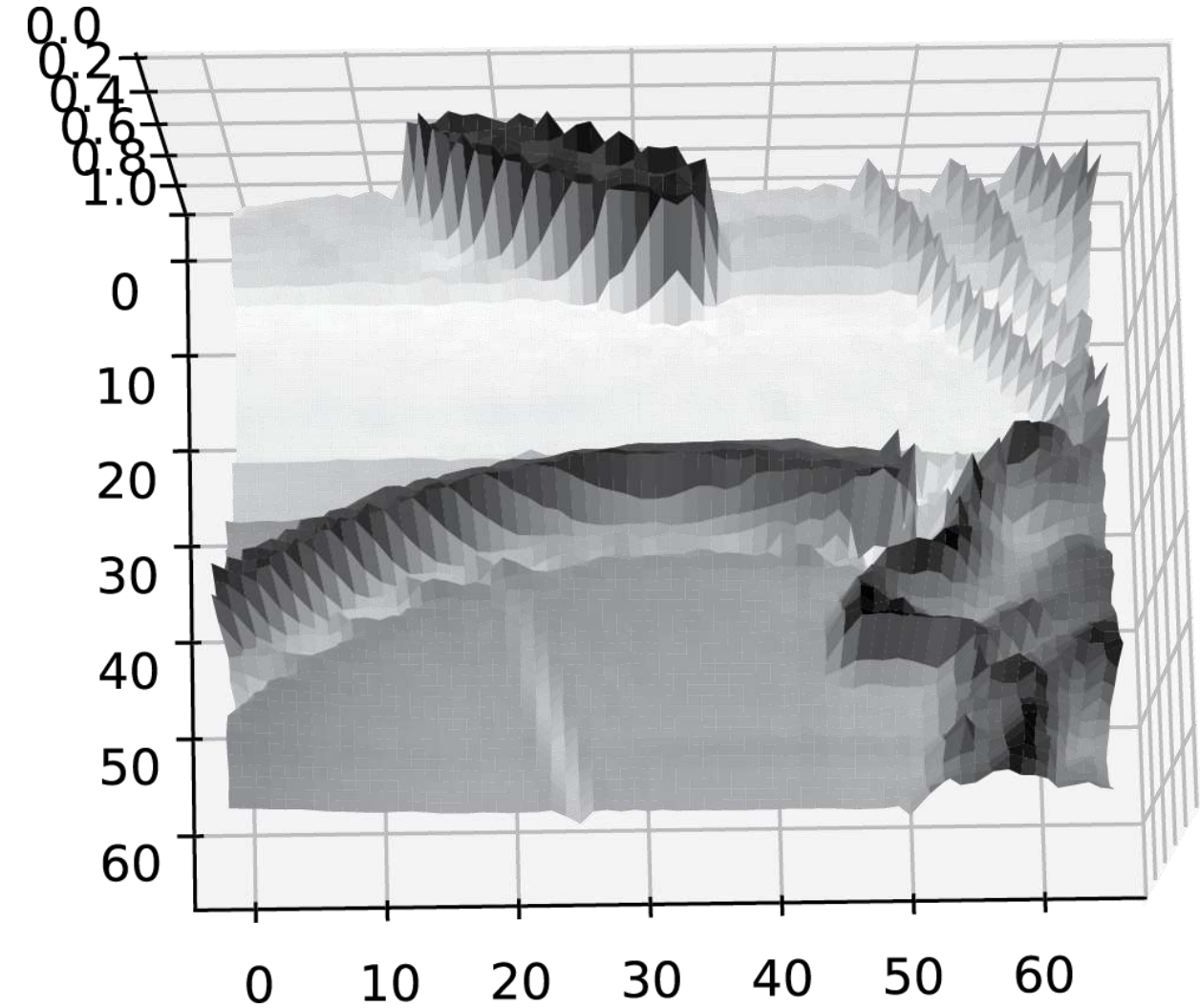
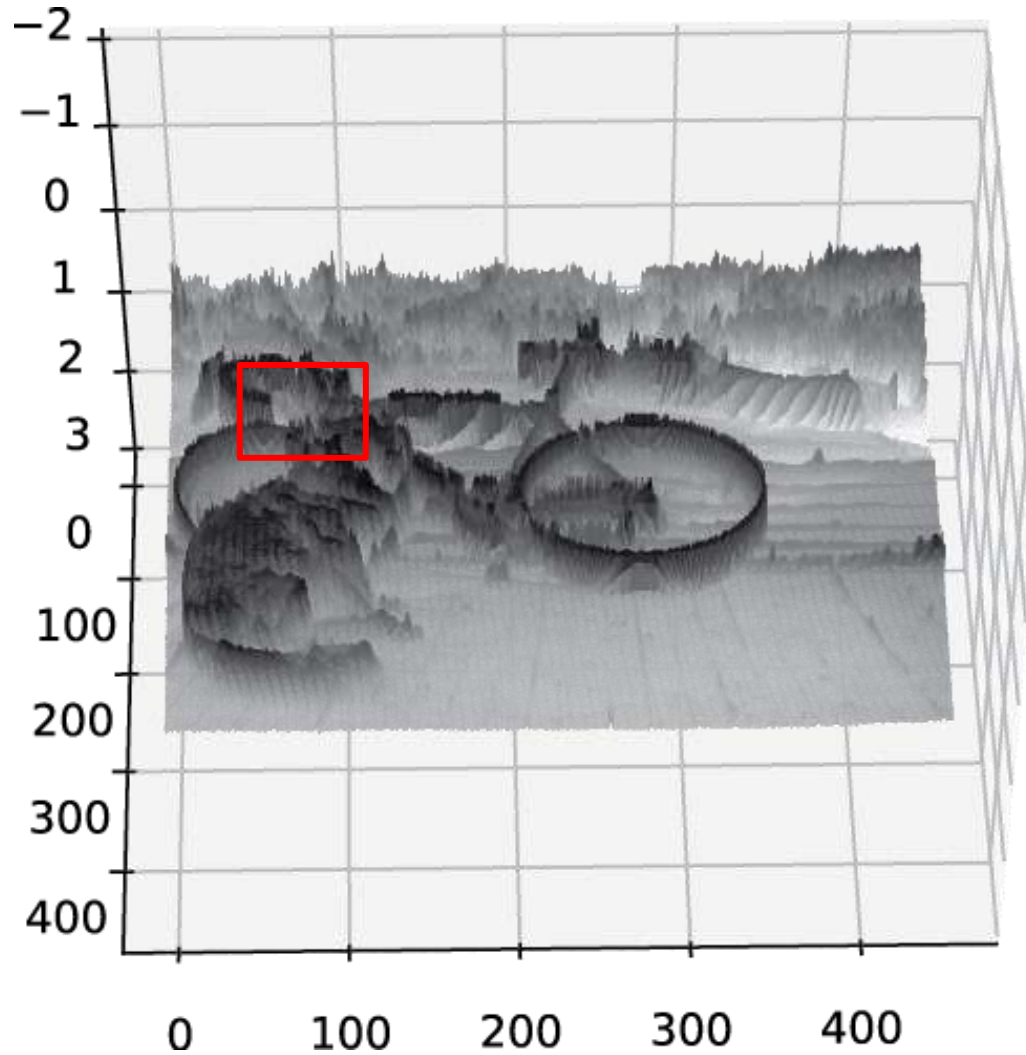




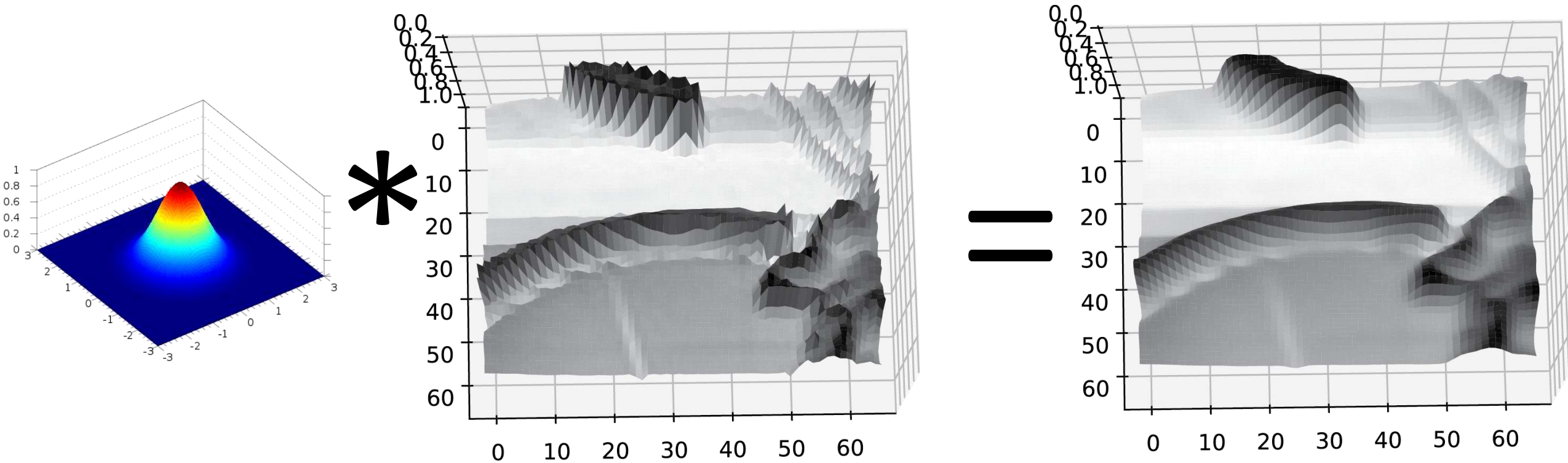
# Today's Agenda

- What can we do with convolutions
- What is an edge – image derivatives
- Sobel filters
- Laplacian filters
- Difference of Gaussian filters
- Canny edge detection

# Images are noisy!



# But we already know how to smooth



# Smooth first, then derivative

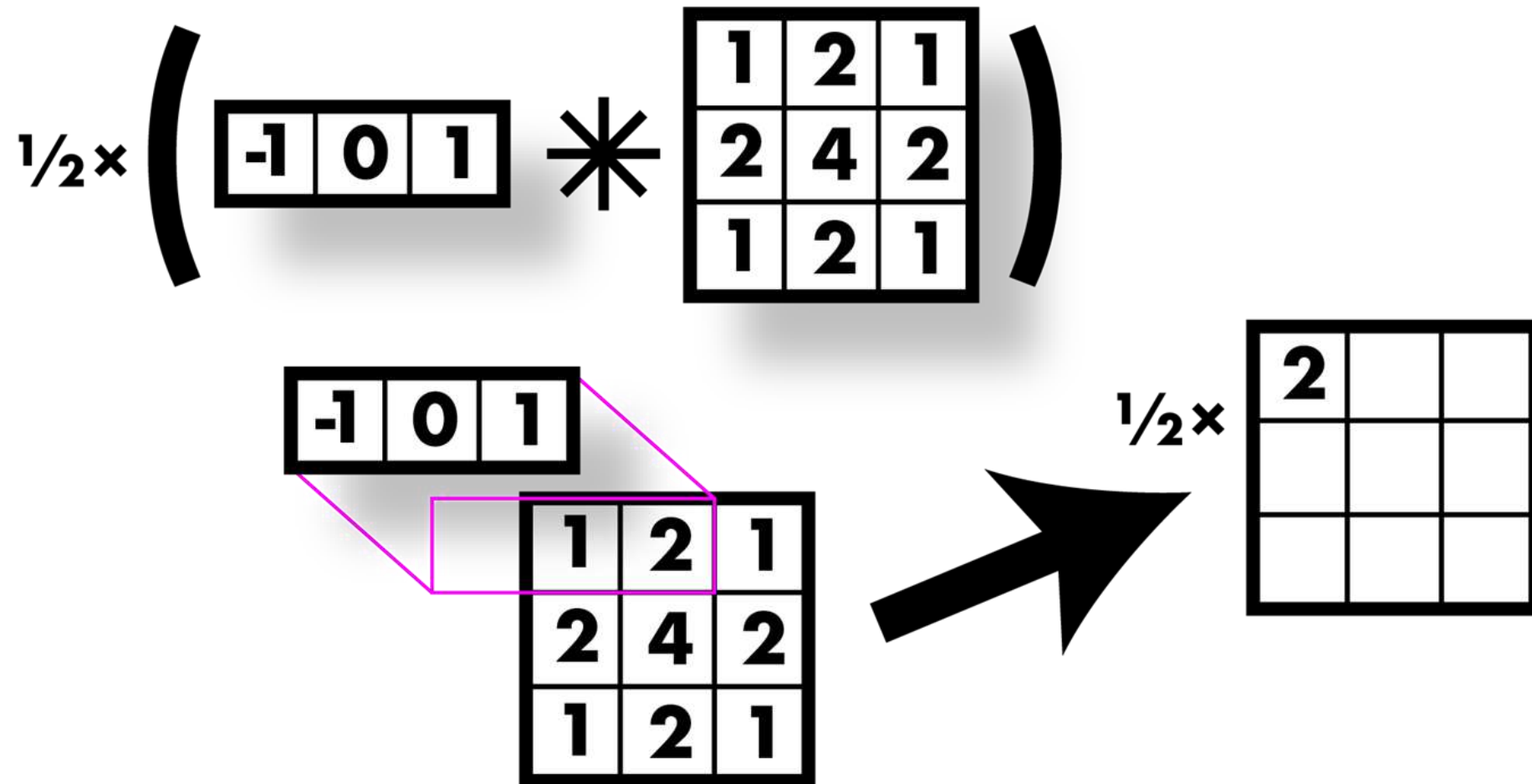
$$\frac{1}{2} \times \begin{pmatrix} -1 & 0 & 1 \end{pmatrix} * \begin{pmatrix} 1 & 2 & 1 \\ 2 & 4 & 2 \\ 1 & 2 & 1 \end{pmatrix} * \begin{pmatrix} \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} \\ \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} \\ \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} \\ \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} \\ \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} \\ \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} \\ \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} & \phantom{0} \end{pmatrix}$$



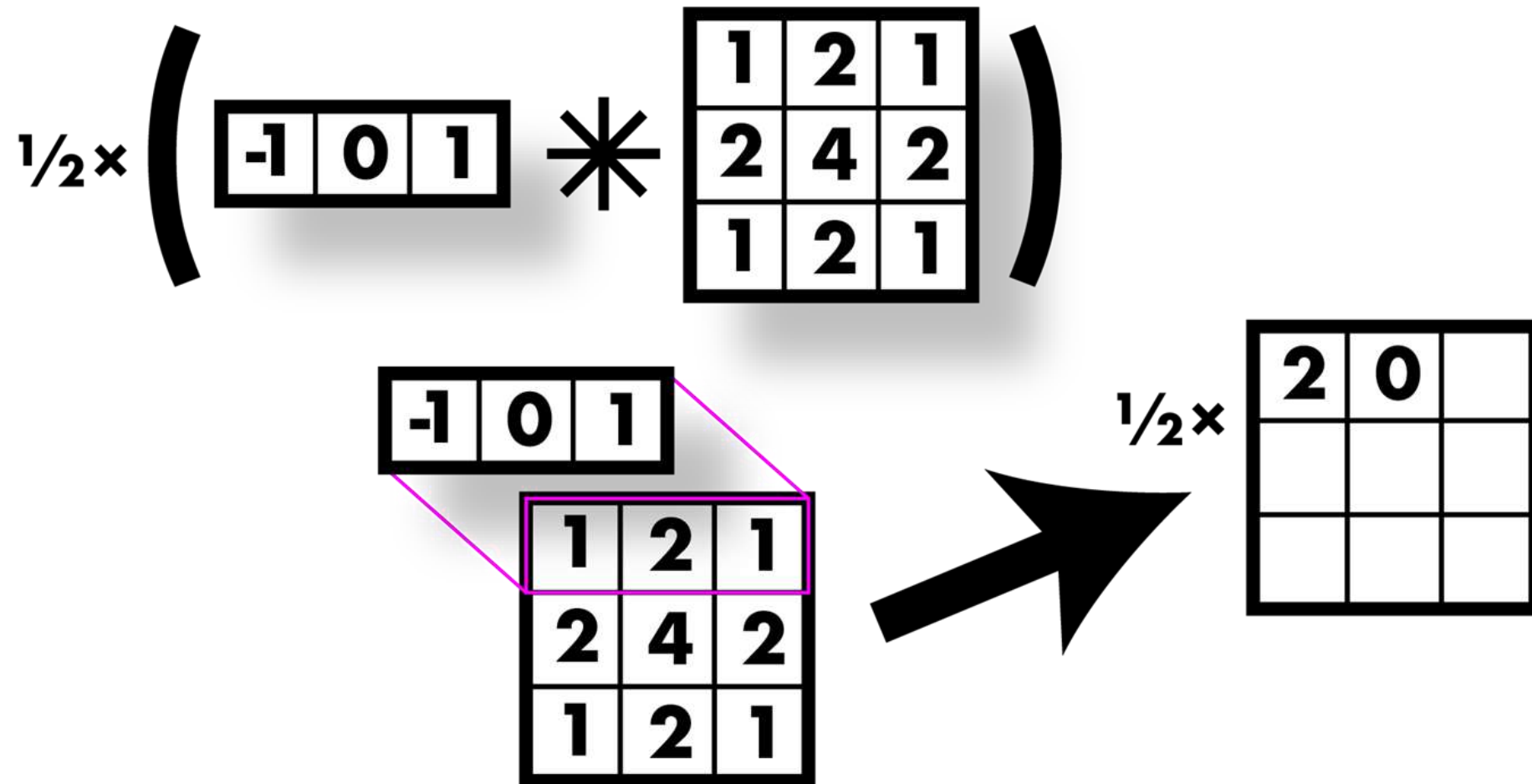
# Smooth first, then derivative

$$\frac{1}{2} \times \left( \begin{bmatrix} -1 & 0 & 1 \end{bmatrix} * \begin{bmatrix} 1 & 2 & 1 \\ 2 & 4 & 2 \\ 1 & 2 & 1 \end{bmatrix} \right) *$$


# Smooth first, then derivative

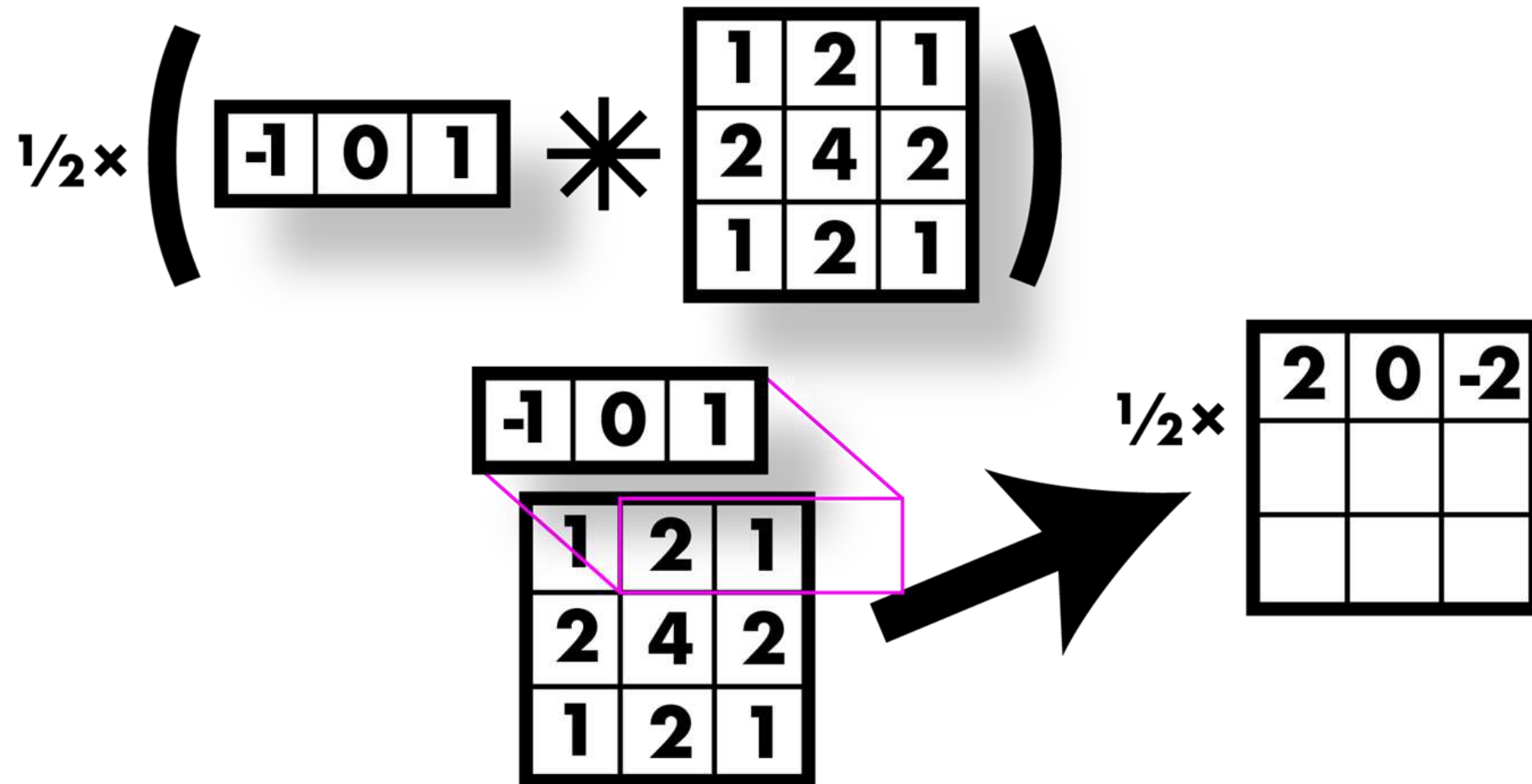


## Smooth first, then derivative





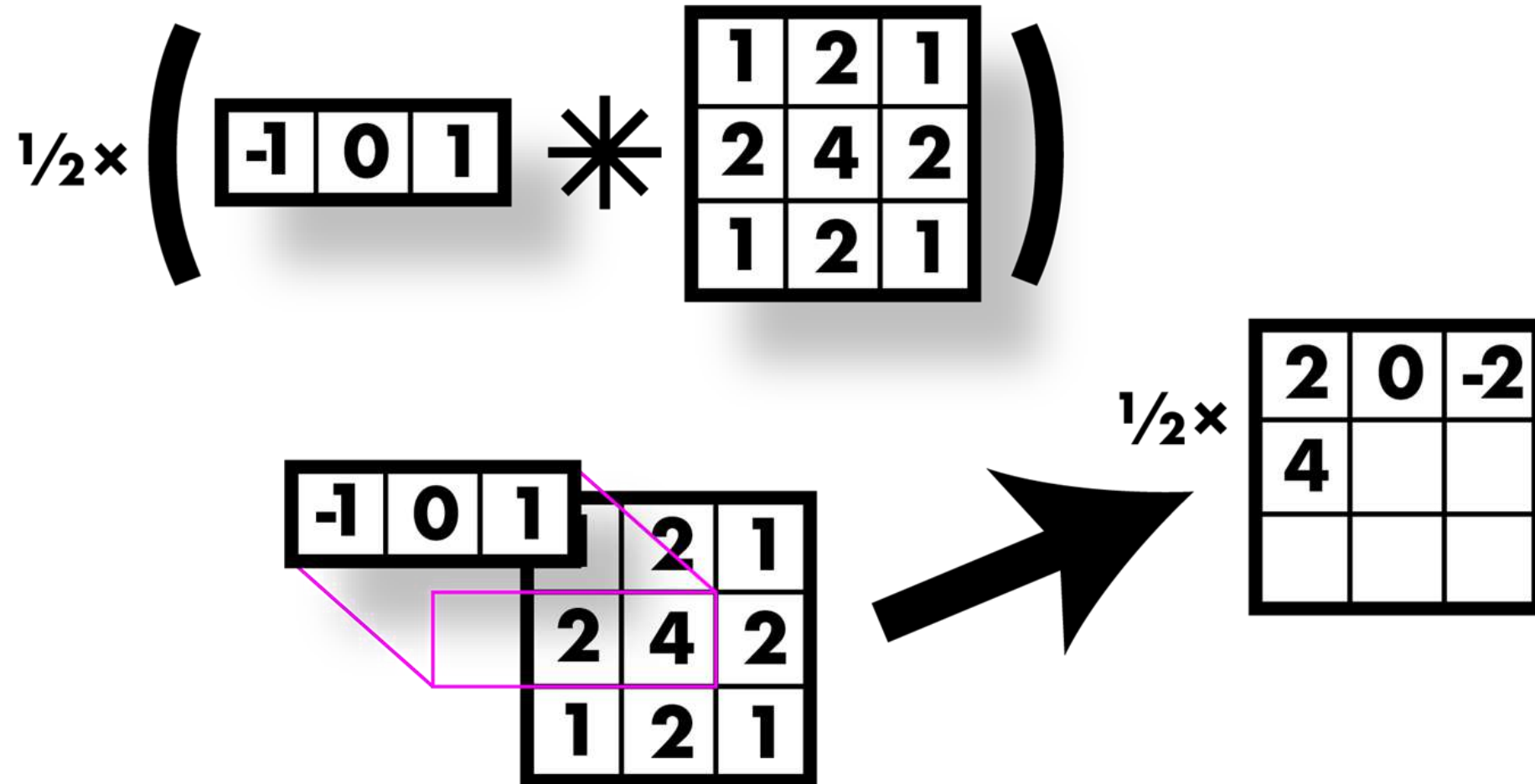
## Smooth first, then derivative





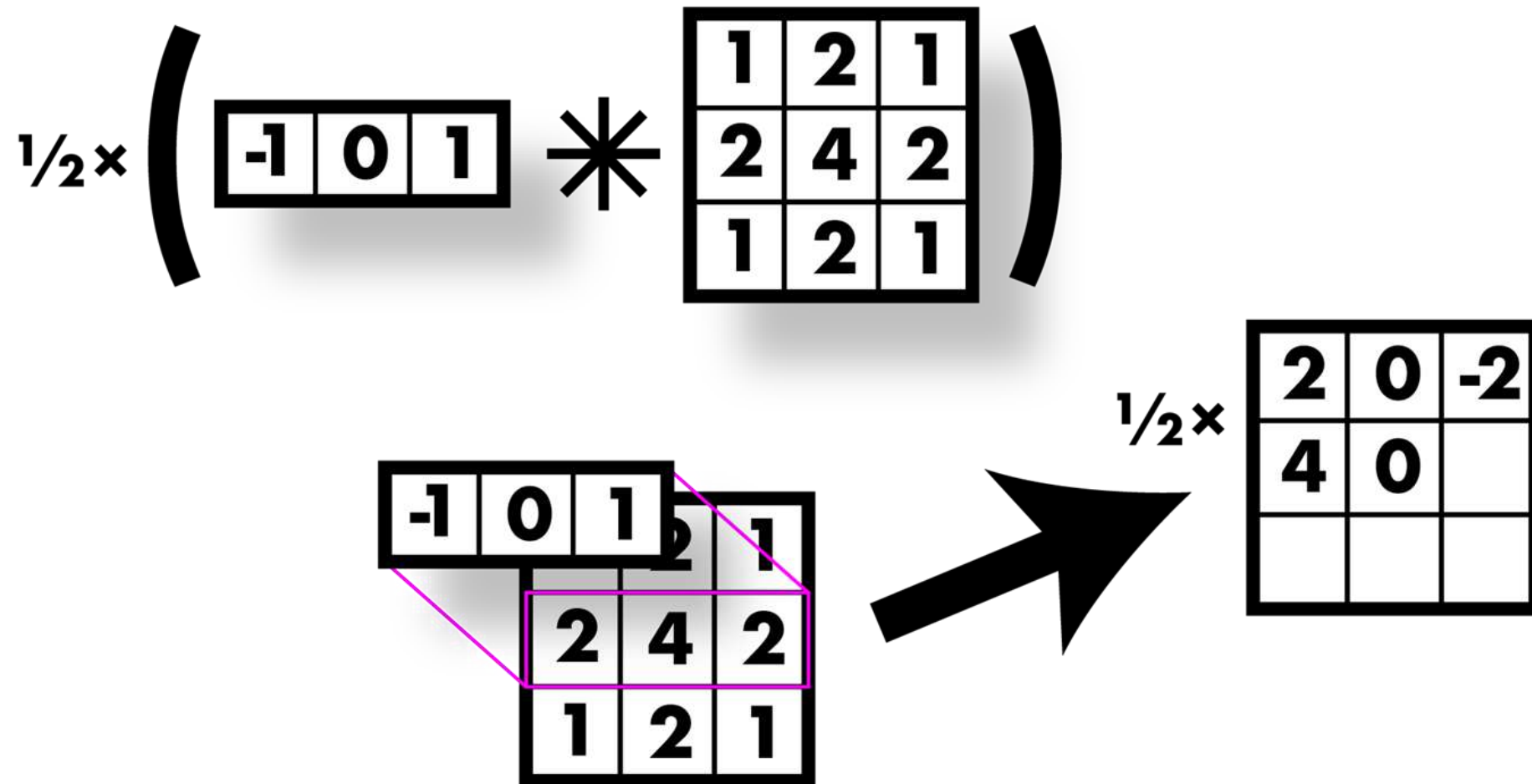


# Smooth first, then derivative

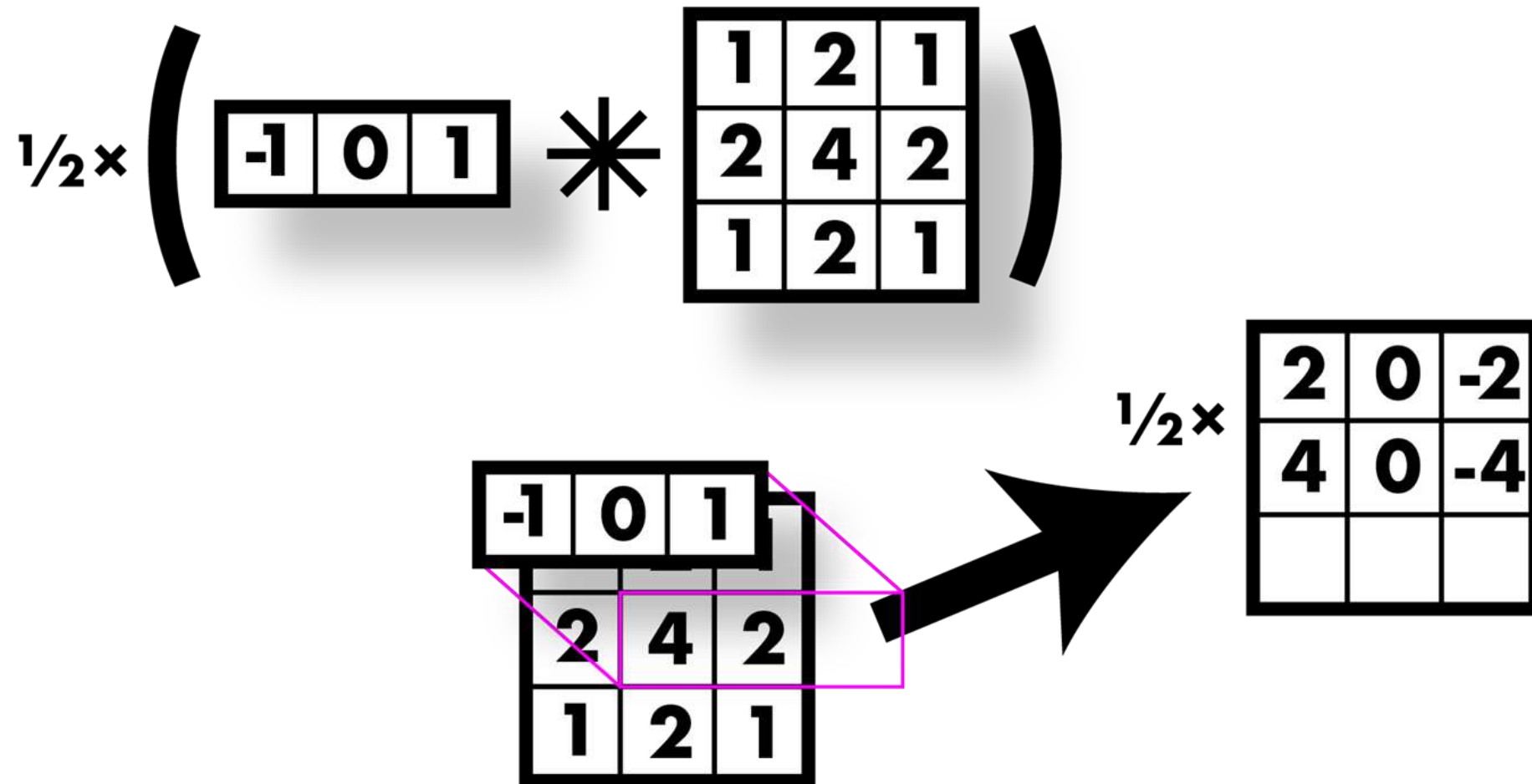




## Smooth first, then derivative

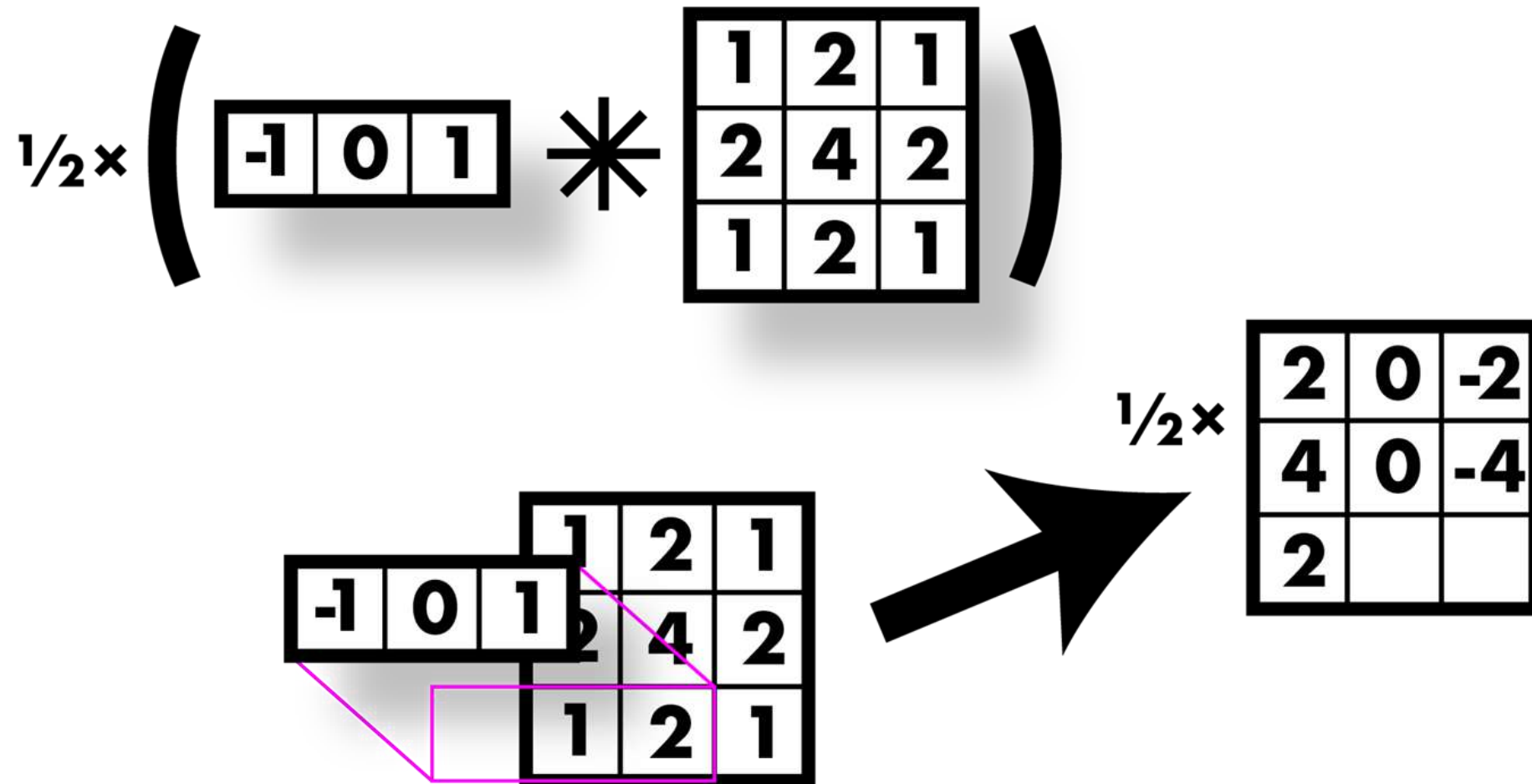


## Smooth first, then derivative





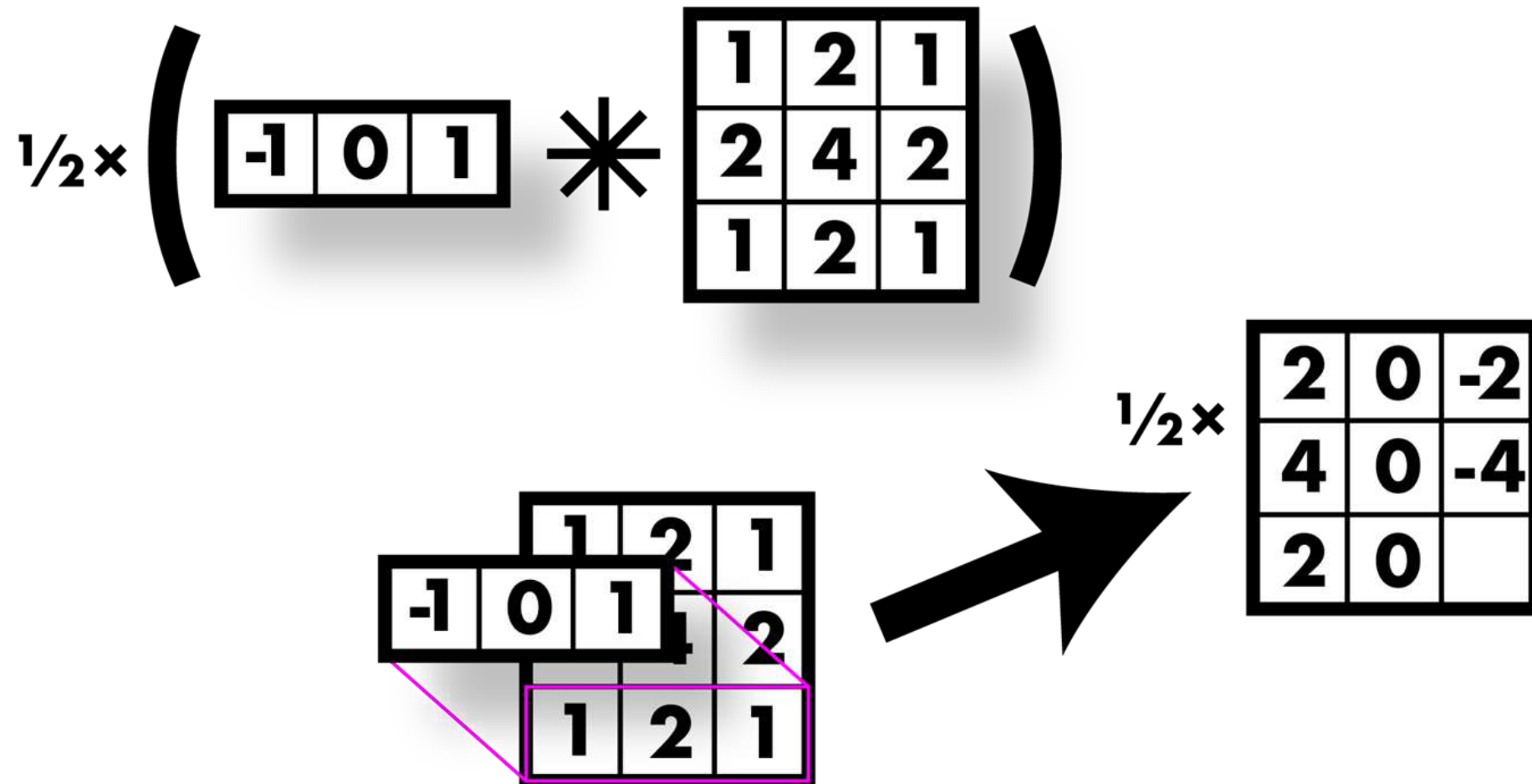
# Smooth first, then derivative





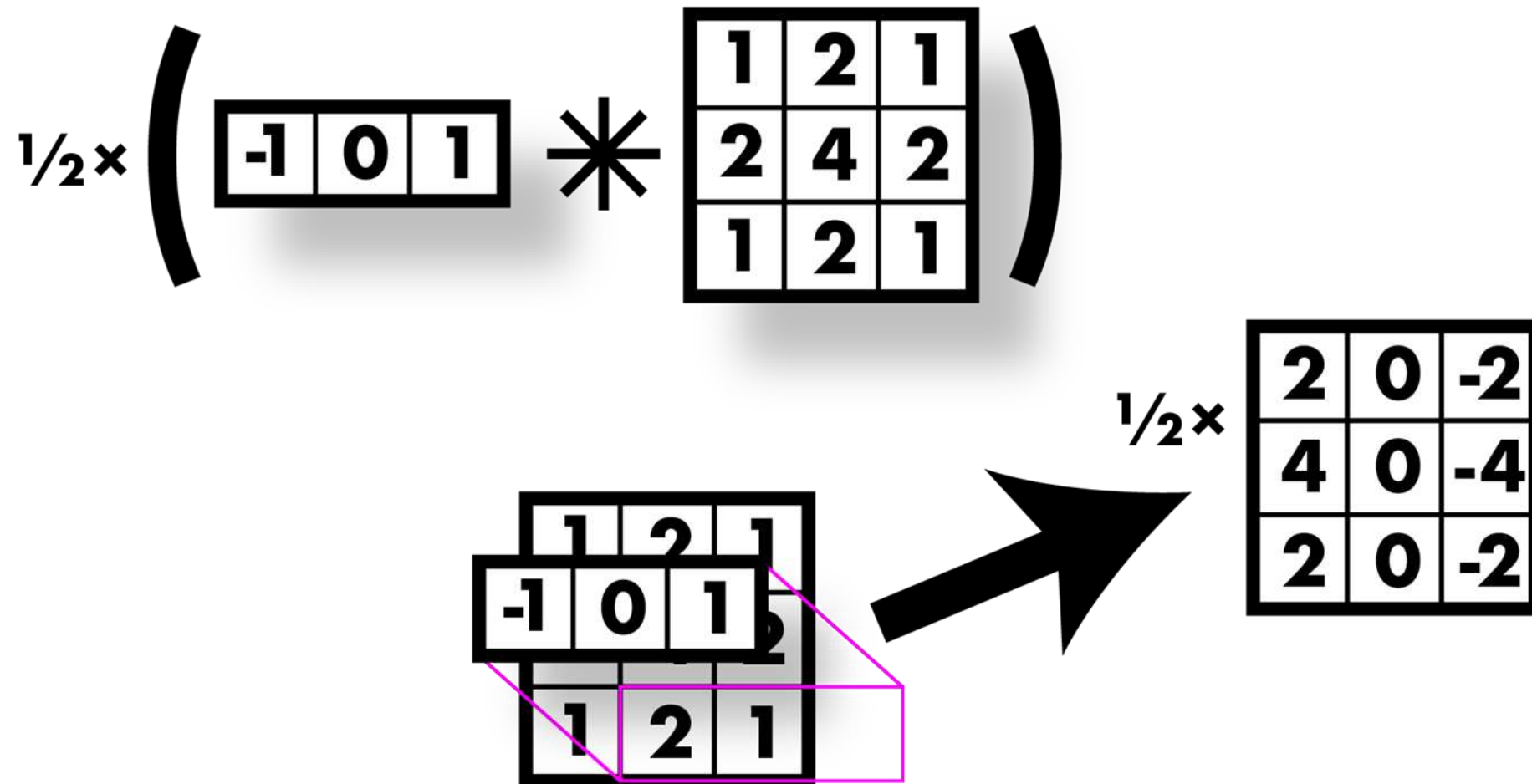


# Smooth first, then derivative



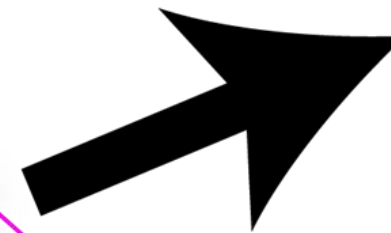
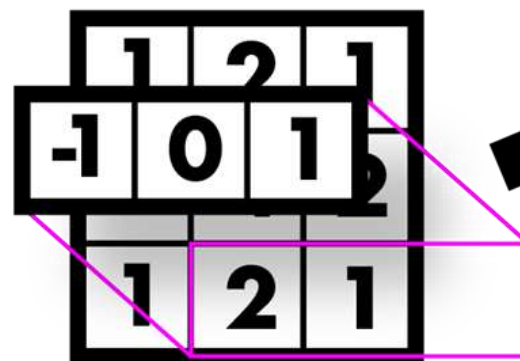
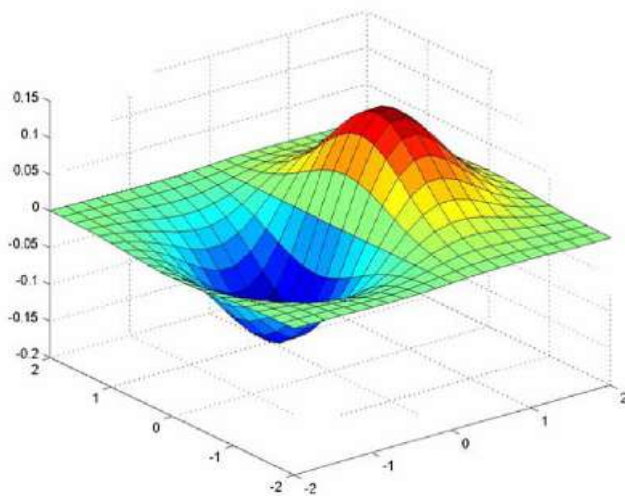


## Smooth first, then derivative



## Sobel filter! Smooth & derivative

$$\frac{1}{2} \times \left( \begin{bmatrix} -1 & 0 & 1 \end{bmatrix} * \begin{bmatrix} 1 & 2 & 1 \\ 2 & 4 & 2 \\ 1 & 2 & 1 \end{bmatrix} \right)$$



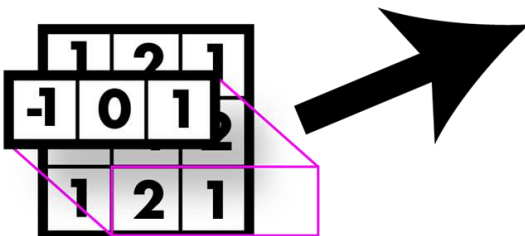
1	0	-1
2	0	-2
1	0	-1

## Image derivatives

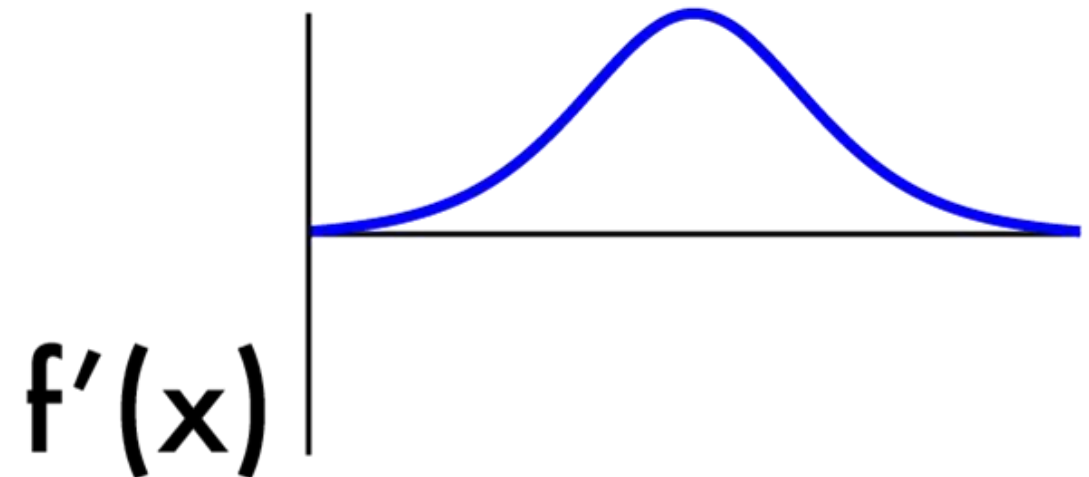
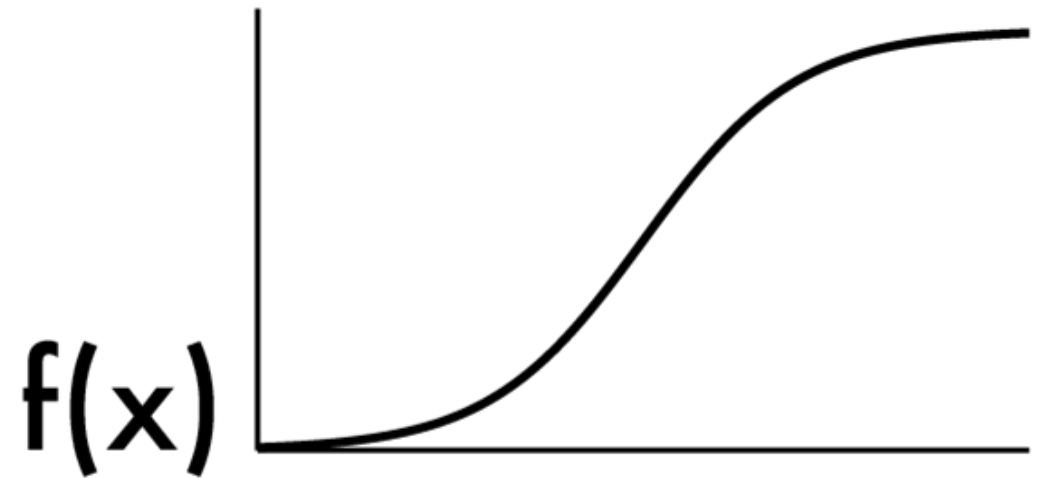
- Recall:

$$f'(a) = \lim_{h \rightarrow 0} \frac{f(a+h) - f(a)}{h}$$

- Want smoothing too!

$$\frac{1}{2} \times \left( \begin{bmatrix} -1 & 0 & 1 \end{bmatrix} * \begin{bmatrix} 1 & 2 & 1 \\ 2 & 4 & 2 \\ 1 & 2 & 1 \end{bmatrix} \right)$$
  


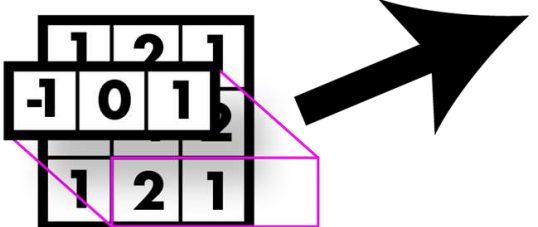
$$\begin{bmatrix} 1 & 0 & -1 \\ 2 & 0 & -2 \\ 1 & 0 & -1 \end{bmatrix}$$

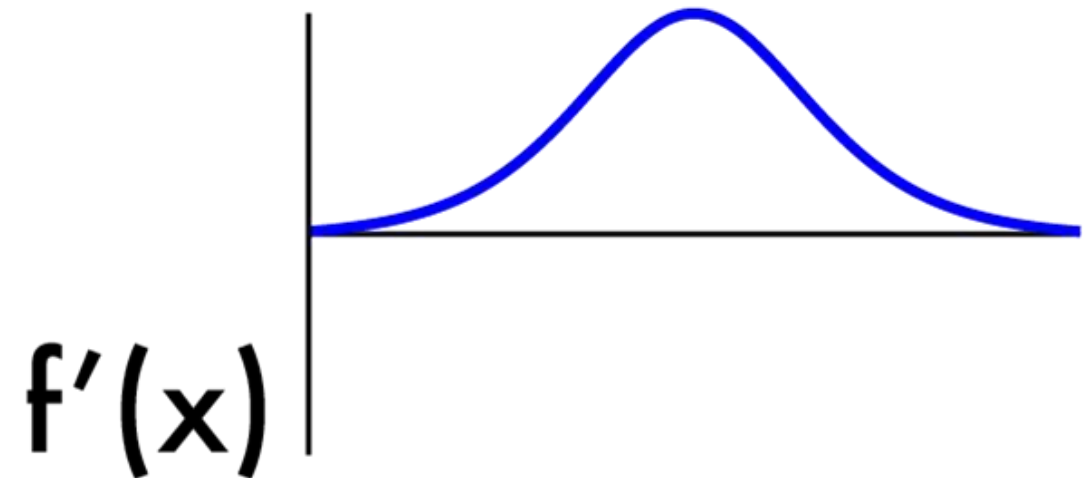
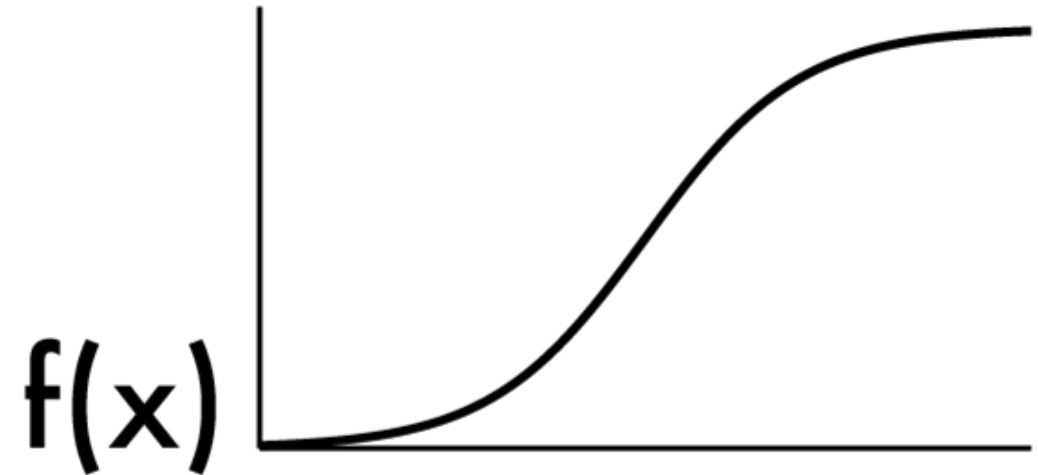




## Finding edges

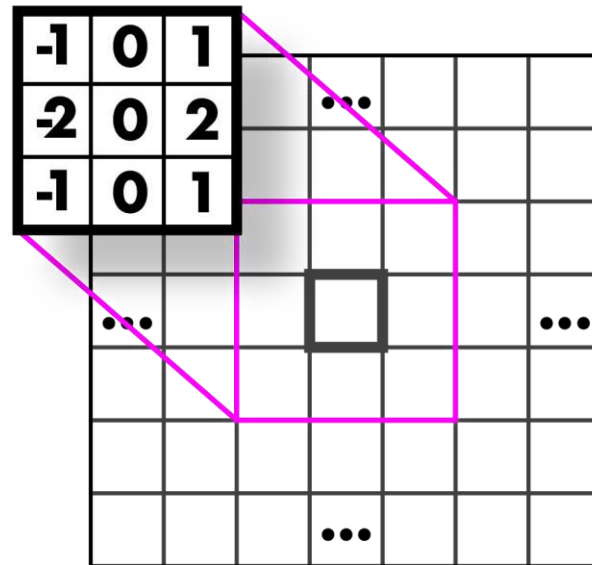
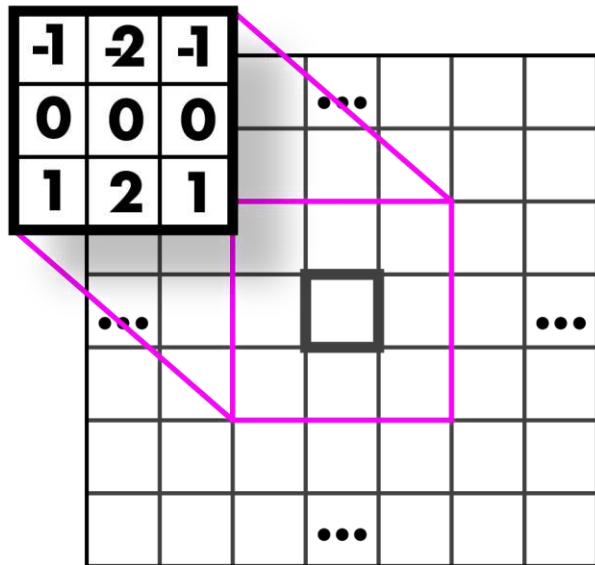
- Could take derivative
- Find high responses
- Sobel filters!
- What about y direction ?

$$\frac{1}{2} \times \left( \begin{bmatrix} -1 & 0 & 1 \end{bmatrix} * \begin{bmatrix} 1 & 2 & 1 \\ 2 & 4 & 2 \\ 1 & 2 & 1 \end{bmatrix} \right)$$
  


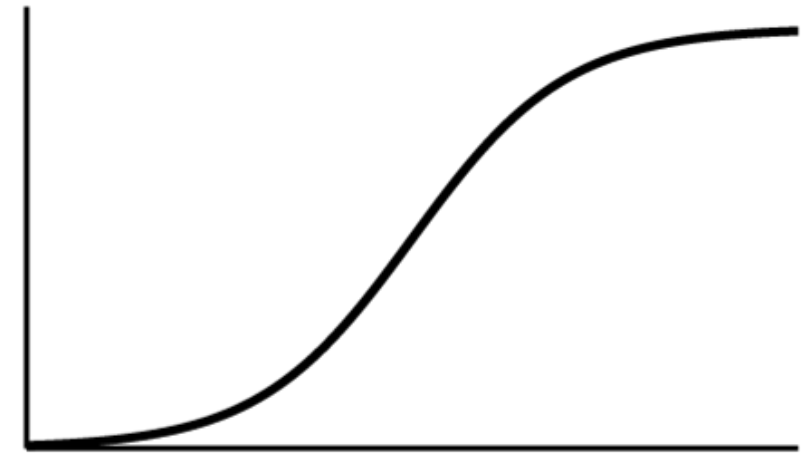


## Finding edges

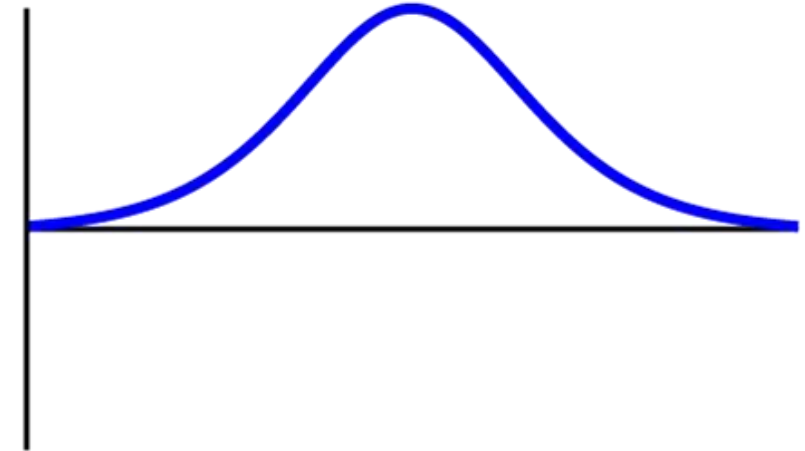
- Could take derivative
- Find high responses
- Sobel filters!
- Let's stop a moment and get some basics



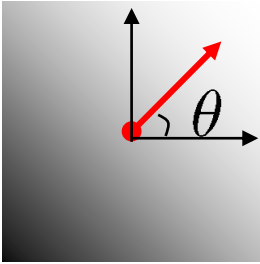
$f(x)$



$f'(x)$



# Simplest image gradient



$$\nabla f = \left[ \frac{\partial f}{\partial x}, \frac{\partial f}{\partial y} \right]$$

$$\frac{\partial f}{\partial x} = f(x + 1, y) - f(x, y)$$

Likewise for  $df/dy$

The **gradient direction** is  $\theta = \tan^{-1} \left( \frac{\partial f}{\partial y} / \frac{\partial f}{\partial x} \right)$

How does this relate to the direction of the edge? -Perpendicular

The *edge strength* is given by the **gradient magnitude**

$$\|\nabla f\| = \sqrt{\left(\frac{\partial f}{\partial x}\right)^2 + \left(\frac{\partial f}{\partial y}\right)^2}$$

## Sobel filters

$$g_x = \begin{bmatrix} 1 & 0 & -1 \\ 2 & 0 & -2 \\ 1 & 0 & -1 \end{bmatrix} \quad g_y = \begin{bmatrix} 1 & 2 & 1 \\ 0 & 0 & 0 \\ -1 & -2 & -1 \end{bmatrix}$$

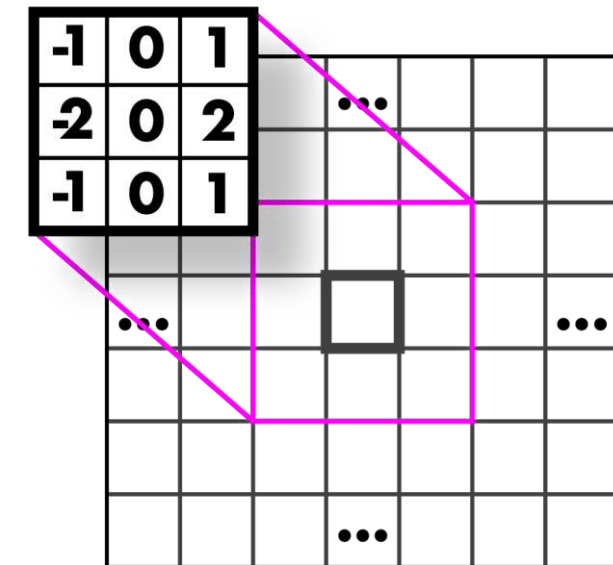
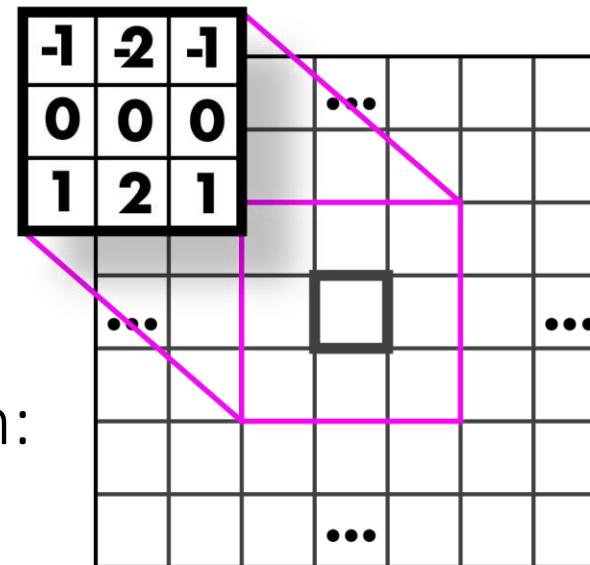
Magnitude:

$$g = \sqrt{g_x^2 + g_y^2}$$

Orientation:

$$\Theta = \tan^{-1} \left( \frac{g_y}{g_x} \right)$$

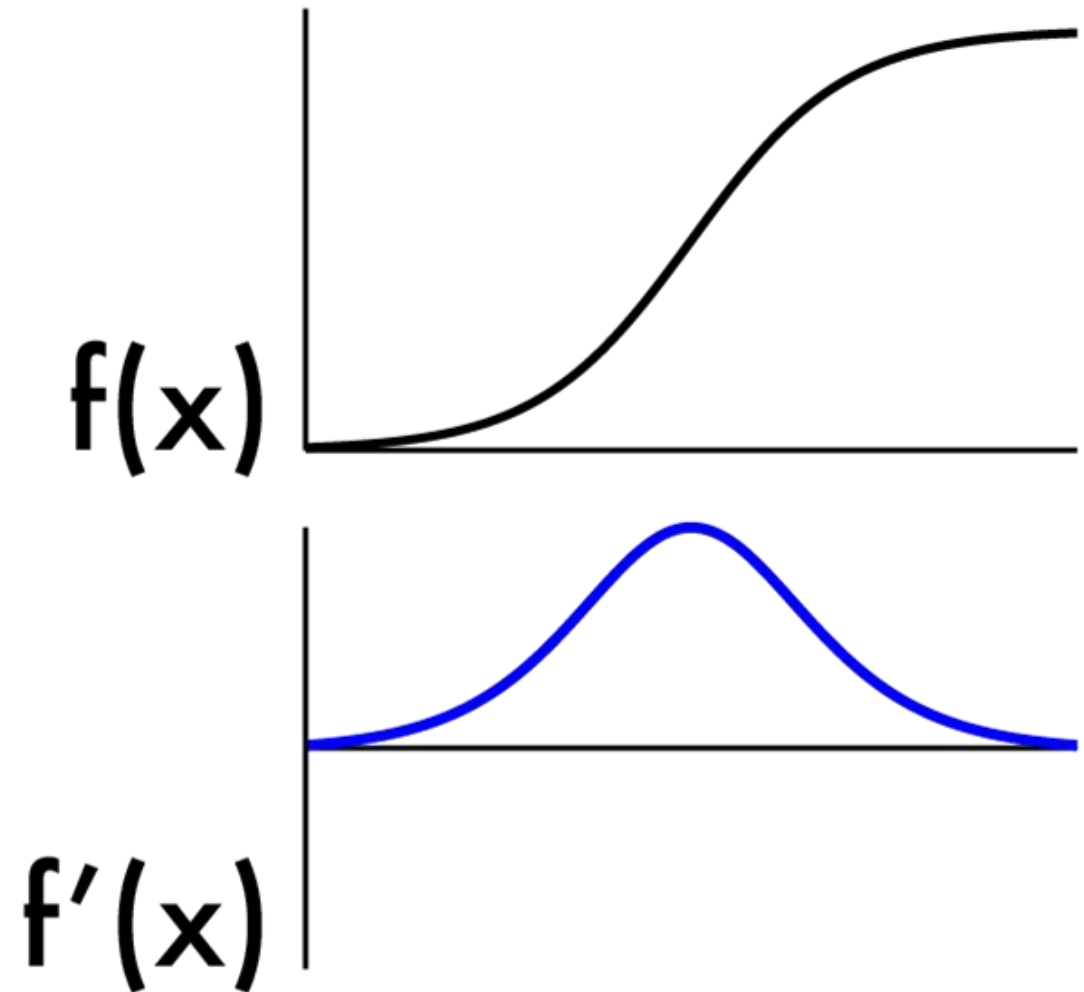
We can change the sign:





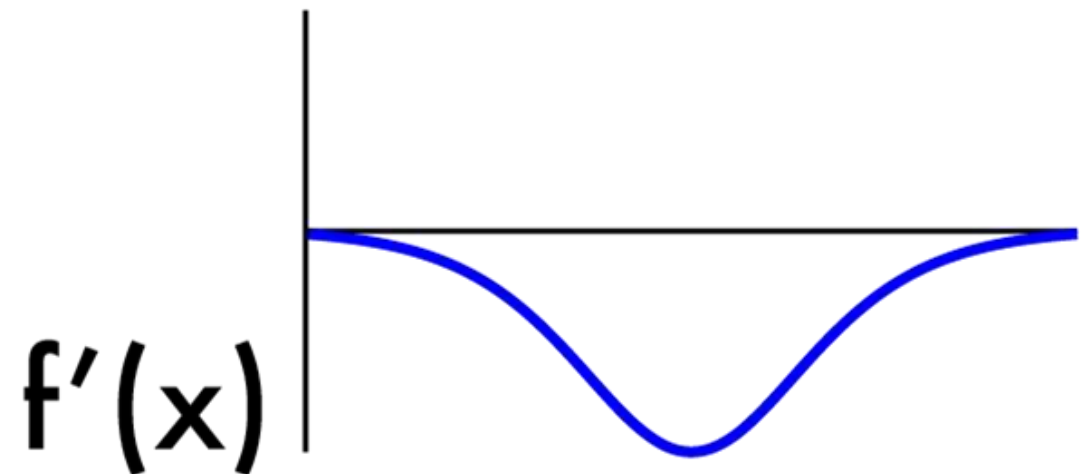
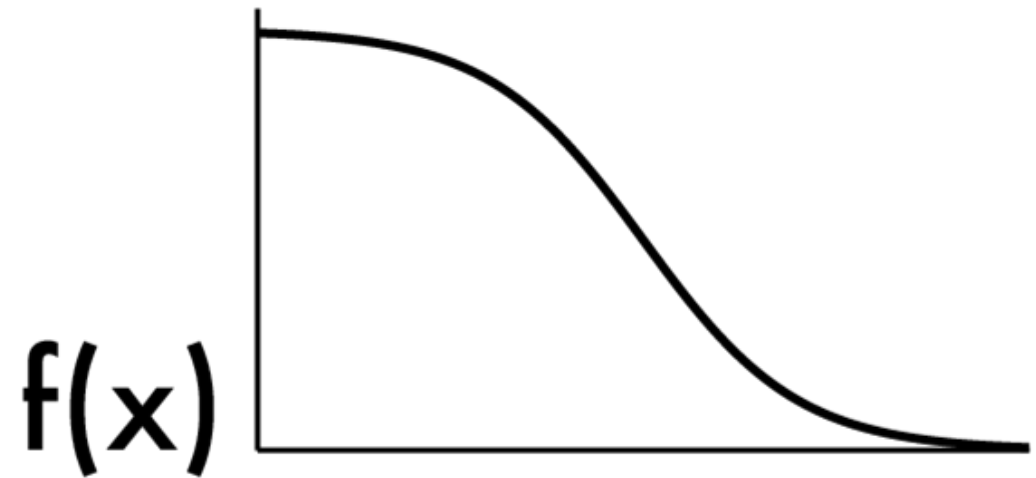
# Finding edges

- Could take derivative
- Find high responses
- Sobel filters!
- But...



# Finding edges

- Could take derivative
- Find high responses
- Sobel filters!
- But...
- Edges go both ways
- Want to find extrema



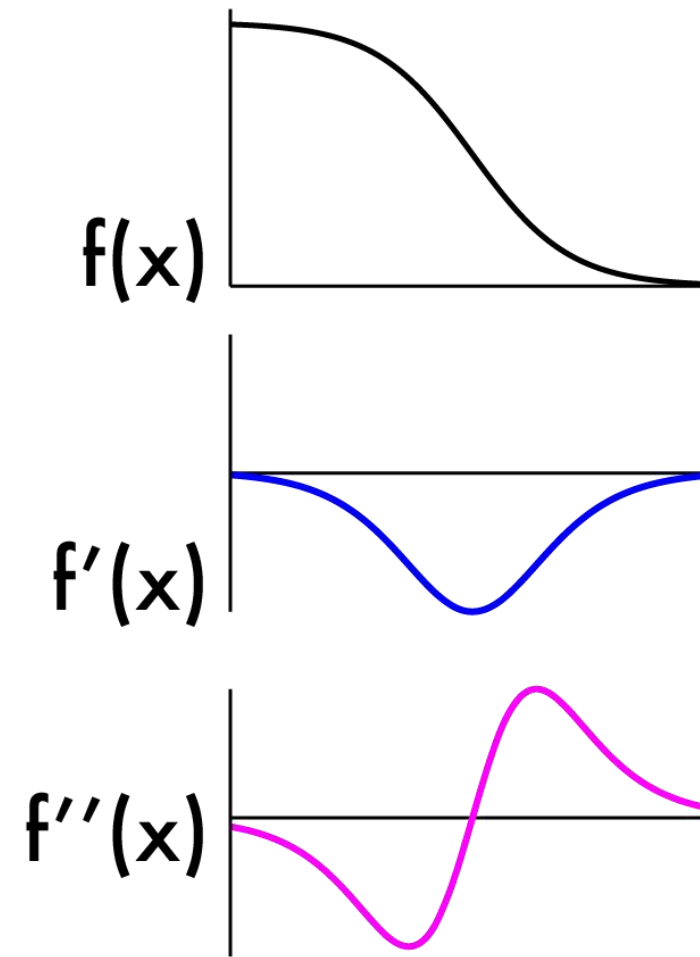
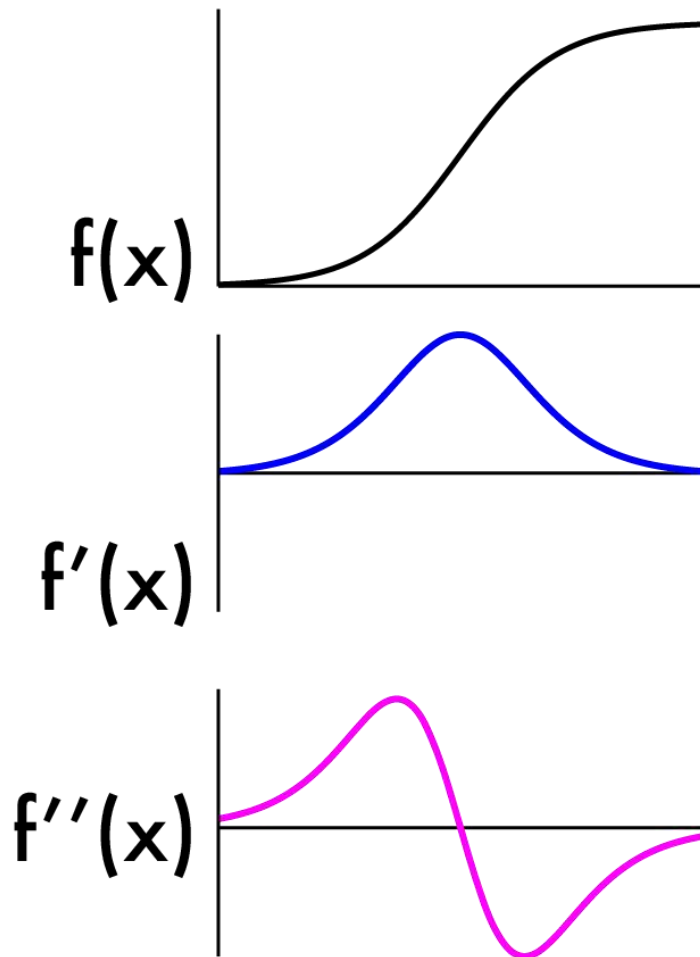
# Today's Agenda

- What can we do with convolutions
- What is an edge – image derivatives
- Sobel filters
- Laplacian filters
- Difference of Gaussian filters
- Canny edge detection



# 2nd derivative!

- Crosses zero at extrema





## Laplacian (2nd derivative)!

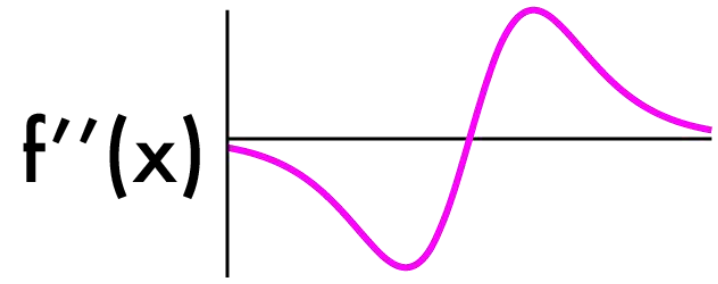
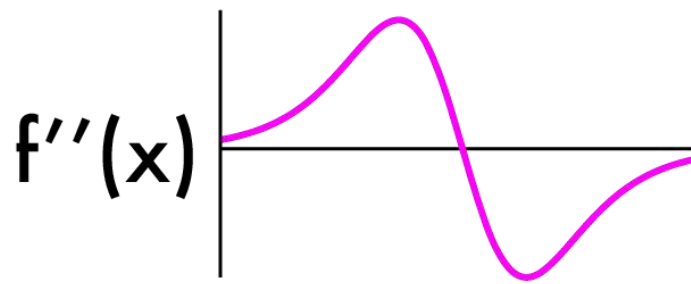
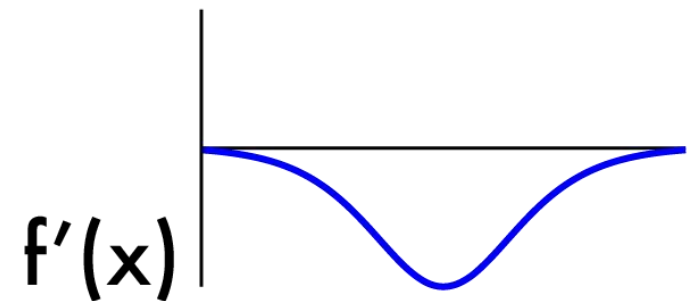
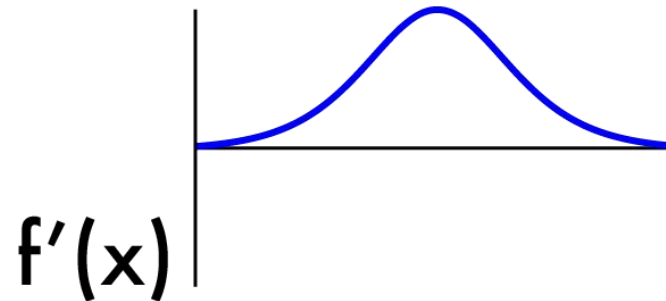
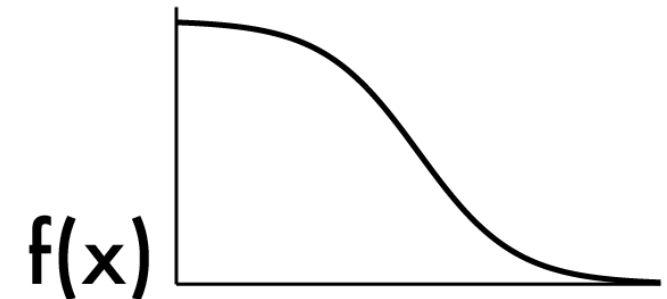
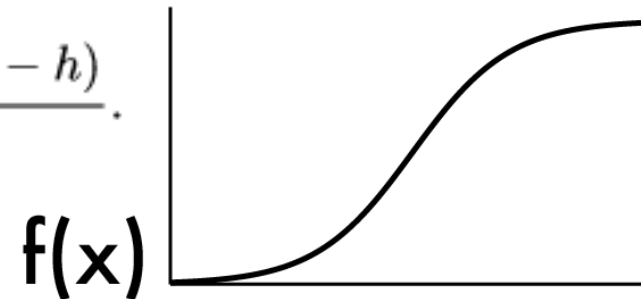
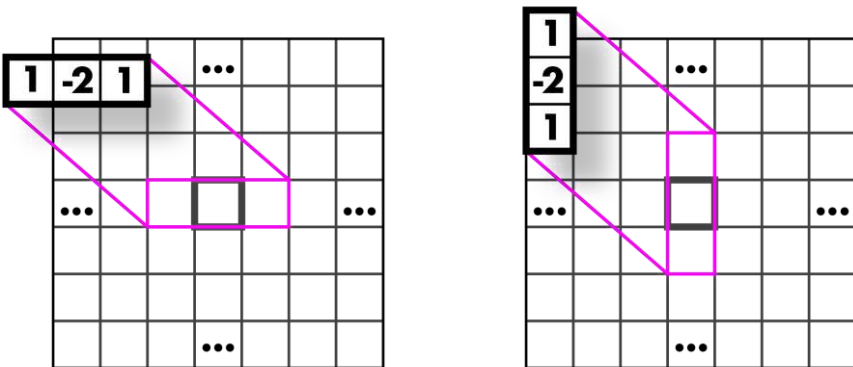
- Crosses zero at extrema
- Recall:

$$f''(x) = \lim_{h \rightarrow 0} \frac{f(x+h) - 2f(x) + f(x-h)}{h^2}$$

- Laplacian:

$$\Delta f = \frac{\partial^2 f}{\partial x^2} + \frac{\partial^2 f}{\partial y^2}$$

- Again, have to estimate  $f''(x)$ :



# Laplacians

- Laplacian:

$$\Delta f = \frac{\partial^2 f}{\partial x^2} + \frac{\partial^2 f}{\partial y^2}$$

# Laplacians

- Laplacian:

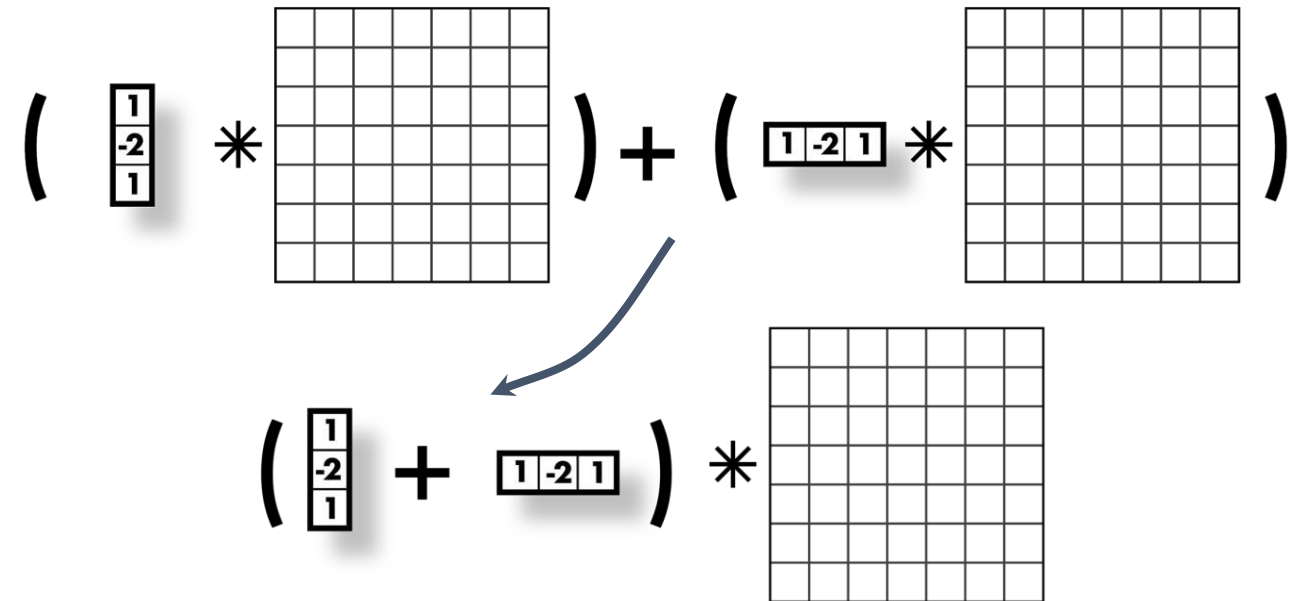
$$\Delta f = \frac{\partial^2 f}{\partial x^2} + \frac{\partial^2 f}{\partial y^2}$$

$$\left( \begin{bmatrix} 1 \\ -2 \\ 1 \end{bmatrix} * \begin{array}{|c|} \hline \square \\ \square \\ \square \\ \square \\ \square \\ \square \\ \square \\ \square \\ \square \\ \square \\ \hline \end{array} \right) + \left( \begin{bmatrix} 1 & -2 & 1 \end{bmatrix} * \begin{array}{|c|} \hline \square \\ \square \\ \square \\ \square \\ \square \\ \square \\ \square \\ \square \\ \square \\ \square \\ \hline \end{array} \right)$$

# Laplacians

- Laplacian:

$$\Delta f = \frac{\partial^2 f}{\partial x^2} + \frac{\partial^2 f}{\partial y^2}$$

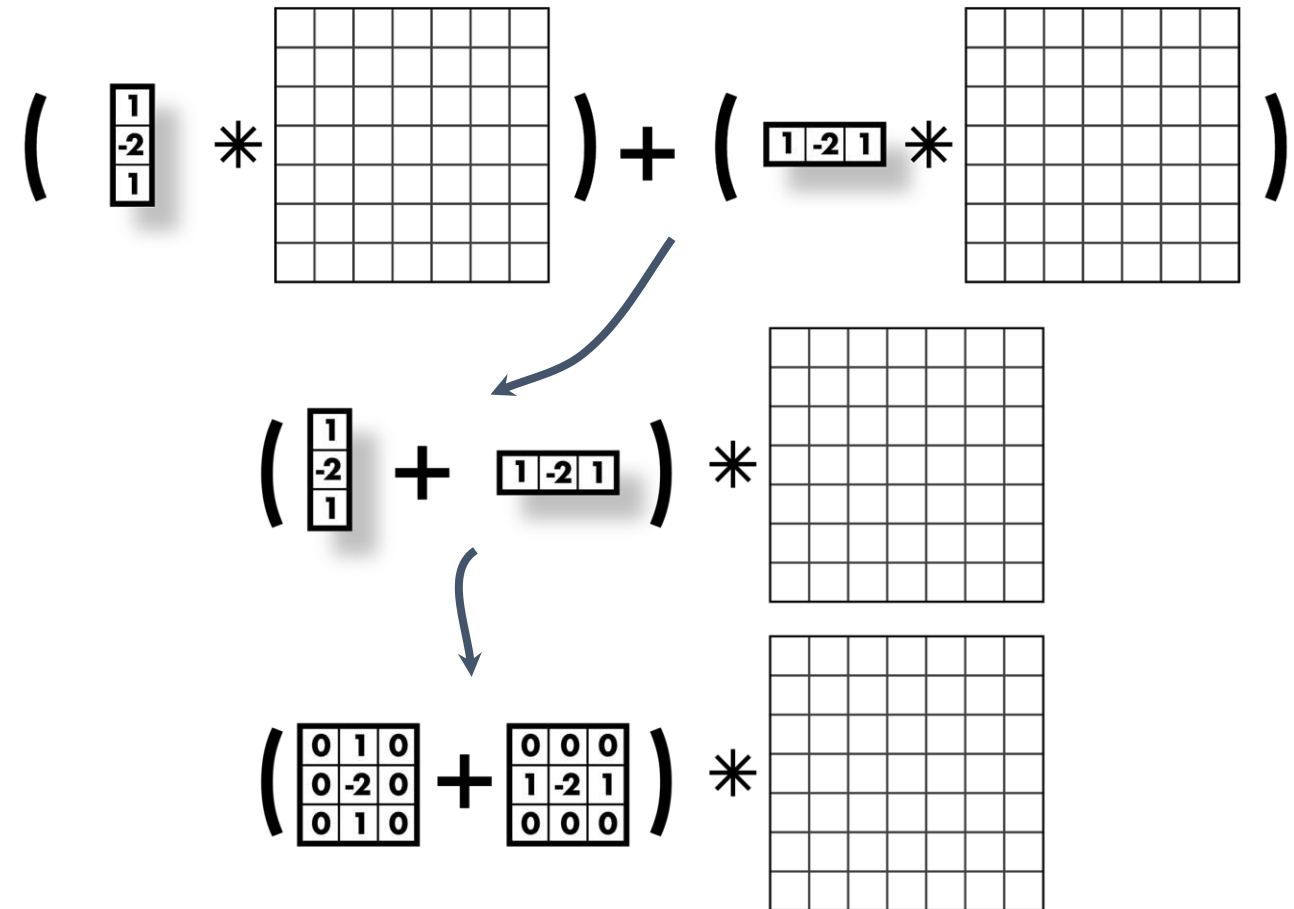




# Laplacians

- Laplacian:

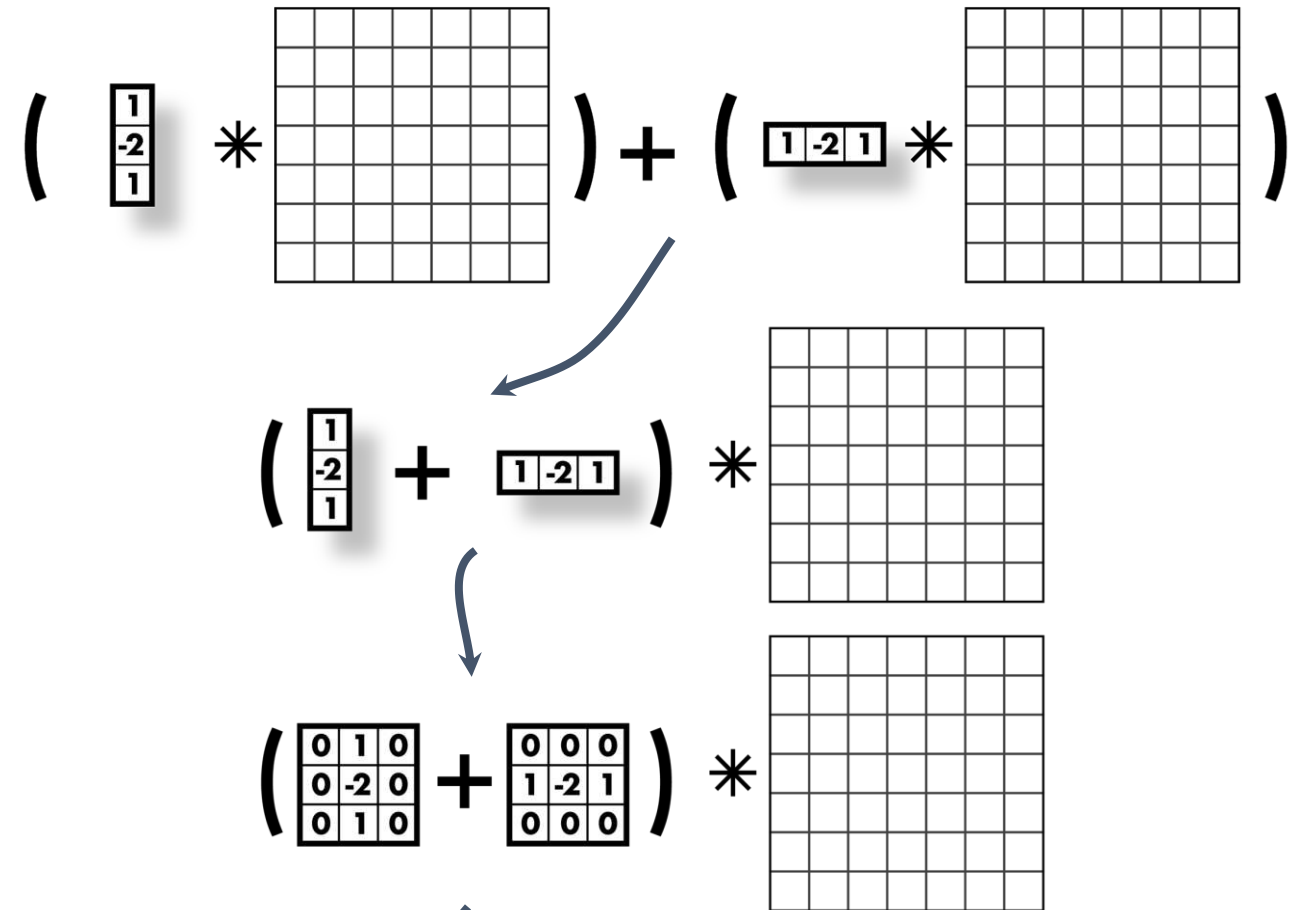
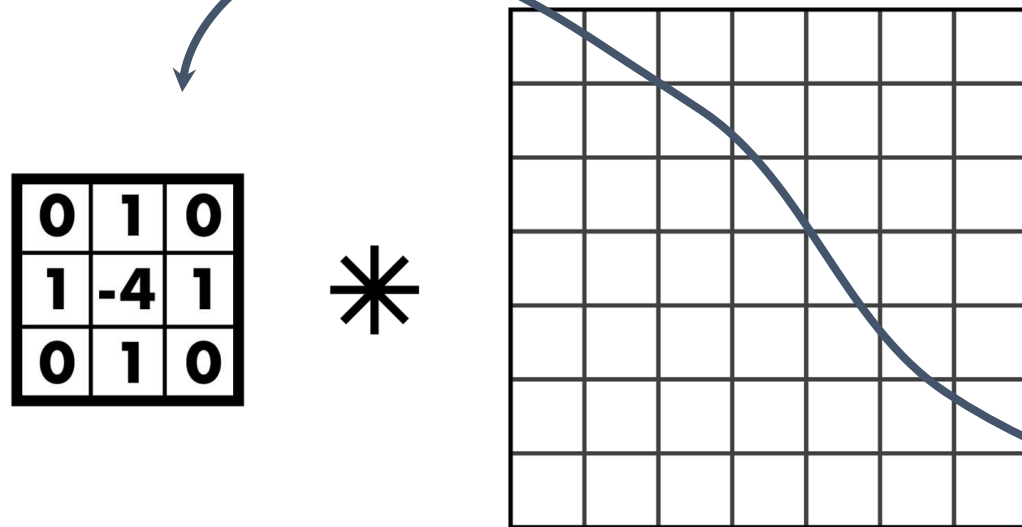
$$\Delta f = \frac{\partial^2 f}{\partial x^2} + \frac{\partial^2 f}{\partial y^2}$$



## Laplacians

- Laplacian:

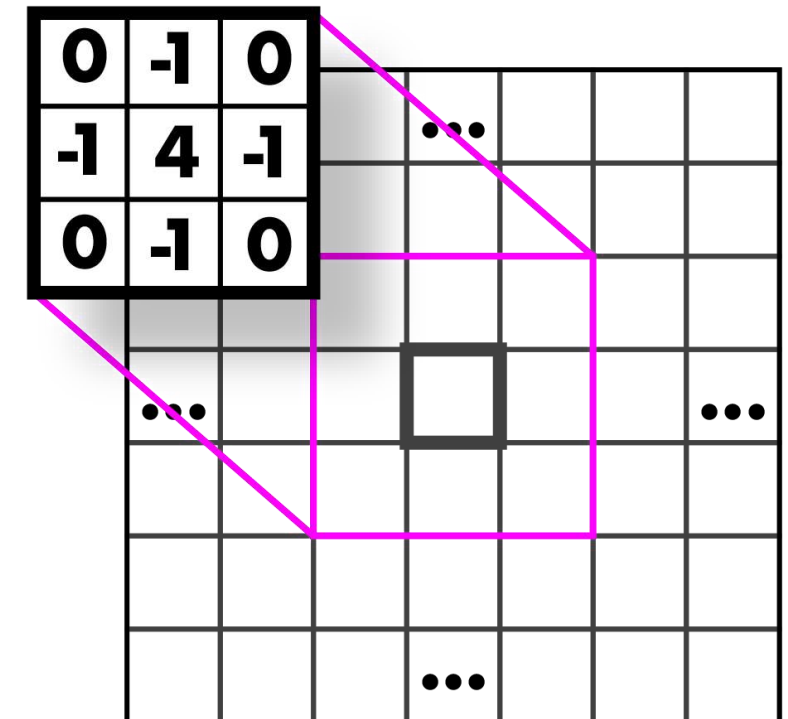
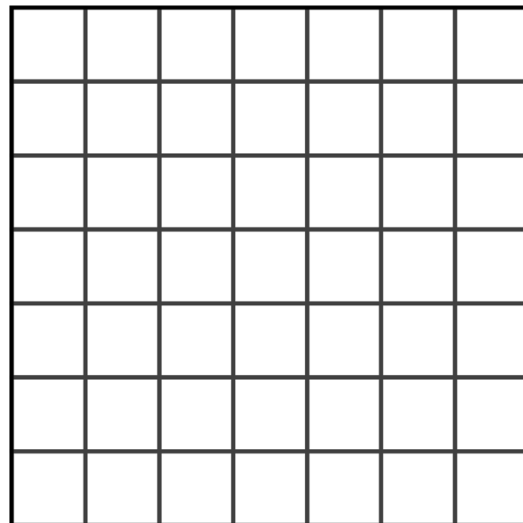
$$\Delta f = \frac{\partial^2 f}{\partial x^2} + \frac{\partial^2 f}{\partial y^2}$$



# Laplacians

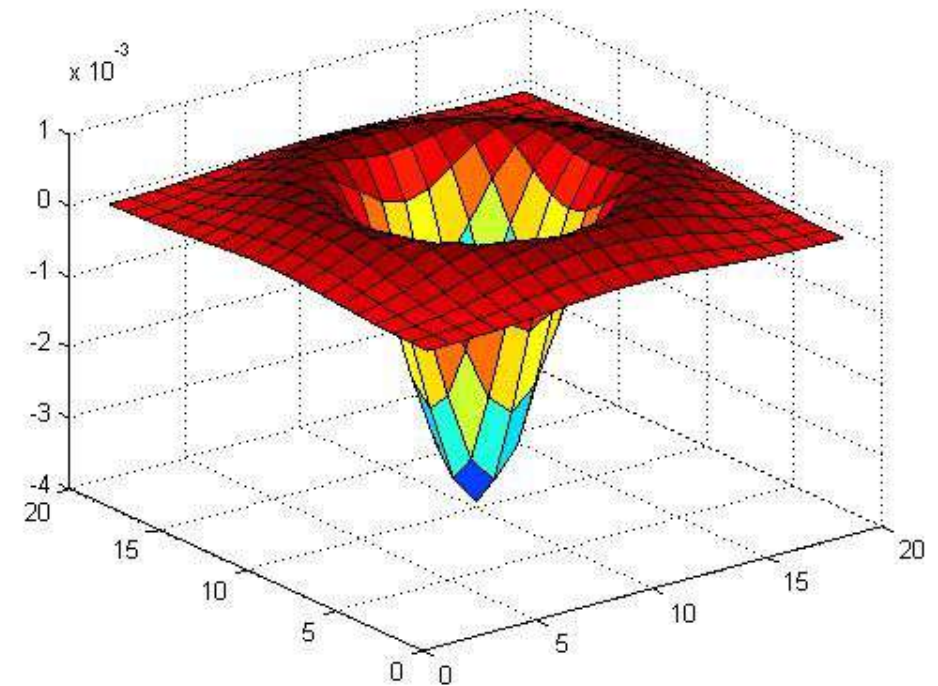
- Laplacian:
  - $\Delta f = \frac{\partial^2 f}{\partial x^2} + \frac{\partial^2 f}{\partial y^2}$
- Negative Laplacian, -4 in middle
- Positive Laplacian --->

0	1	0
1	-4	1
0	1	0



# Laplacians also sensitive to noise

- Again, use gaussian smoothing
- Can just use one kernel since convs commute
- Laplacian of Gaussian, LoG
- Can get good approx. with 5x5 - 9x9 kernels





# Today's Agenda

- What can we do with convolutions
- What is an edge – image derivatives
- Sobel filters
- Laplacian filters
- Difference of Gaussian filters
- Canny edge detection

# Another edge detector

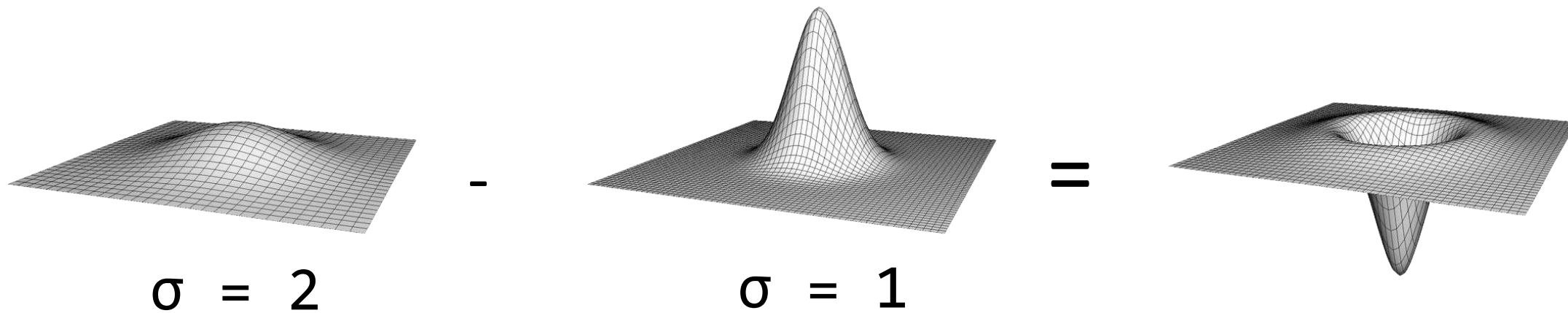
- Image is a function
  - Has high frequency and low frequency components
  - Think in terms of fourier transform
- Edges are high frequency changes
- Maybe we want to find edges of a specific size (i.e. specific frequency)

# Difference of Gaussian (DoG)

- Gaussian is a low pass filter
- Strongly reduce components with frequency  $f > \sigma$
- $(g * I)$  low frequency components
- $I - (g * I)$  high frequency components
- $g(\sigma_1) * I - g(\sigma_2) * I$ 
  - Components in between these frequencies
- $g(\sigma_1) * I - g(\sigma_2) * I = [g(\sigma_1) - g(\sigma_2)] * I$

# Difference of Gaussian (DoG)

$$- g(\sigma_1) * I - g(\sigma_2) * I = [g(\sigma_1) - g(\sigma_2)] * I$$

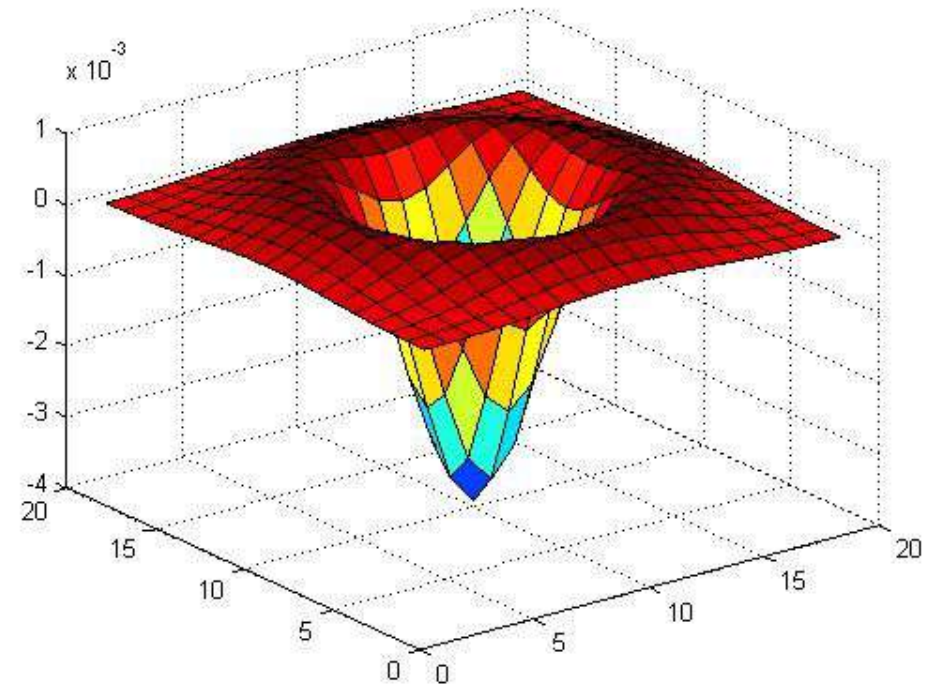
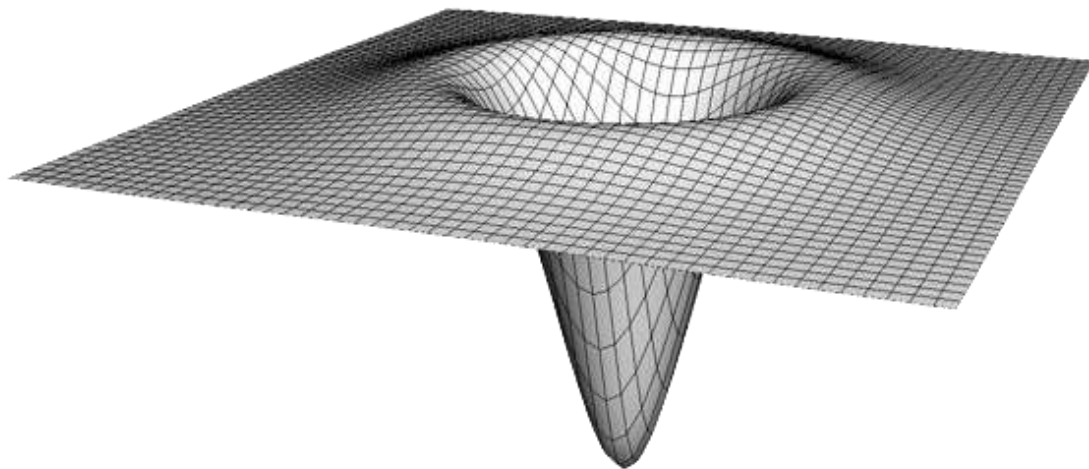






# Difference of Gaussian (DoG)

- $g(\sigma_1) * I - g(\sigma_2) * I = [g(\sigma_1) - g(\sigma_2)] * I$
- This looks a lot like our LoG!
- (not actually the same but similar)





## DoG (1 - 0)





## DoG (2 - 1)





## DoG (3 - 2)





## DoG (4 - 3)

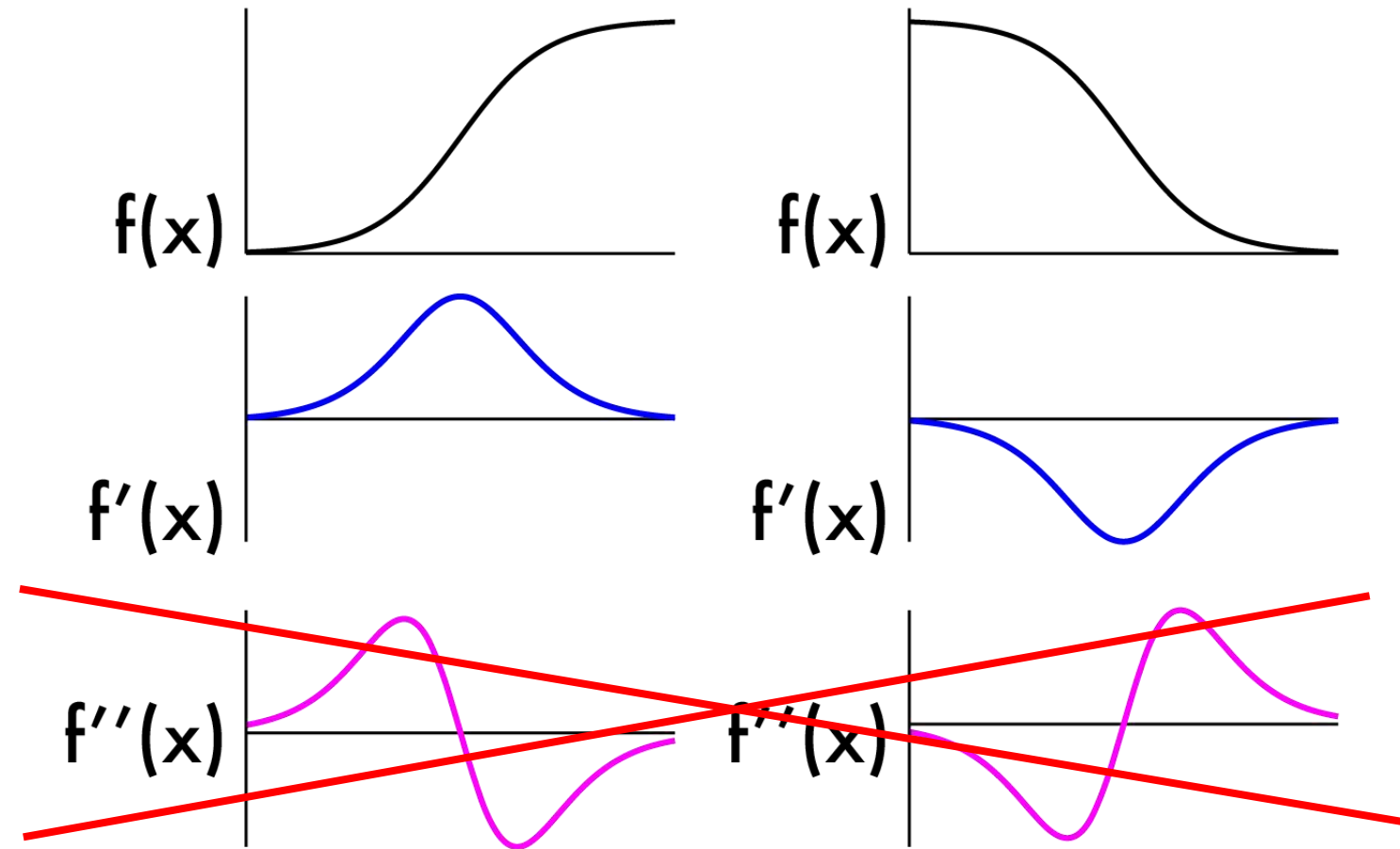


# Today's Agenda

- What can we do with convolutions
- What is an edge – image derivatives
- Sobel filters
- Laplacian filters
- Difference of Gaussian filters
- Canny edge detection

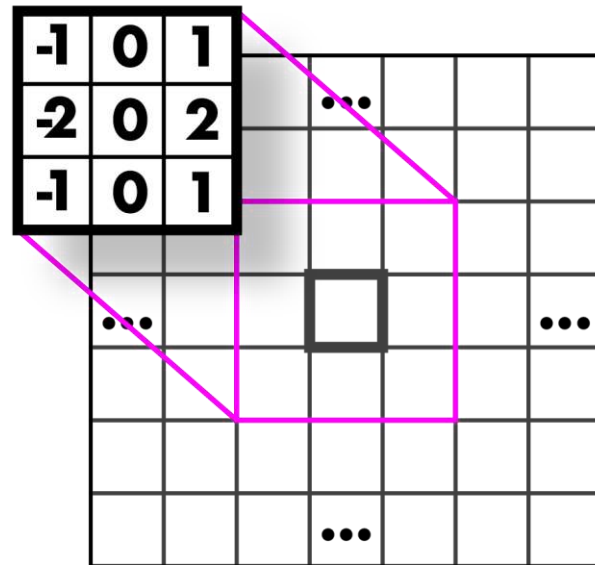
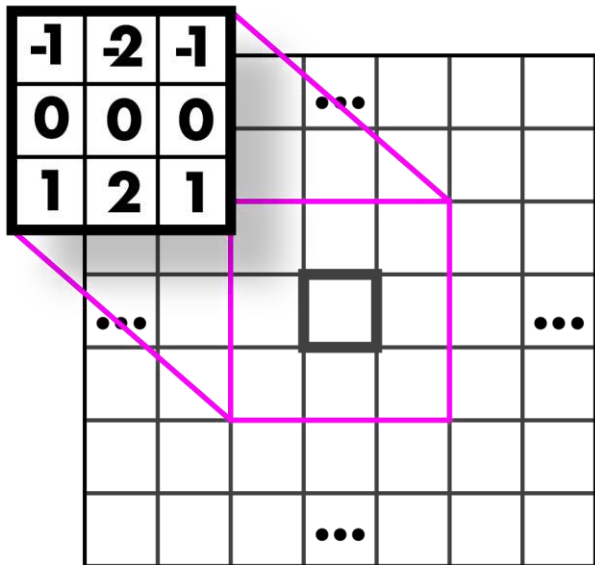
# Another approach: gradient magnitude

- Don't need 2nd derivatives
- Just use magnitude of gradient



## Another approach: gradient magnitude

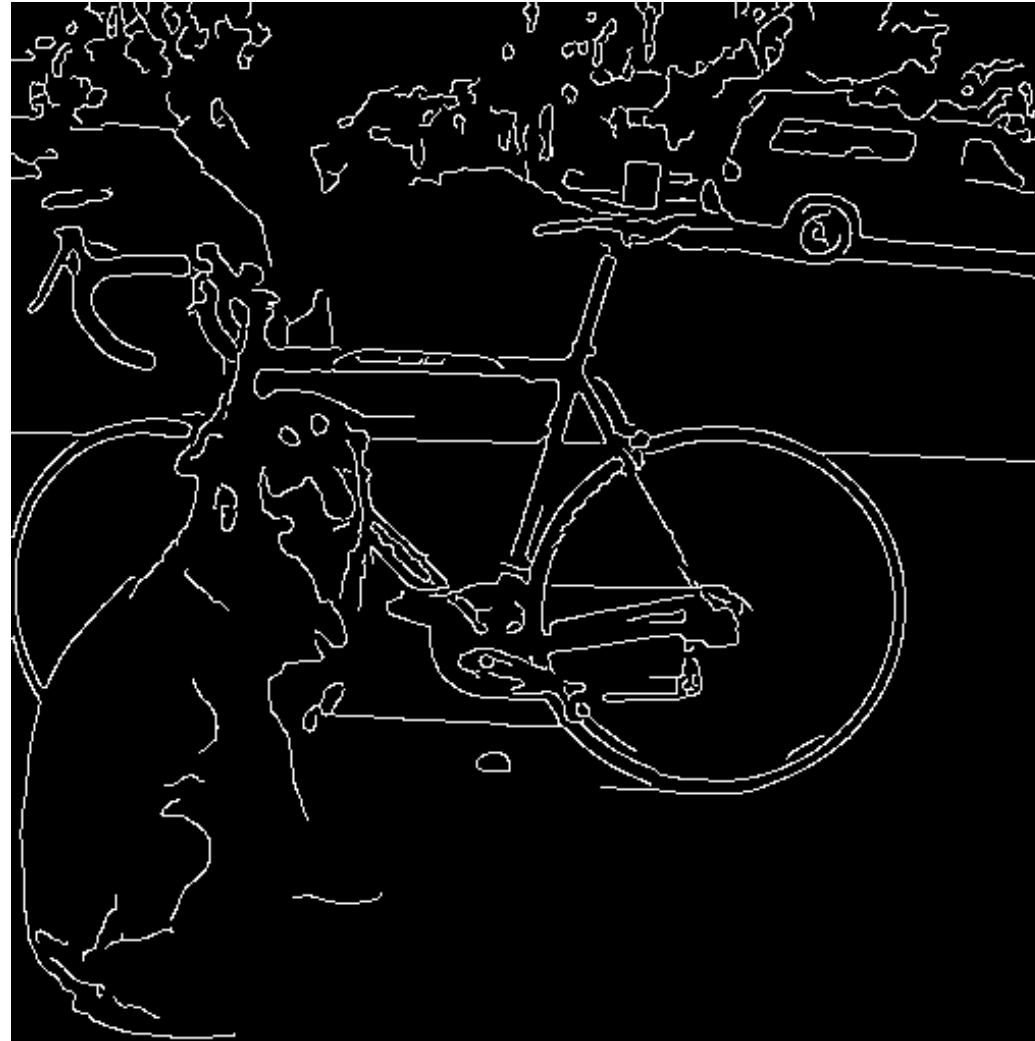
- Don't need 2nd derivatives
- Just use magnitude of gradient
- Are we done? No!







## What we really want: line drawing



# Canny Edge Detection

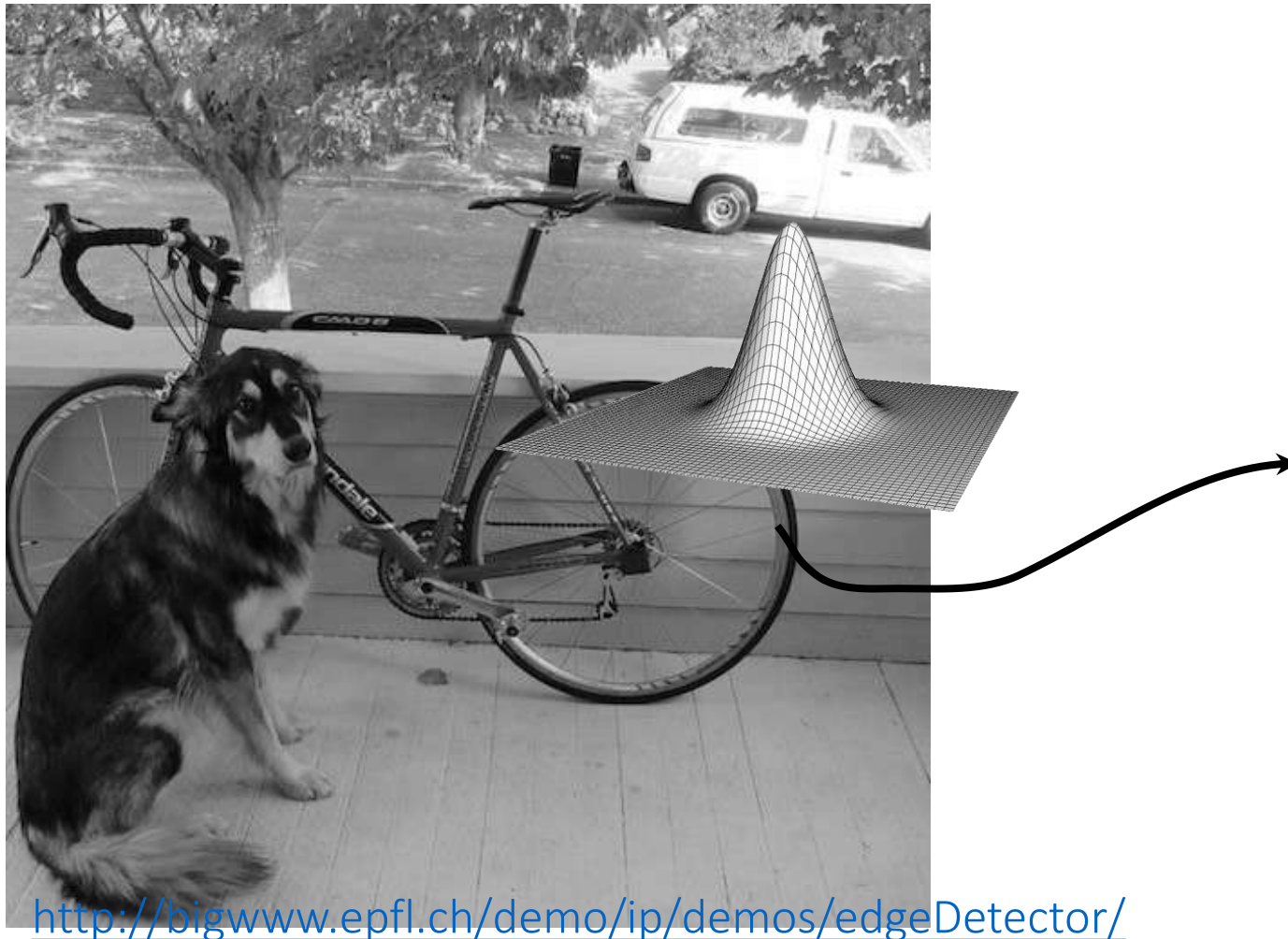
## Algorithm:

- Smooth image (only want “real” edges, not noise)
- Calculate gradient direction and magnitude
- Non-maximum suppression perpendicular to edge
- Threshold into strong, weak, no edge
- Connect together components

<http://bigwww.epfl.ch/demo/ip/demos/edgeDetector/>

# Smooth image

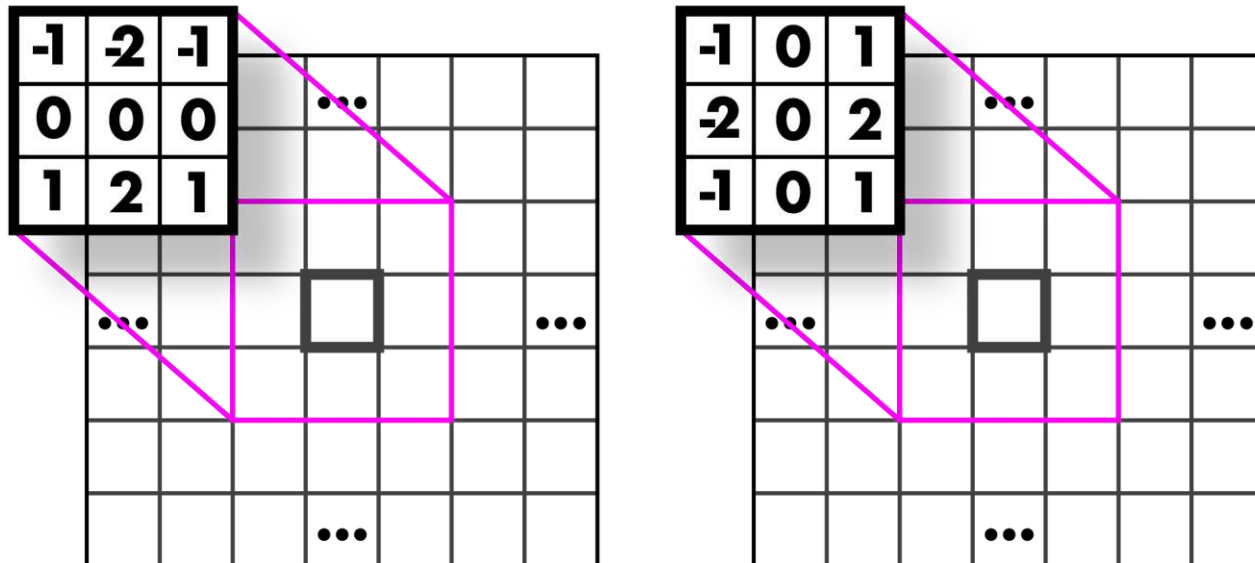
- You know how to do this, gaussians!





# Gradient magnitude and direction

- Sobel filter



<http://bigwww.epfl.ch/demo/ip/demos/edgeDetector/>

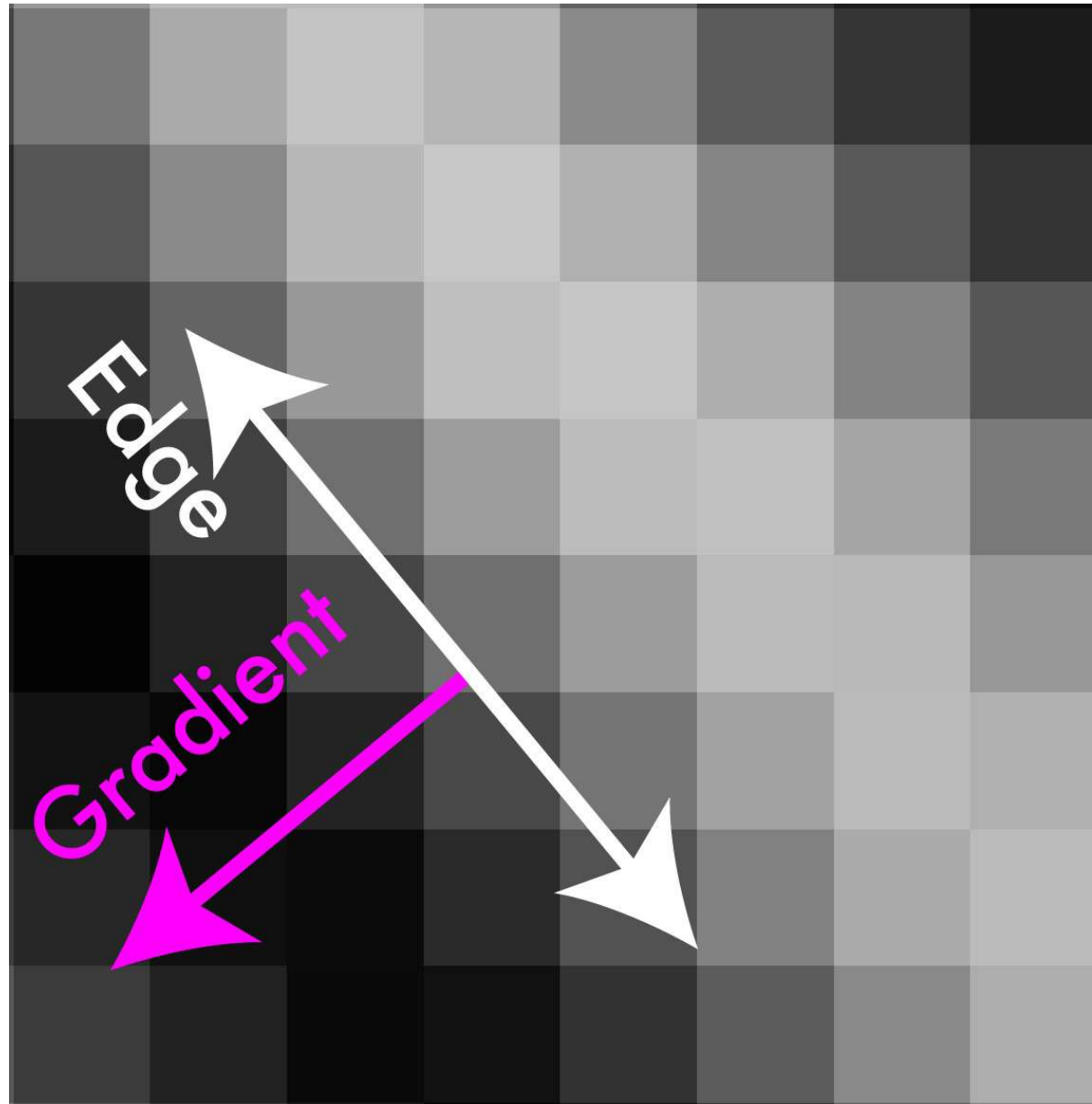
# Non-maximum suppression

- Want single pixel edges, not thick blurry lines
- Need to check nearby pixels
- See if response is highest

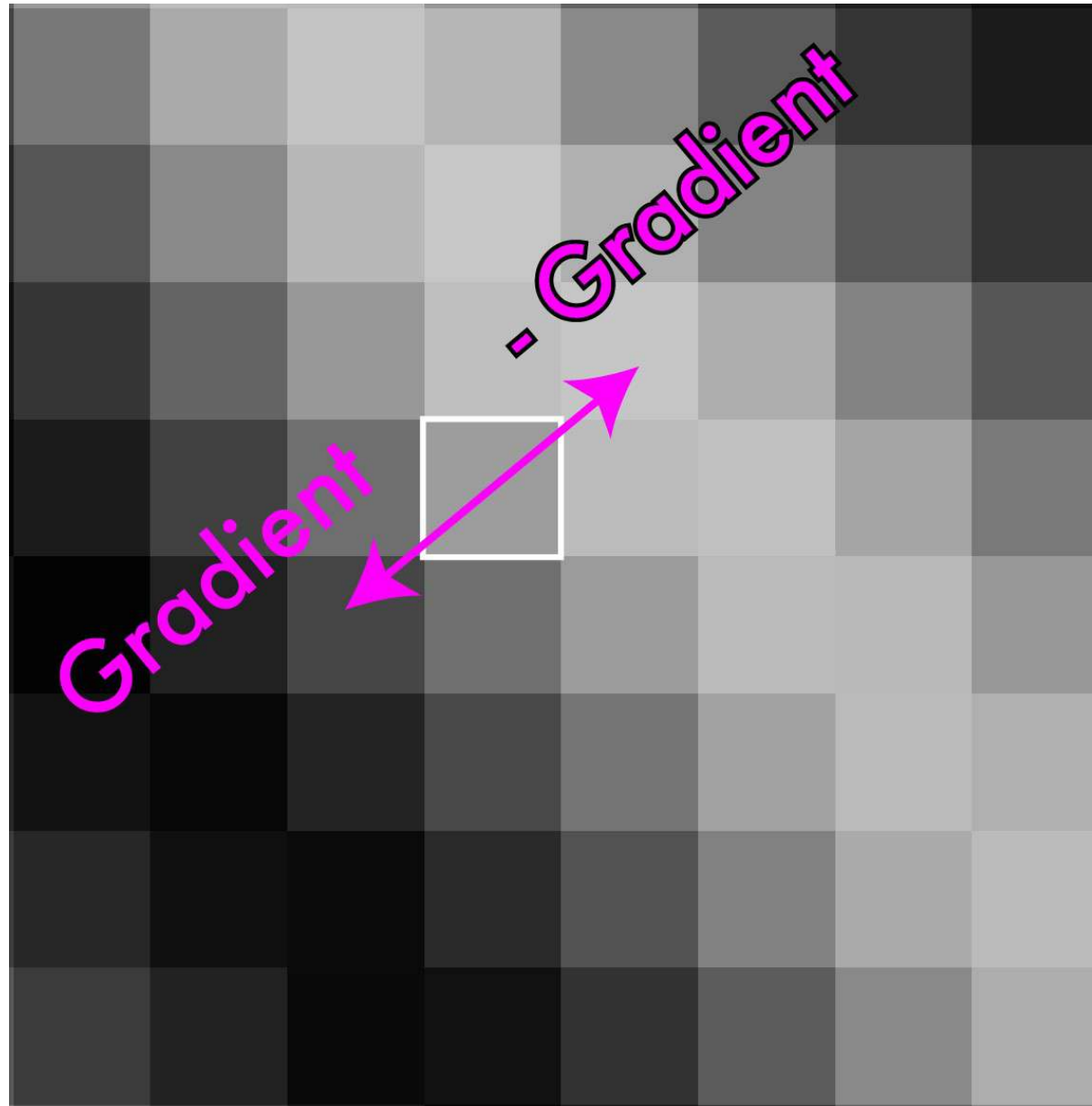


<http://bigwww.epfl.ch/demo/ip/demos/edgeDetector/>

## Non-maximum suppression

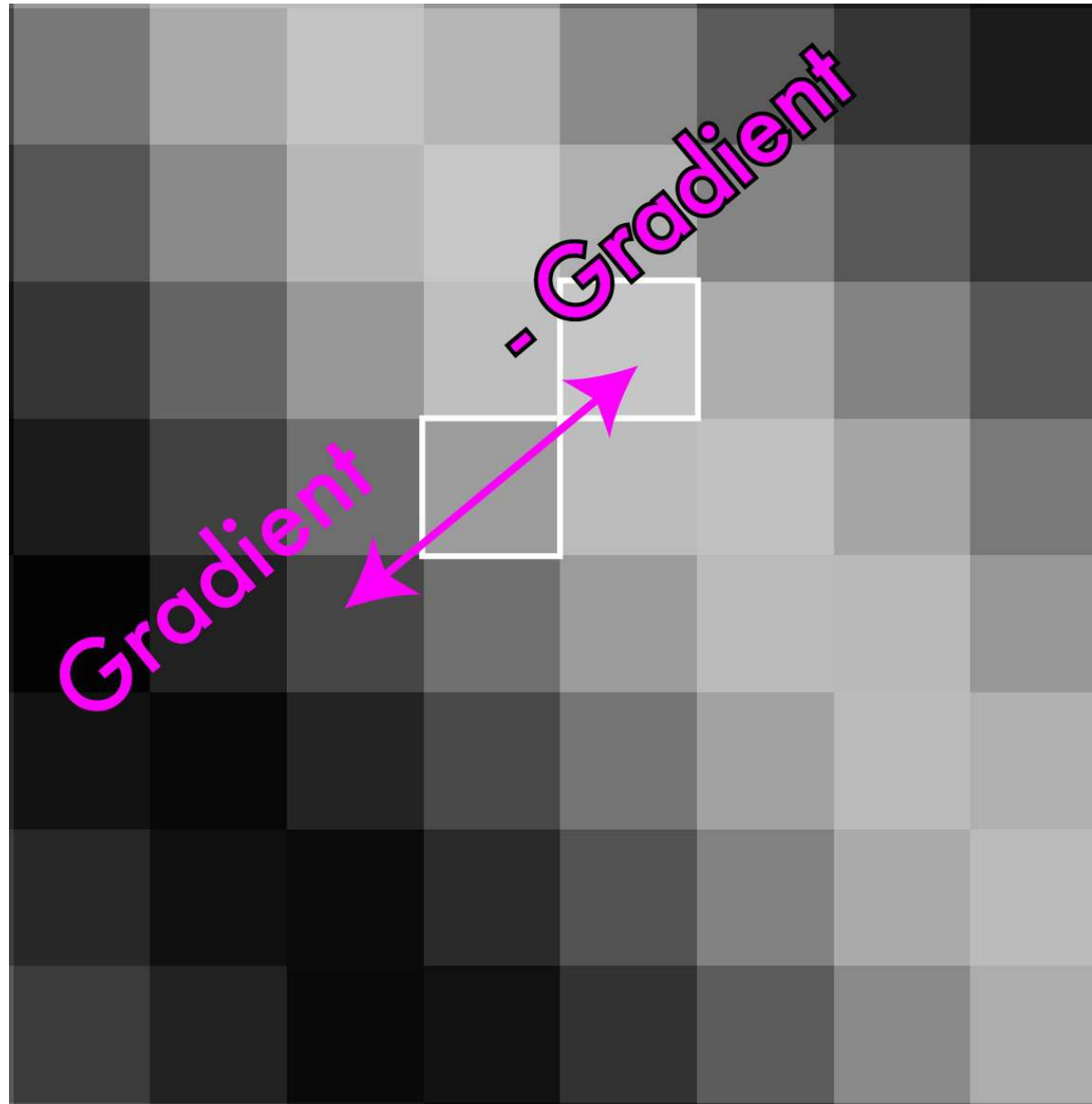


## Non-maximum suppression

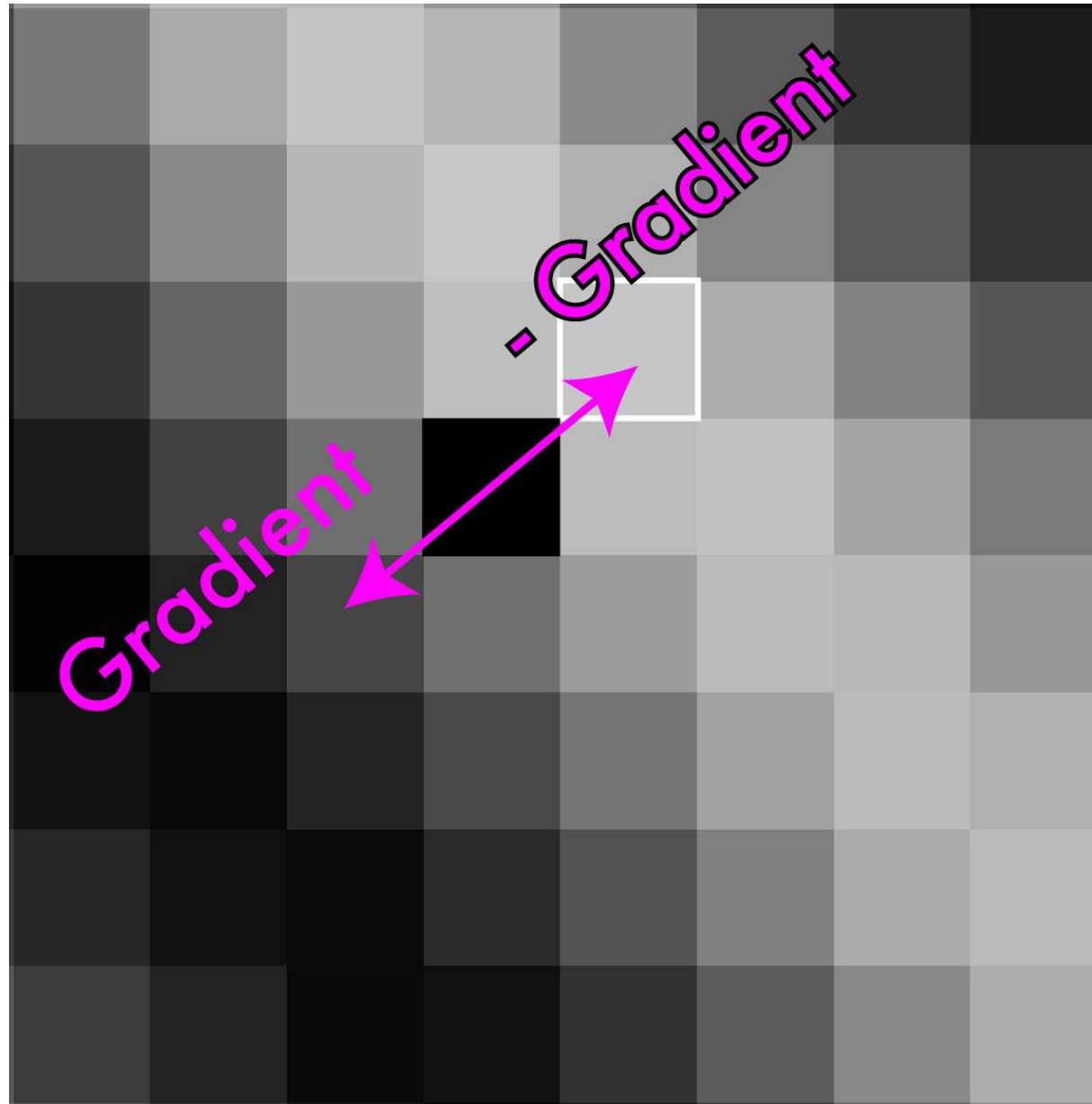




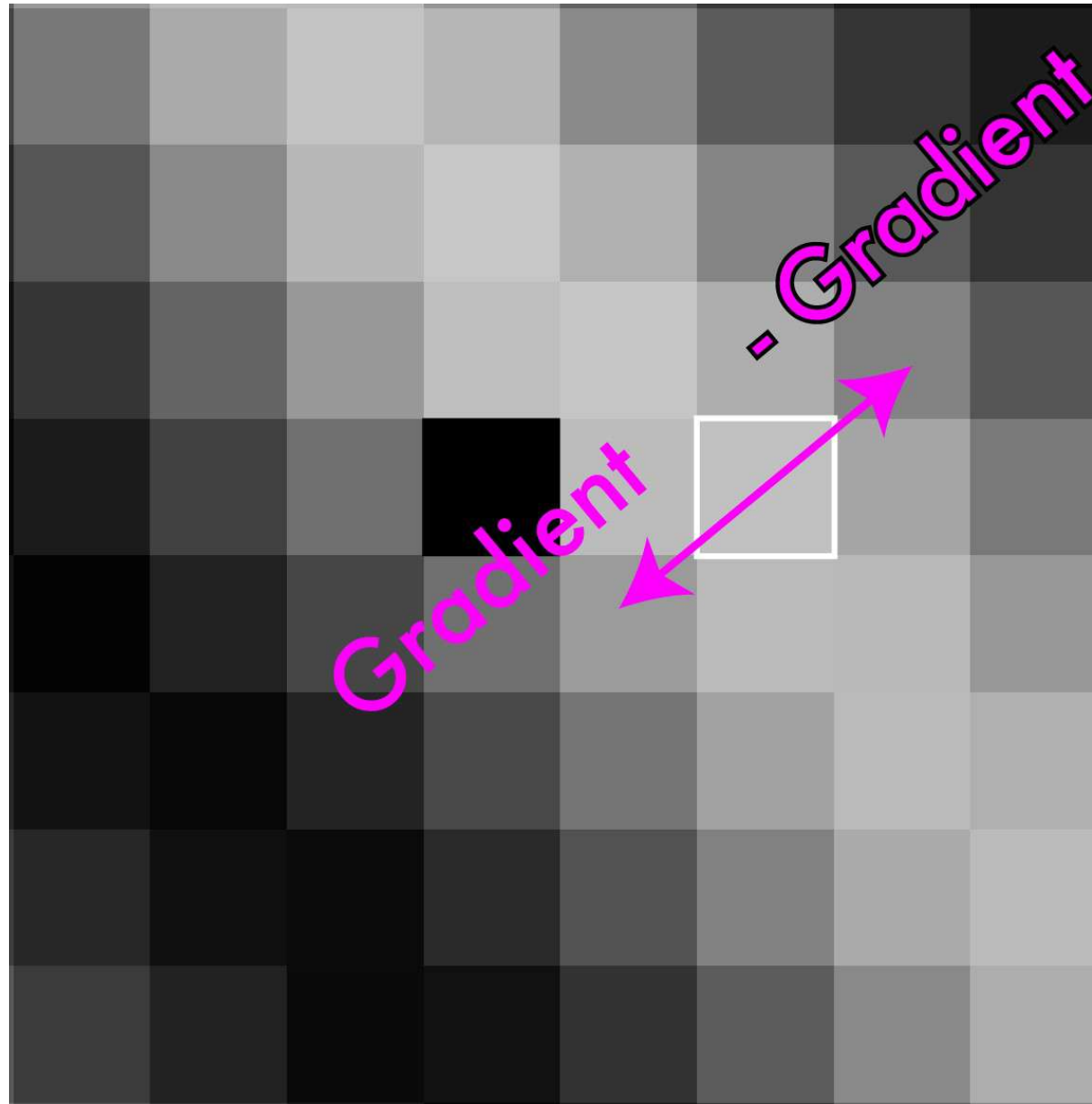
## Non-maximum suppression



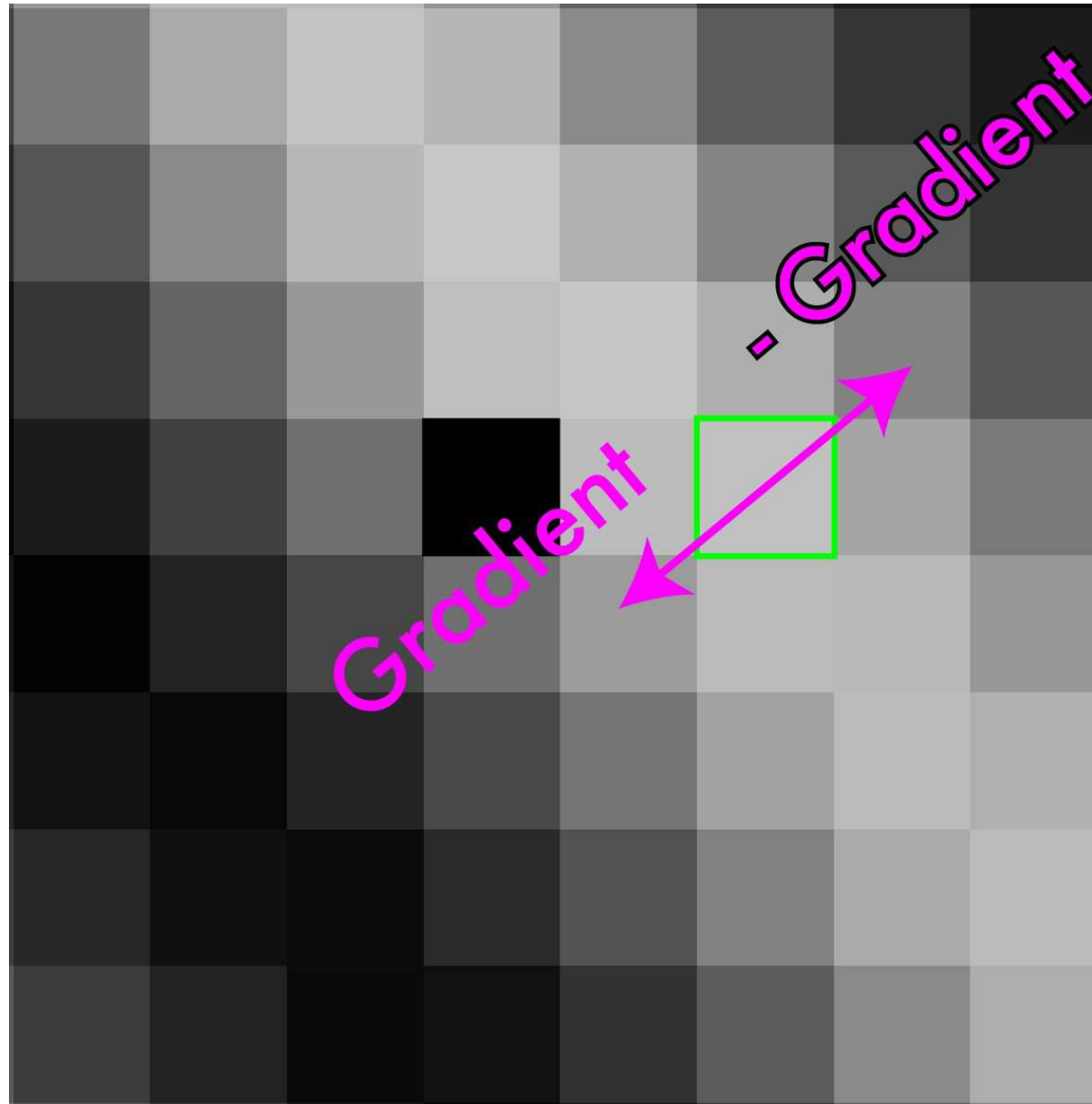
## Non-maximum suppression



## Non-maximum suppression

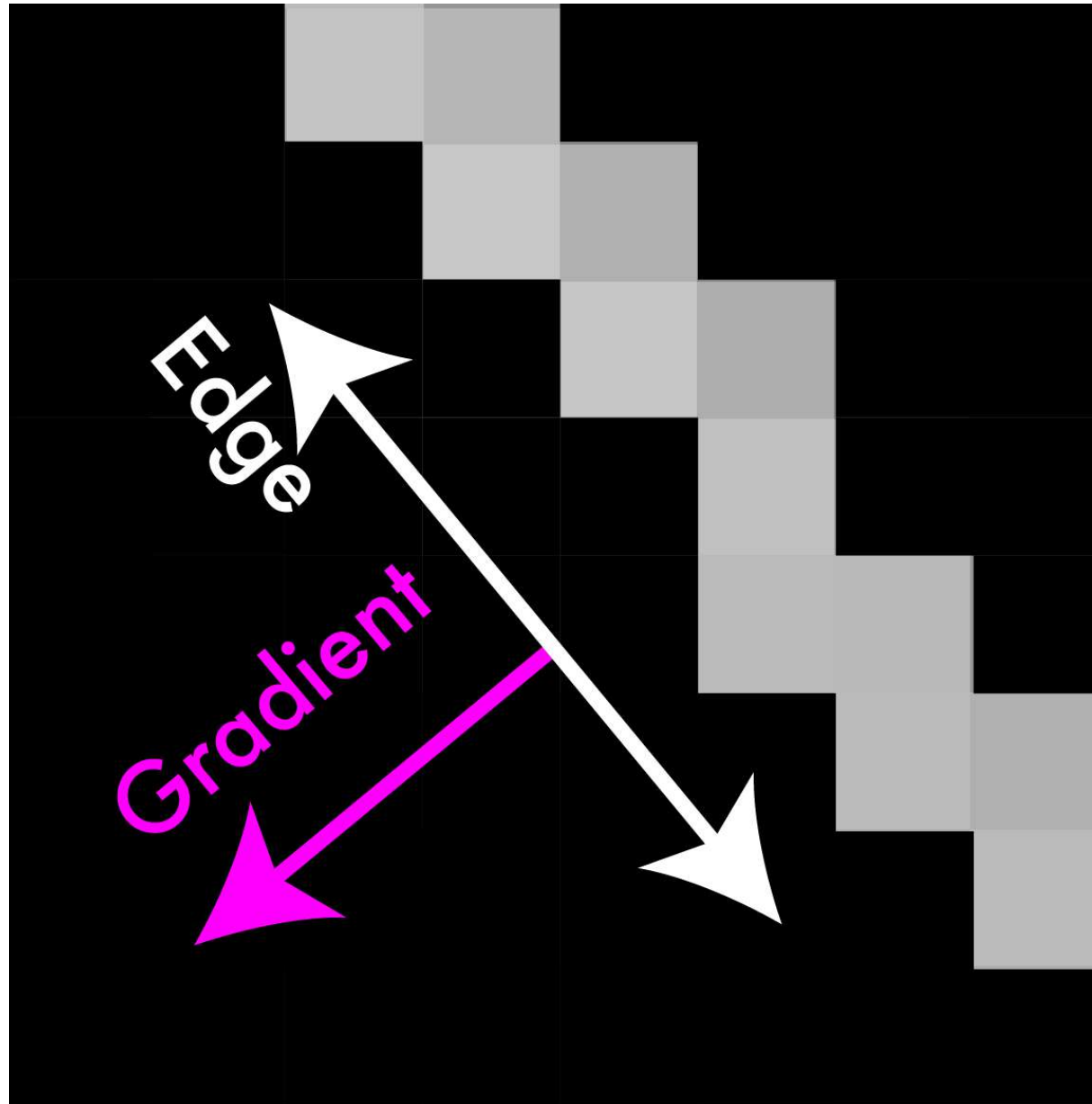


# Non-maximum suppression





## Non-maximum suppression



## Non-maximum suppression



<http://bigwww.epfl.ch/demo/ip/demos/edgeDetector/>

# Threshold edges

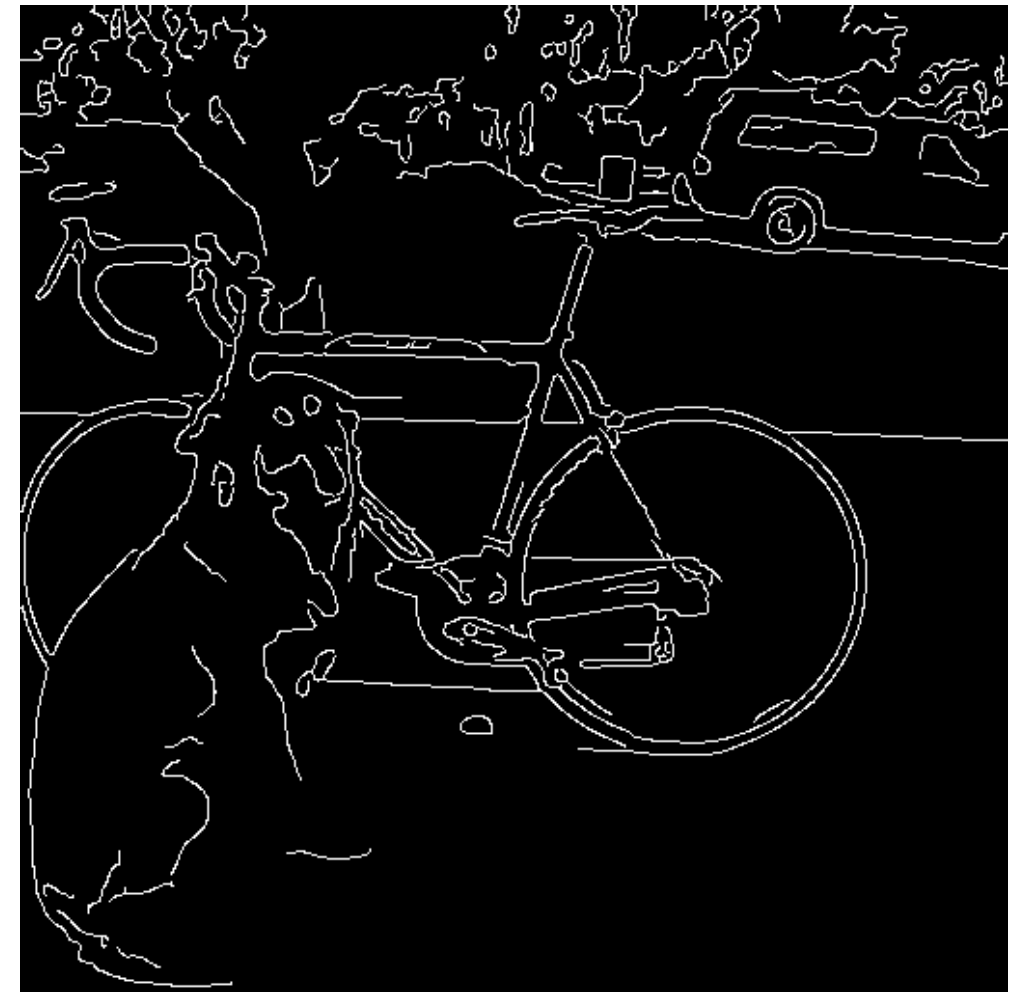
- Still some noise
- Only want strong edges
- 2 thresholds, 3 cases
  - $R > T$ : strong edge
  - $R < T$  but  $R > t$ : weak edge
  - $R < t$ : no edge
- Why two thresholds?



<http://bigwww.epfl.ch/demo/ip/demos/edgeDetector/>

# Connect them up!

- Strong edges are edges!
- Weak edges are edges iff they connect to strong
- Look in some neighborhood (usually 8 closest)



<http://bigwww.epfl.ch/demo/ip/demos/edgeDetector/>



# Canny Edge Detection

## Algorithm:

- Smooth image (only want “real” edges, not noise)
- Calculate gradient direction and magnitude
- Non-maximum suppression perpendicular to edge
- Threshold into strong, weak, no edge
- Connect together components
- Tunable: Sigma, thresholds

<http://bigwww.epfl.ch/demo/ip/demos/edgeDetector/>

## Canny Edge Detection



<http://bigwww.epfl.ch/demo/ip/demos/edgeDetector/>

**MAI4CAREU**

Master programmes in Artificial  
Intelligence 4 Careers in Europe



# Thank you.





University of Cyprus – MSc Artificial Intelligence

# MAI644 – COMPUTER VISION

## Lecture 7: Features – Corners

**Melinos Averkiou**

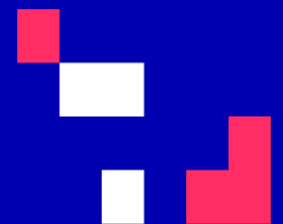
CYENS Centre of Excellence

University of Cyprus - Department of Computer Science

[m.averkiou@cyens.org.cy](mailto:m.averkiou@cyens.org.cy)



**CYENS**  
CENTRE OF EXCELLENCE





# Last time

- What can we do with convolutions
- What is an edge – image derivatives
- Sobel filters
- Laplacian filters
- Difference of Gaussian filters
- Canny edge detection

# Today's Agenda

- Features
- Self-difference
- Harris corner detection

**[material based on Joseph Redmon's course]**



# Today's Agenda

- Features
- Self-difference
- Harris corner detection



# What else is there?



So what can we do with these convolutions anyway?

- Blurring
- Sharpening
- Edges
- Features
- Derivatives
- Super-resolution
- Classification
- Detection
- Image captioning
- ...

# Features!

- Highly descriptive local regions
- Ways to describe those regions
- Useful for:
  - Matching
  - Recognition
  - Detection

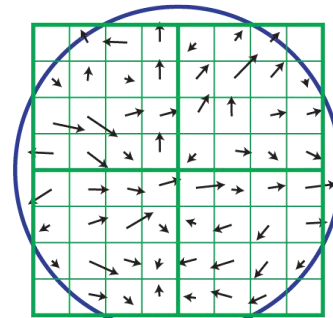
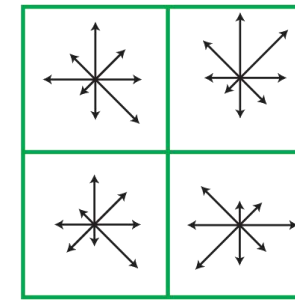
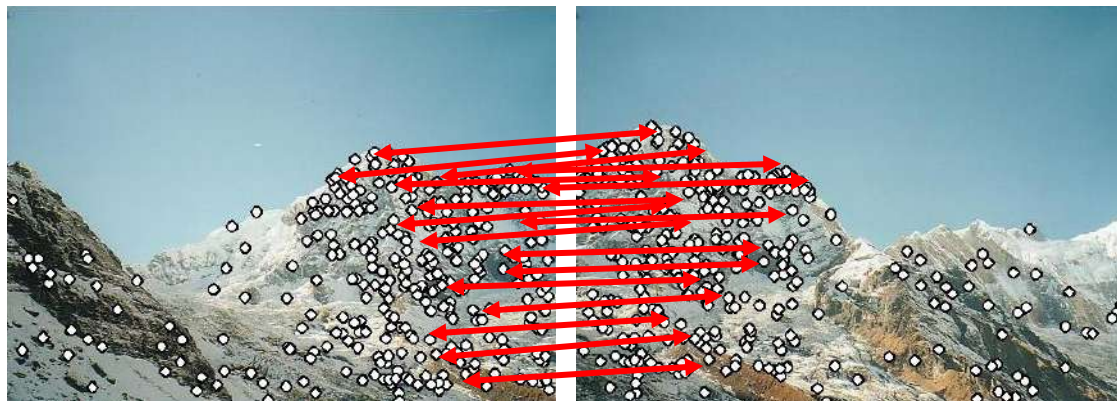
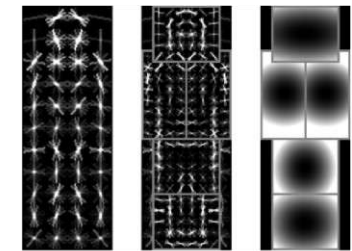
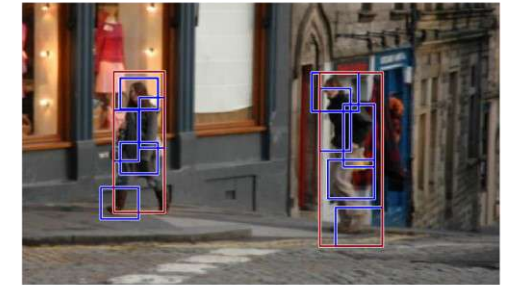


Image gradients



Keypoint descriptor



# What makes a good feature?

- Want to find patches in image that are useful or have some meaning
- For objects, want a patch that is common to that object but not in general
- For panorama stitching, want patches that we can find easily in another image of same place
- Good features are unique!
  - Can find the “same” feature easily
  - Not mistaken for “different” features

# Application - How to create a panorama

- Say we are stitching images to create a panorama
- Want patches in one image to match to patches in another image
- Hopefully only match one spot



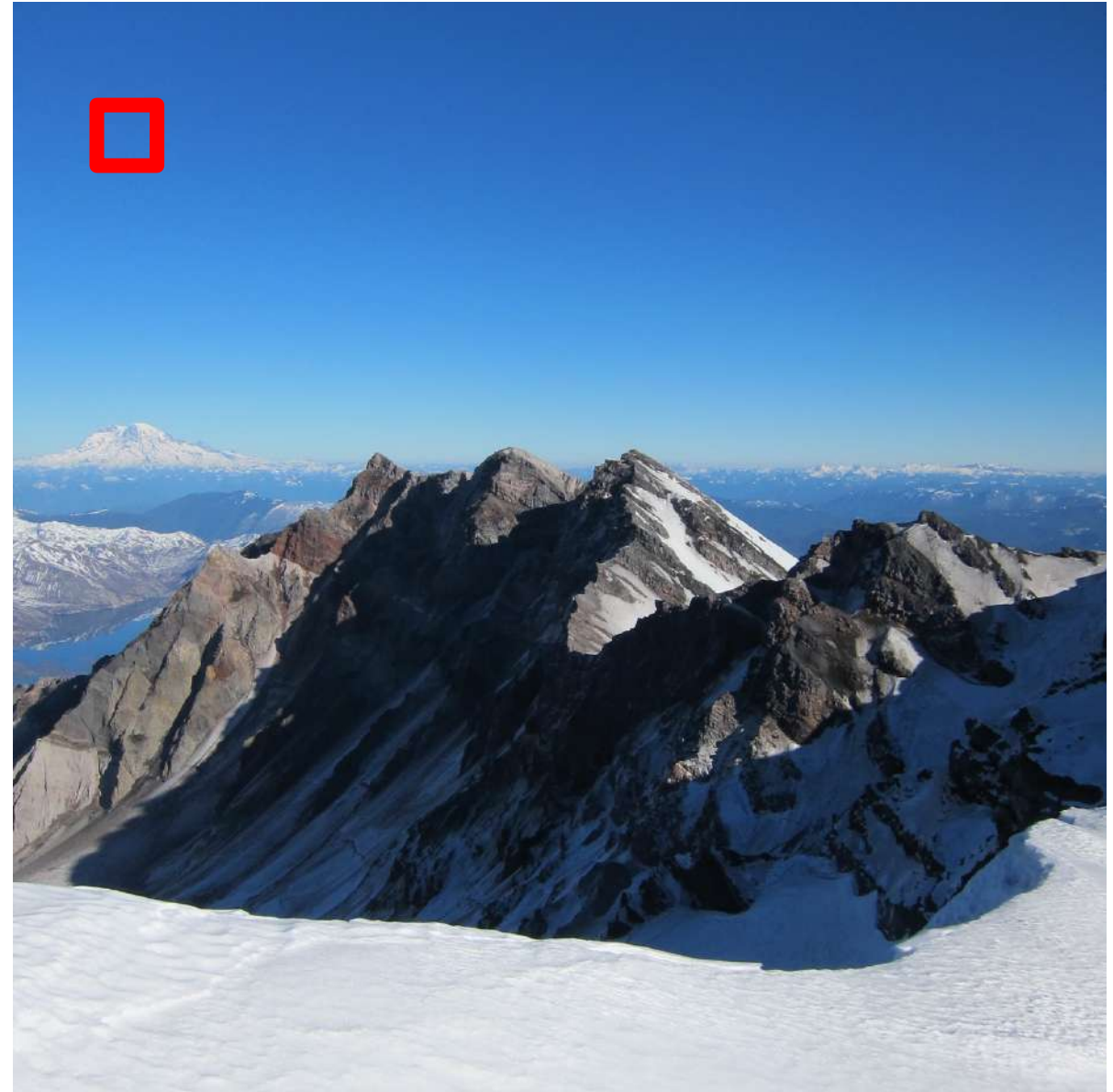


# How close are two patches?

- (weighted) Summed squared difference
- Images I, J
- $\sum_{x,y} w(x,y) (I(x,y) - J(x,y))^2$

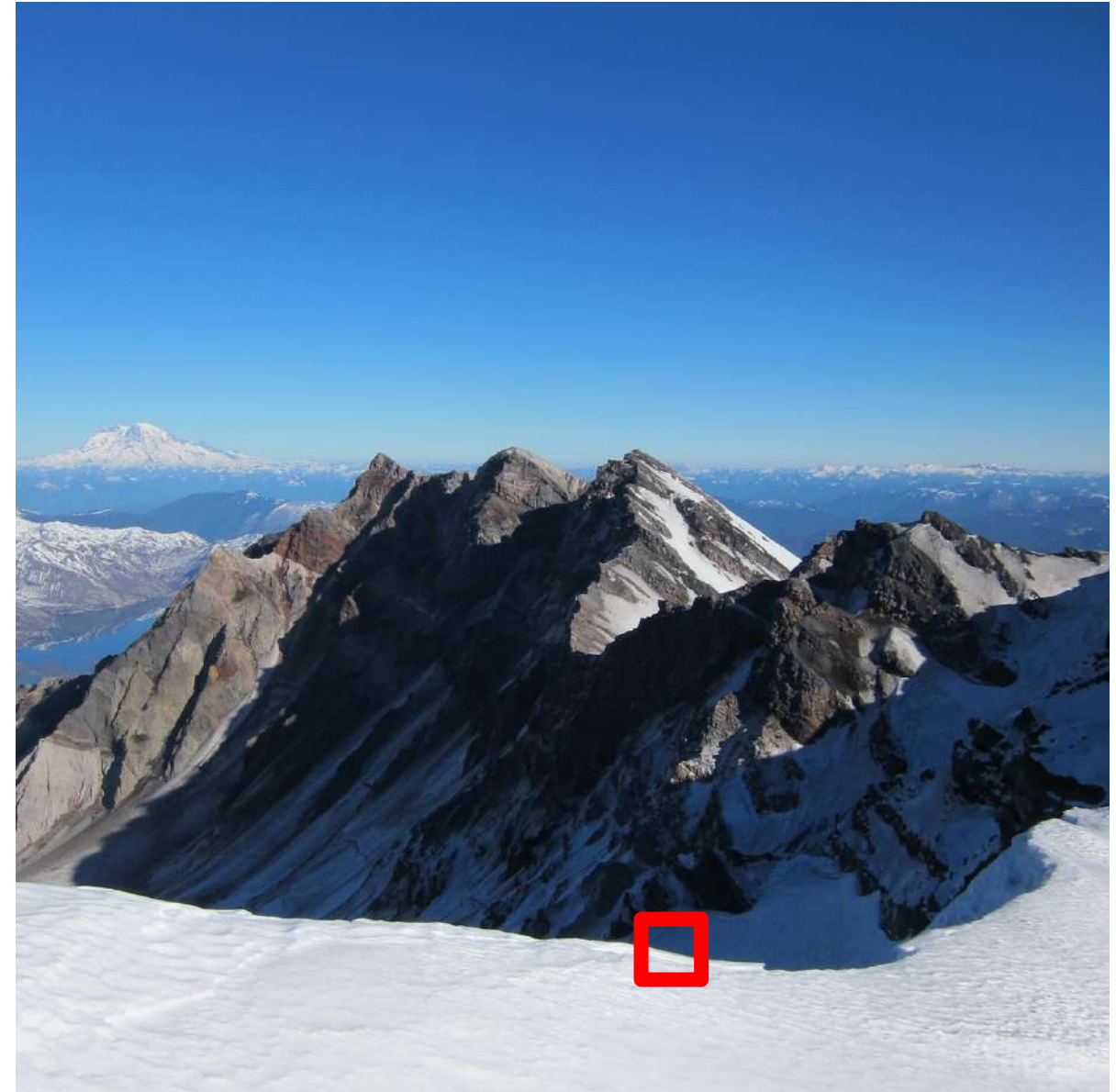
# How can we find unique patches?

- Sky: bad
  - Very little variation
  - Could match any other sky



# How can we find unique patches?

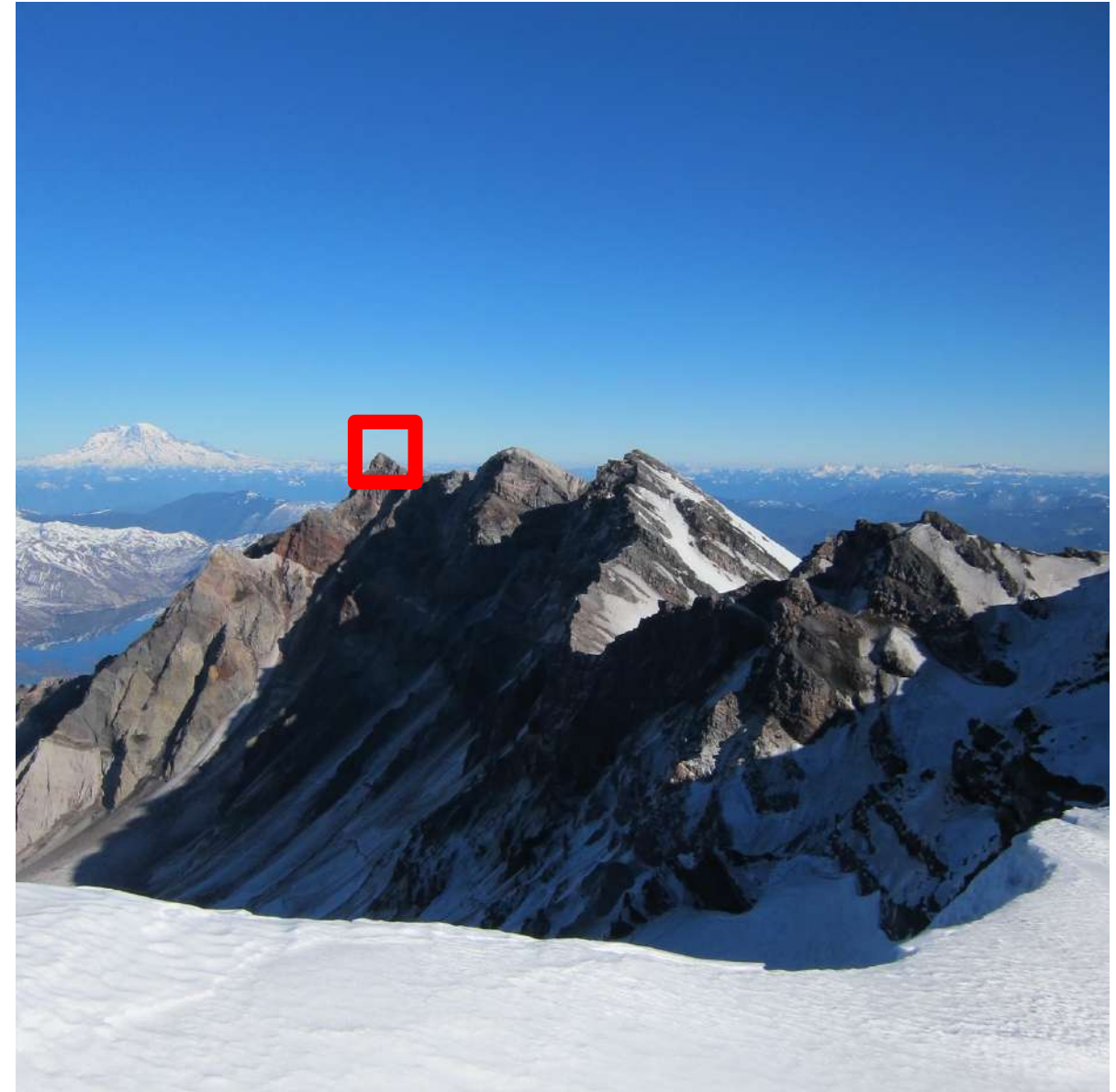
- Sky: bad
  - Very little variation
  - Could match any other sky
- Edge: ok
  - Variation in one direction
  - Could match other patches along same edge





# How can we find unique patches?

- Sky: bad
  - Very little variation
  - Could match any other sky
- Edge: ok
  - Variation in one direction
  - Could match other patches along same edge
- Corners: good!
  - Only one alignment matches





# Today's Agenda

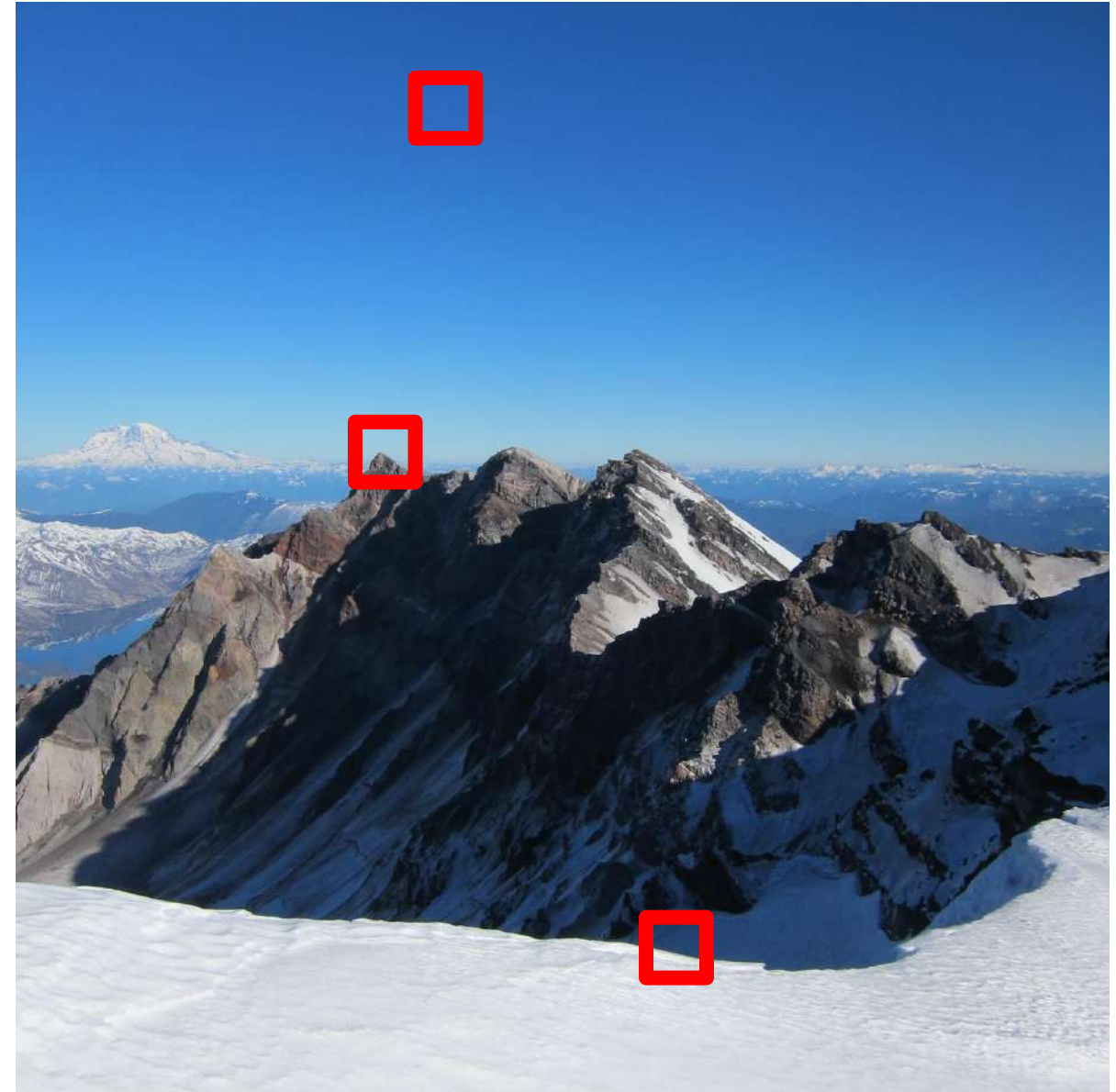
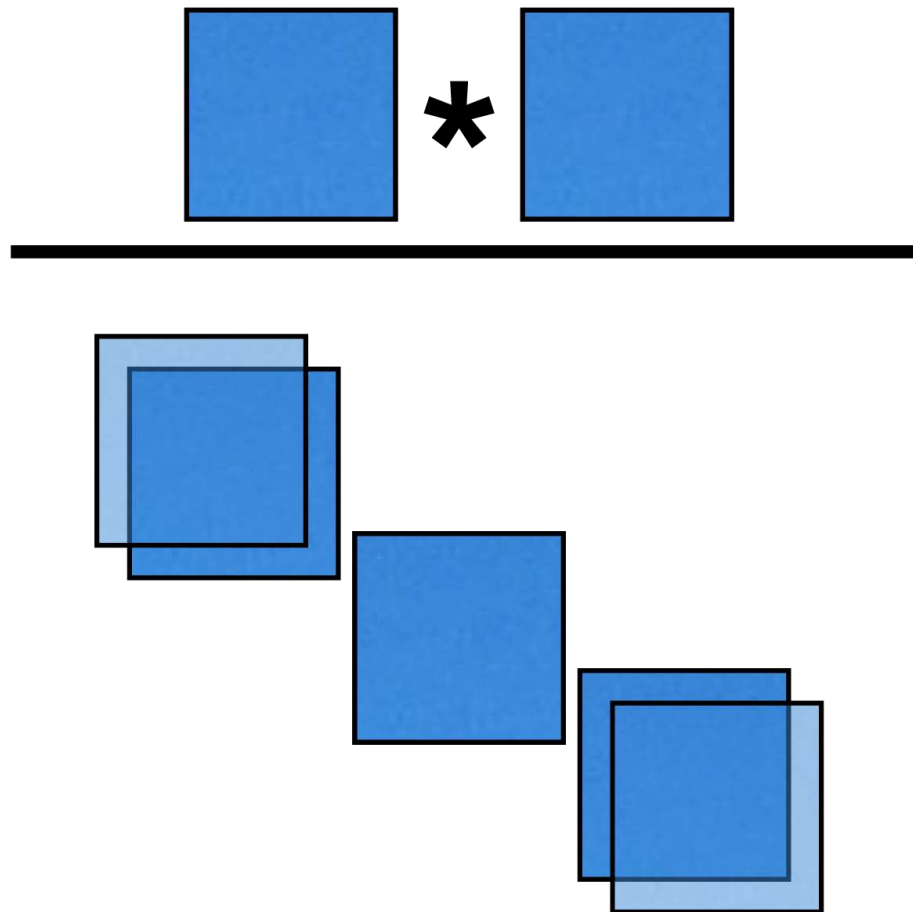
- Features
- Self-difference
- Harris corner detection

# How can we find unique patches?

- Want a patch that is unique in the image
- Can calculate distance between patch and every other patch, a lot of computation
- Instead, we could think about auto-correlation:
  - How well does image match shifted version of itself?
- Autocorrelation:  $\sum_{\mathbf{d}} \sum_{x,y} w(x,y) ( I(x+\mathbf{d}_x, y+\mathbf{d}_y) - I(x,y) )^2$
- Measure of self-difference (how much am I not myself?)

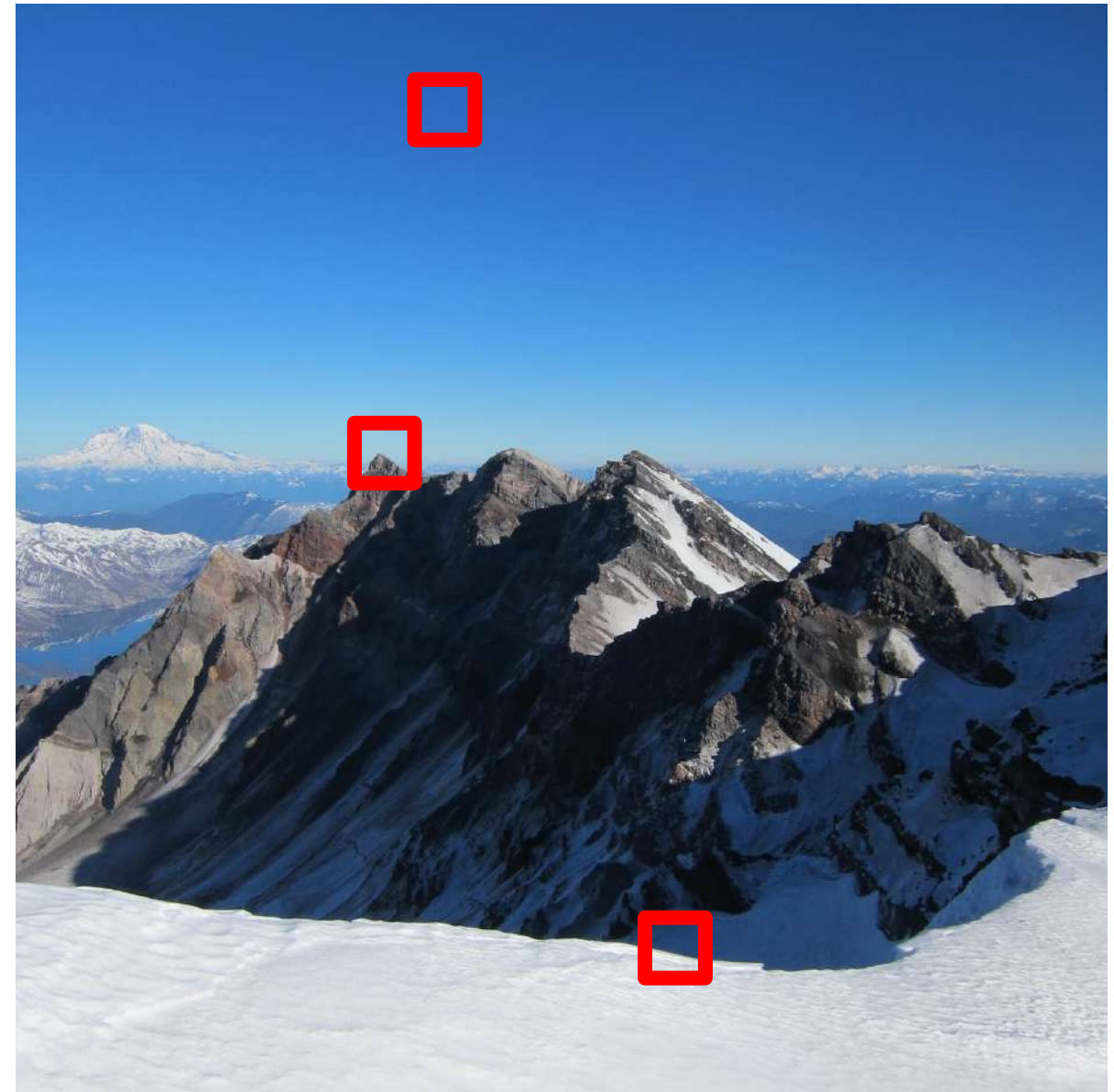
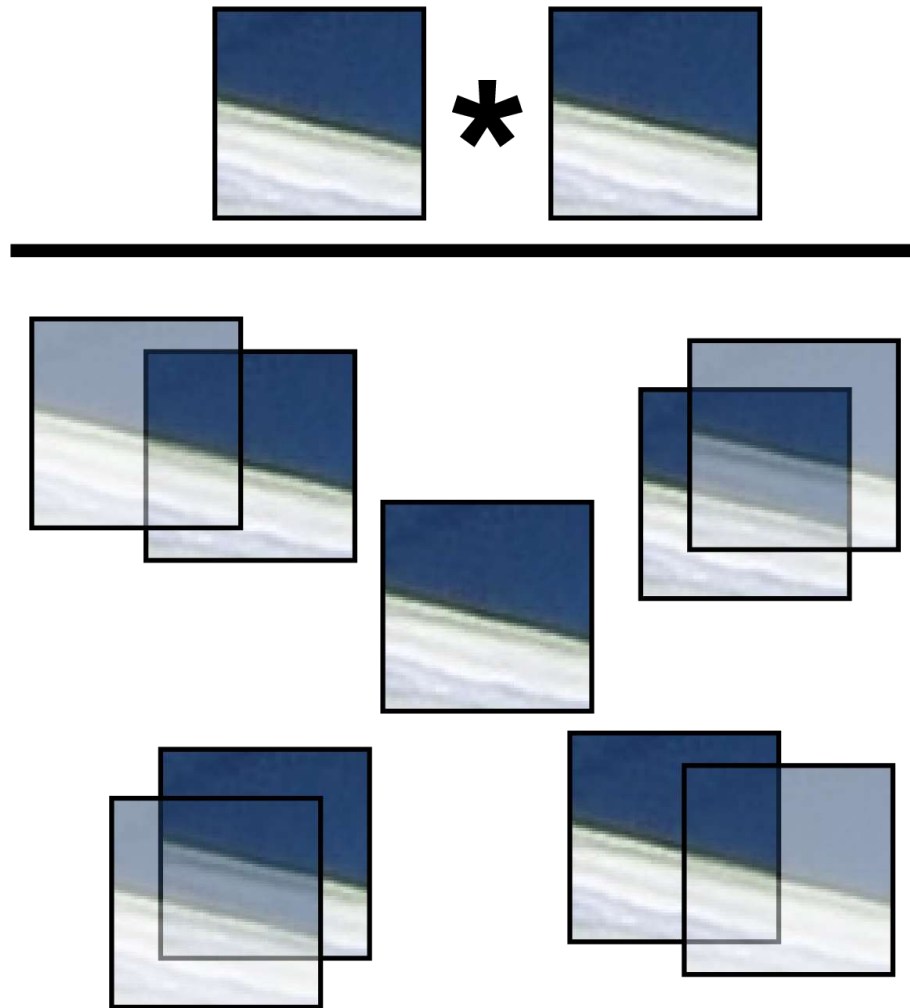
# Self-difference

Sky: low everywhere



## Self-difference

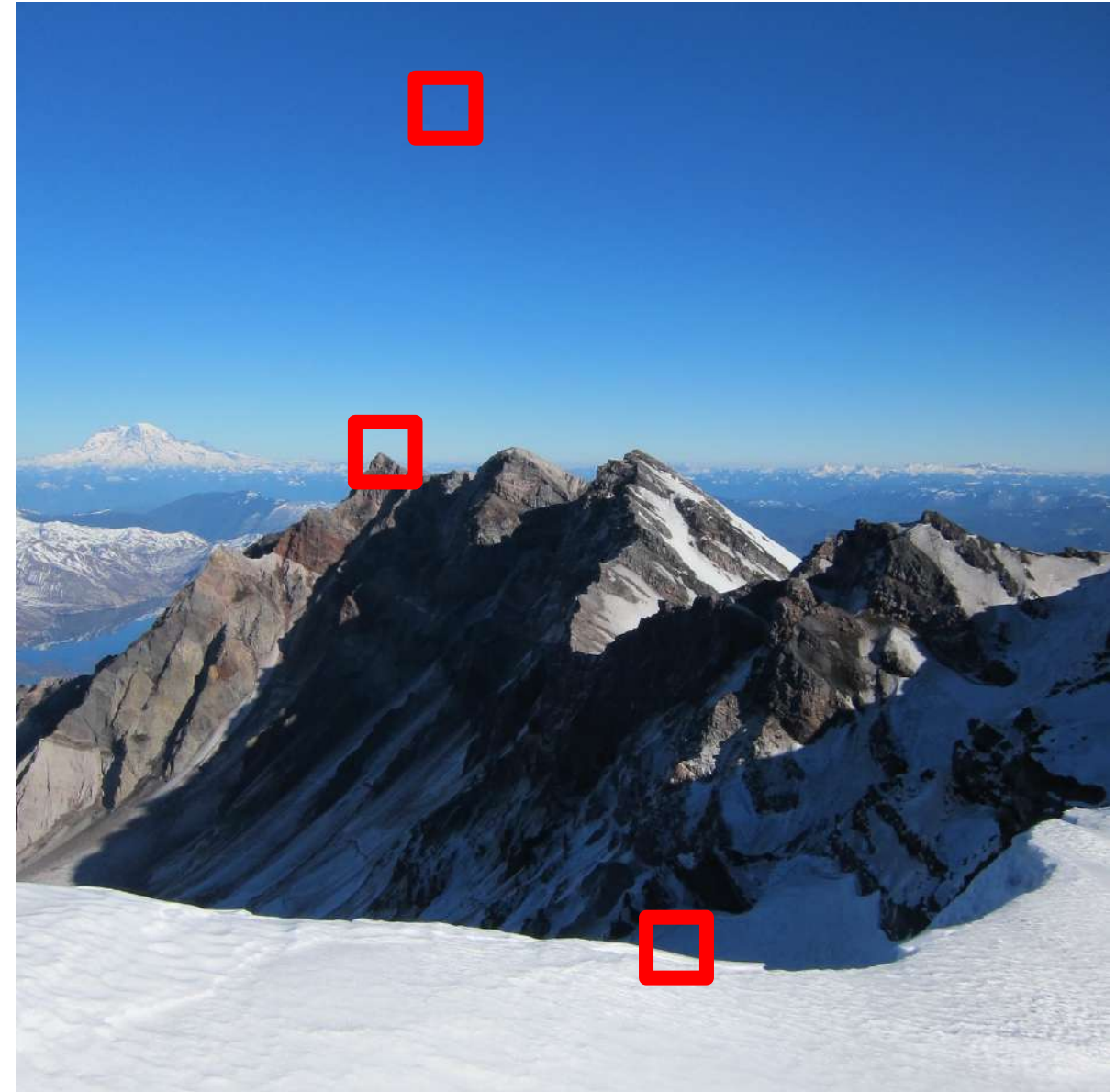
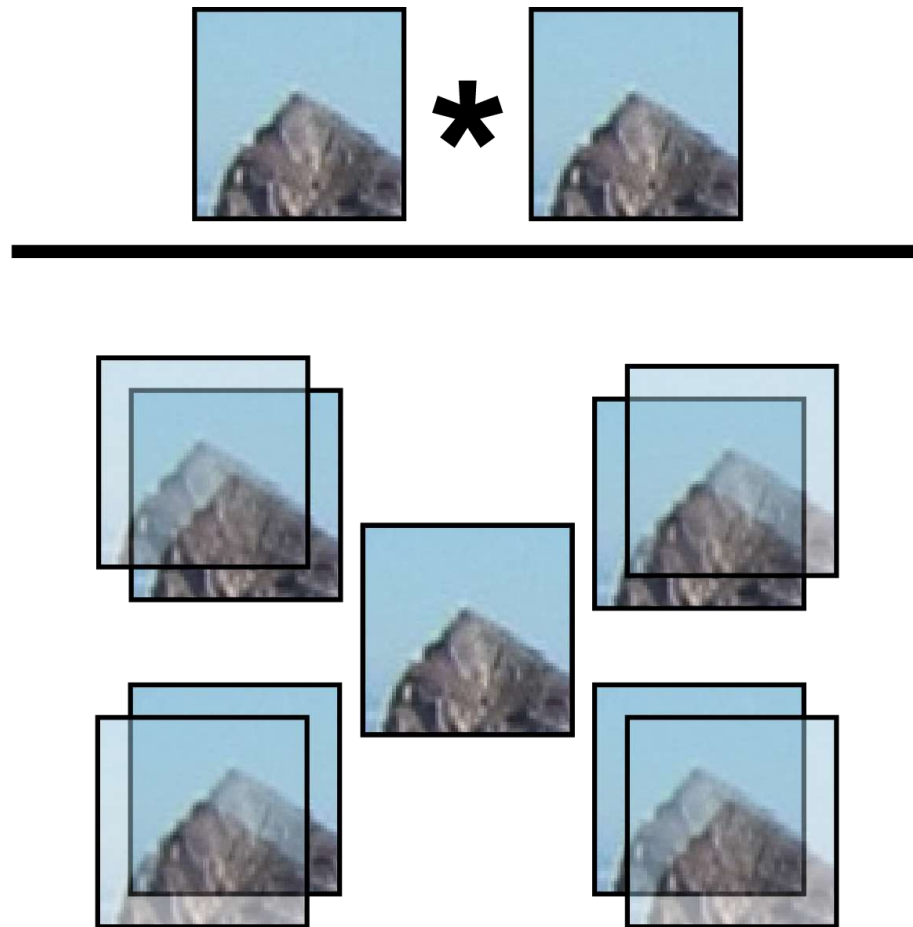
Edge: low along edge





## Self-difference

Corner: mostly high

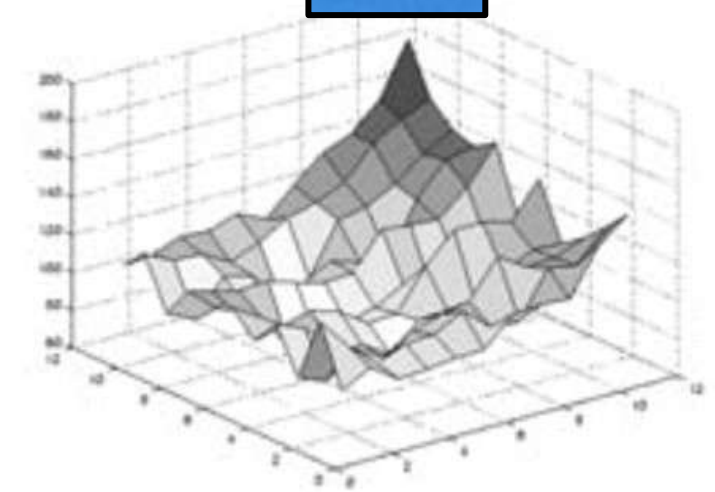
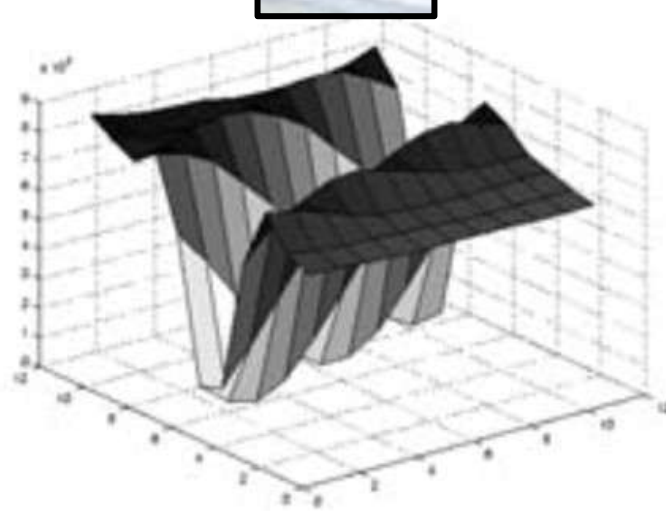
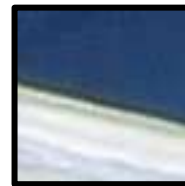
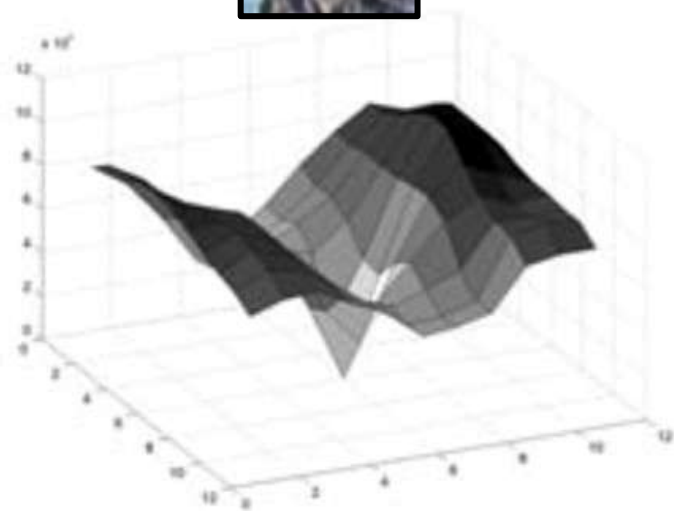


## Self-difference

Sky: low everywhere

Edge: low along edge

Corner: mostly high



# Today's Agenda

- Features
- Self-difference
- Harris corner detection

# Self-difference is still expensive

- $\sum_{\mathbf{d}} \sum_{x,y} w(x,y) ( I(x+\mathbf{d}_x, y+\mathbf{d}_y) - I(x,y) )^2$
- Lots of summing
- Need an approximation
  - Use Taylor Series expansion of the image function (see Szeliski book Section 7.1.1)
  - $I(x+\mathbf{d}_x, y+\mathbf{d}_y) \approx I(x,y) + \mathbf{d}_x I_x(x,y) + \mathbf{d}_y I_y(x,y)$



# Approximate self-difference

- Look at nearby gradients  $I_x$  and  $I_y$
- If gradients are mostly zero, not a lot going on
  - Low self-difference
- If gradients are mostly in one direction, edge
  - Still low self-difference
- If gradients are in two(ish) directions, corner!
  - High self-difference, good patch!

# Approximate self-difference

- How do we tell what's going on with gradients?
- Eigenvectors/values!
- Need structure matrix for patch, just a weighted sum of nearby gradient information

$$\begin{array}{|c|c|} \hline \sum_i w_i l_x(i) l_x(i) & \sum_i w_i l_x(i) l_y(i) \\ \hline \sum_i w_i l_x(i) l_y(i) & \sum_i w_i l_y(i) l_y(i) \\ \hline \end{array}$$

- Not as complex as it looks, weighted sum of gradients near pixel

[https://en.wikipedia.org/wiki/Structure\\_tensor](https://en.wikipedia.org/wiki/Structure_tensor)

# Structure matrix

- Weighted sum of gradient information

$$\begin{vmatrix} \sum_i w_i l_x(i) l_x(i) & \sum_i w_i l_x(i) l_y(i) \\ \sum_i w_i l_x(i) l_y(i) & \sum_i w_i l_y(i) l_y(i) \end{vmatrix}$$

- Can use Gaussian weighting
- Eigenvectors/values of this matrix summarize the distribution of the gradients nearby
- $\lambda_1$  and  $\lambda_2$  are eigenvalues
  - $\lambda_1$  and  $\lambda_2$  both small: no gradient
  - $\lambda_1 \gg \lambda_2$  or  $\lambda_1 \ll \lambda_2$  : gradient in one direction
  - $\lambda_1$  and  $\lambda_2$  both large: multiple gradient directions, corner

# Estimating eigenvalues

- Use determinant and trace:
  - $\det(S) = \lambda_1 * \lambda_2$
  - $\text{trace}(S) = \lambda_1 + \lambda_2$
- Response function (R score):  $\det(S) - \kappa \text{trace}(S)^2 = \lambda_1 \lambda_2 - \kappa (\lambda_1 + \lambda_2)^2$
- If this score is large, both  $\lambda_1$  and  $\lambda_2$  are large



# Harris Corner Detector

- Calculate derivatives  $I_x$  and  $I_y$
- Calculate 3 measures  $I_x I_x$ ,  $I_y I_y$ ,  $I_x I_y$
- Calculate weighted sums
  - Want a (weighted) sum of nearby pixels, guess what this is?
  - Gaussian – or weights could just be ones
- Estimate response based on R score
- Non-max suppression in a neighborhood

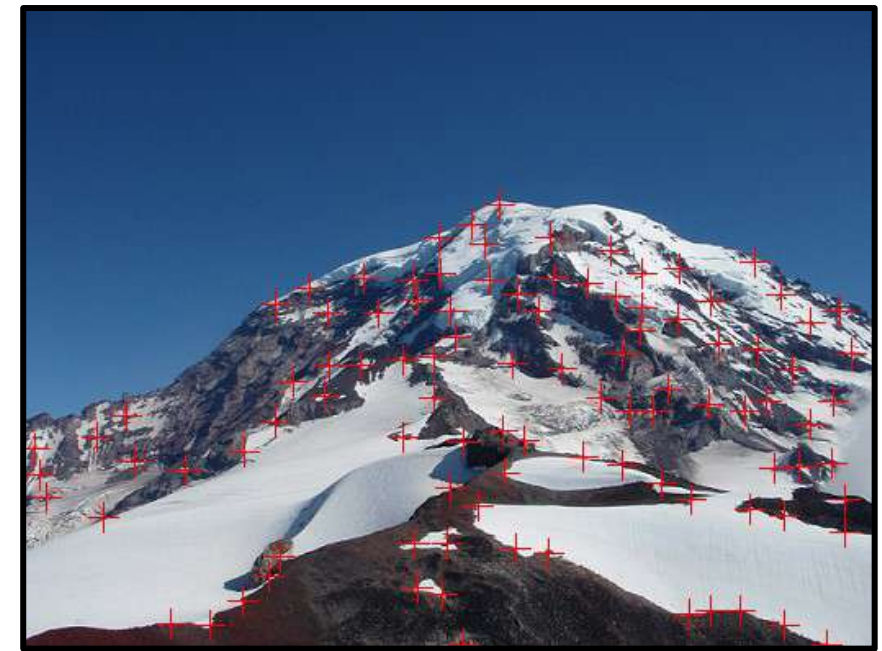






# Ok, we found corners, now what?

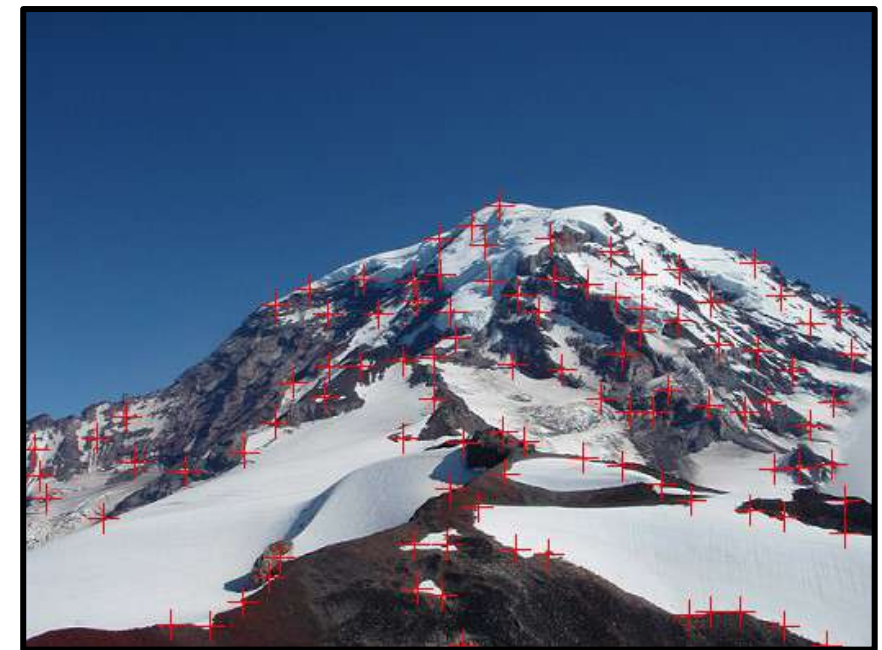
- Need to match image patches to each other
- Need to figure out transform between images





# Ok, we found corners, now what?

- Need to match image patches to each other
  - What is a patch? How do we look for matches? Pixels?
- Need to figure out transform between images
  - How can we transform images?
  - How do we solve for this transformation given matches?



**MAI4CAREU**

Master programmes in Artificial  
Intelligence 4 Careers in Europe



# Thank you.





University of Cyprus – MSc Artificial Intelligence

# MAI644 – COMPUTER VISION

## Lecture 8: Feature Descriptors and Image Transforms

**Melinos Averkiou**

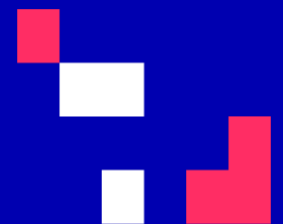
CYENS Centre of Excellence

University of Cyprus - Department of Computer Science

[m.averkiou@cyens.org.cy](mailto:m.averkiou@cyens.org.cy)



**CYENS**  
CENTRE OF EXCELLENCE



## Last time

- Features
- Self-difference
- Harris corner detection



# Today's Agenda

- Basic feature descriptor and matching
- Histogram of Oriented Gradients
- SIFT
- Image transformations
- Estimate transformations

**[material based on Joseph Redmon's course]**

# Today's Agenda

- Basic feature descriptor and matching
- Histogram of Oriented Gradients
- SIFT
- Image transformations
- Estimate transformations

# Ok, we found corners, now what?

- Need to match image patches to each other
- Need to figure out transform between images



# Ok, we found corners, now what?

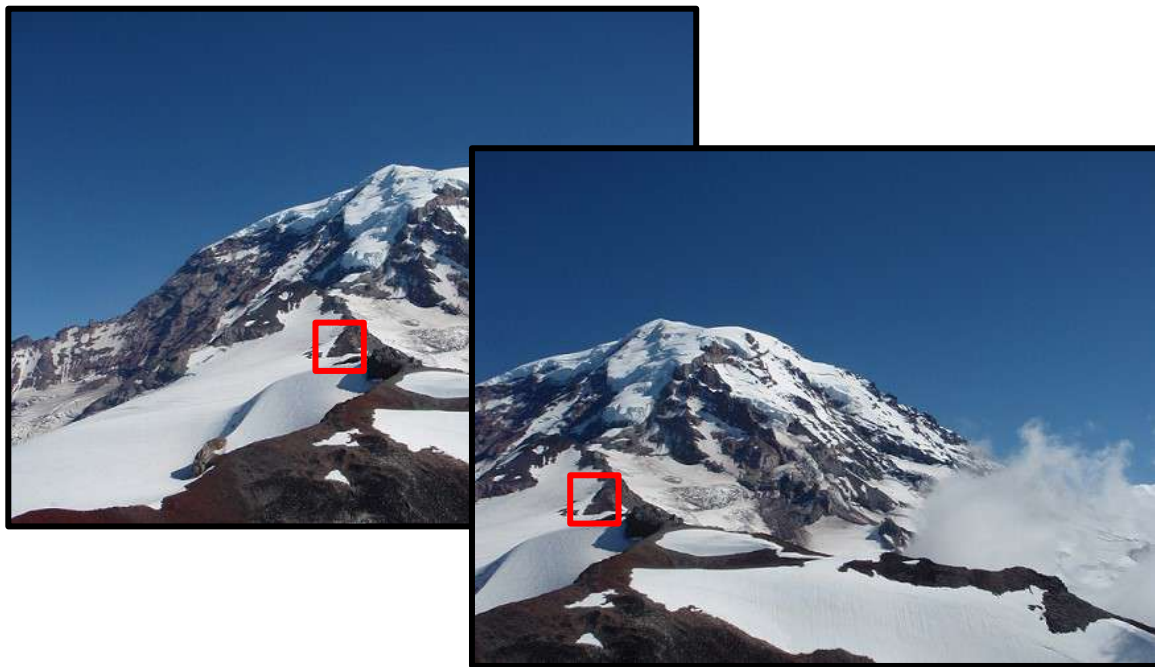
- Need to match image patches to each other
  - What is a match? How do we look for matches? Pixels?
- Need to figure out transform between images





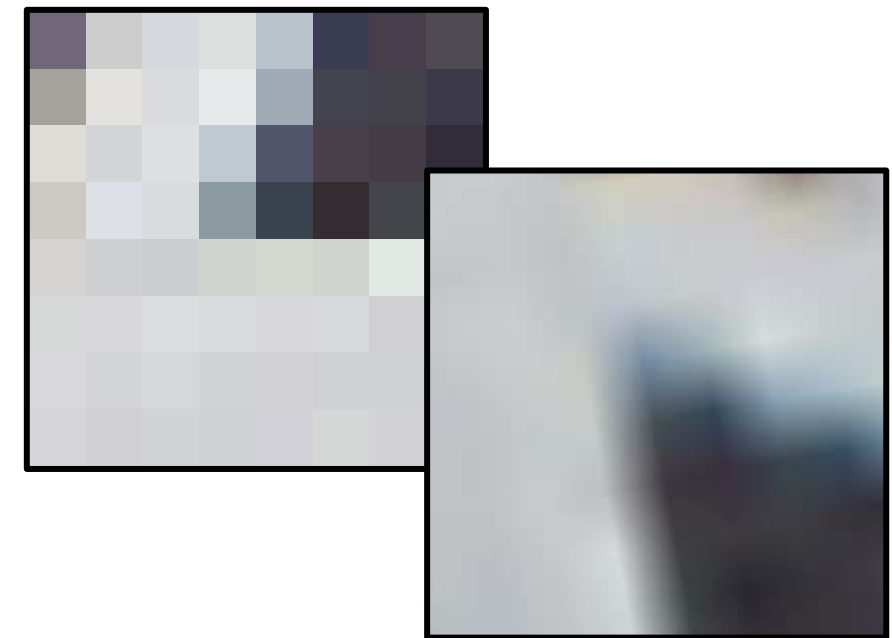
# Matching patches: descriptors!

- We want a way to represent an image patch
- Can be very simple, just pixels!
- Finding matching patch is easy, distance metric:
  - $\sum_{x,y} (I(x,y) - J(x,y))^2$
  - What problems are there with just using pixels?



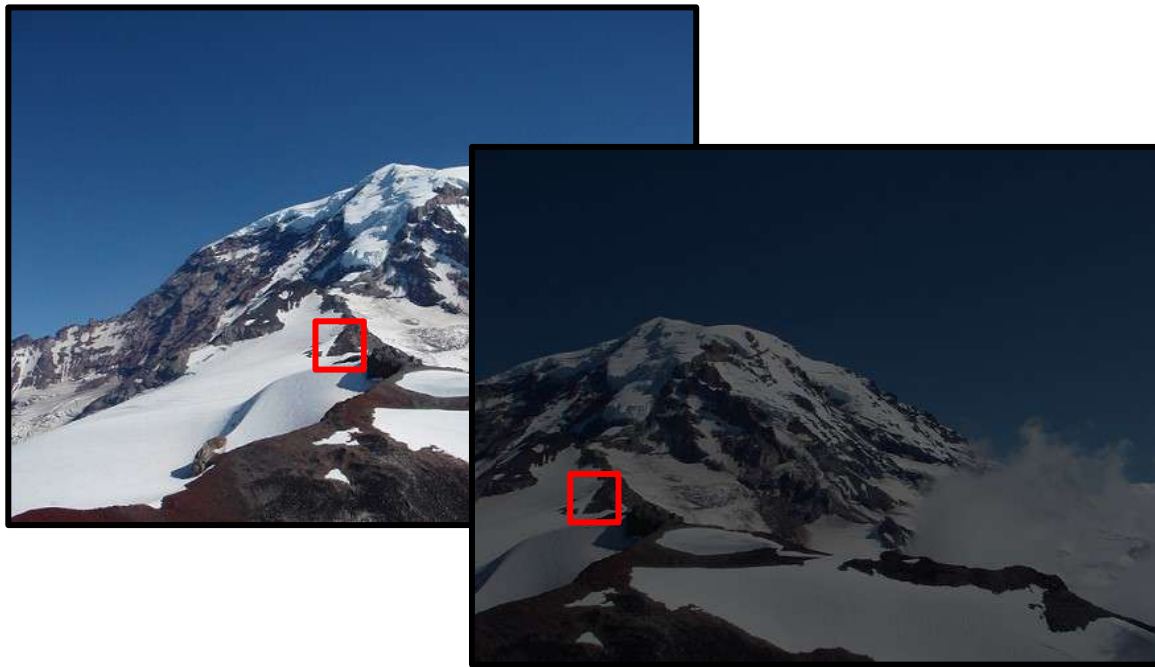
# Matching patches: descriptors!

- We want a way to represent an image patch
- Can be very simple, just pixels!
- Finding matching patch is easy, distance metric:
  - $\sum_{x,y} (I(x,y) - J(x,y))^2$
  - Not invariant to some image transformations (e.g. rotation, scaling) !



# Matching patches: descriptors!

- We want a way to represent an image patch
- Can be very simple, just pixels!
- Finding matching patch is easy, distance metric:
  - $\sum_{x,y} (I(x,y) - J(x,y))^2$
  - Not invariant to lighting changes !



# Matching patches: descriptors!

- We want feature descriptors invariant to lighting and image transforms !
- Descriptors can be more complex
  - Gradient information
  - How much context?



# Today's Agenda

- Basic feature descriptor and matching
- Histogram of Oriented Gradients
- SIFT
- Image transformations
- Estimate transformations

# Histogram of Oriented Gradients (HOG)

- By Dalal and Triggs 2005
- Better image descriptor
- Not reliant on magnitude, just direction
  - Invariant to some lighting changes
- They used it to train an SVM to recognize people

Binning from <http://aishack.in/tutorials/sift-scale-invariant-feature-transform-features/>

# Histogram of Oriented Gradients (HOG)

Steps to calculate HOG Feature Descriptor

1. Compute gradients
2. Bin gradient directions to create histogram
3. Normalize histograms of gradients

# Histogram of Oriented Gradients (HOG)

## Steps to calculate HOG Feature Descriptor

### 1. Compute gradients

Gaussian smoothing (experimented with various  $\sigma$ ), followed by a derivative filter

- $\sigma = 0$  , i.e., no smoothing gave best results
- 1D filter  $[-1, 0, 1]$  gave best results



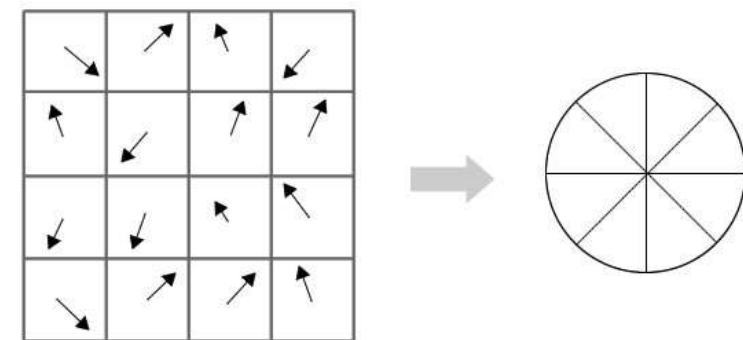
# Histogram of Oriented Gradients (HOG)

Steps to calculate HOG Feature Descriptor

2. Bin gradient directions to create histogram

Split image into 8x8 'cells' and compute histogram for each cell

- Unsigned gradients, i.e.,  $\theta=0-180$  degrees gave best results
- 9 bins gave best results



# Histogram of Oriented Gradients (HOG)

Steps to calculate HOG Feature Descriptor

3. Normalize histograms of gradients

Gather overlapping 'cells' into 'blocks', concatenate histograms and normalize

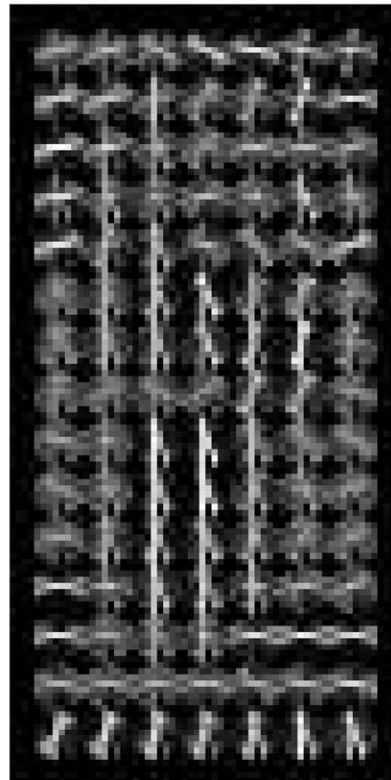
- 16x16 blocks of 4 (2x2) cells gave best results
- L2 Normalization gave best results

# Histogram of Oriented Gradients (HOG)

For each training image of 64x128 there are 7x15 blocks, so the overall descriptor is  $7 \times 15 \times 36 = 3780$  dimensions



Training image



HOG descriptor of the image visualized for each 16x16 block

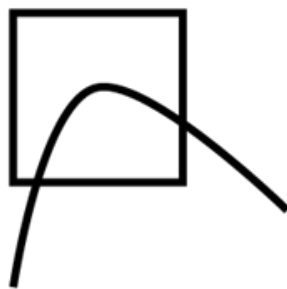


Descriptor weighted by the SVM weights

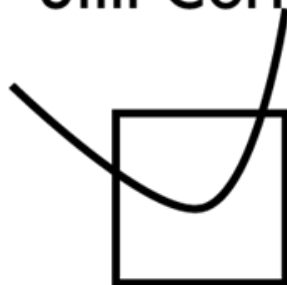
# This is as good as it gets ?

- Not so fast...
- Harris has some issues:
  - Corner detection is rotation invariant
  - Harris not invariant to scale
- Descriptors are also hard
  - Just looking at pixels is not rotation invariant!
  - HOG also not rotation invariant

Corner



Still Corner



Not Corner





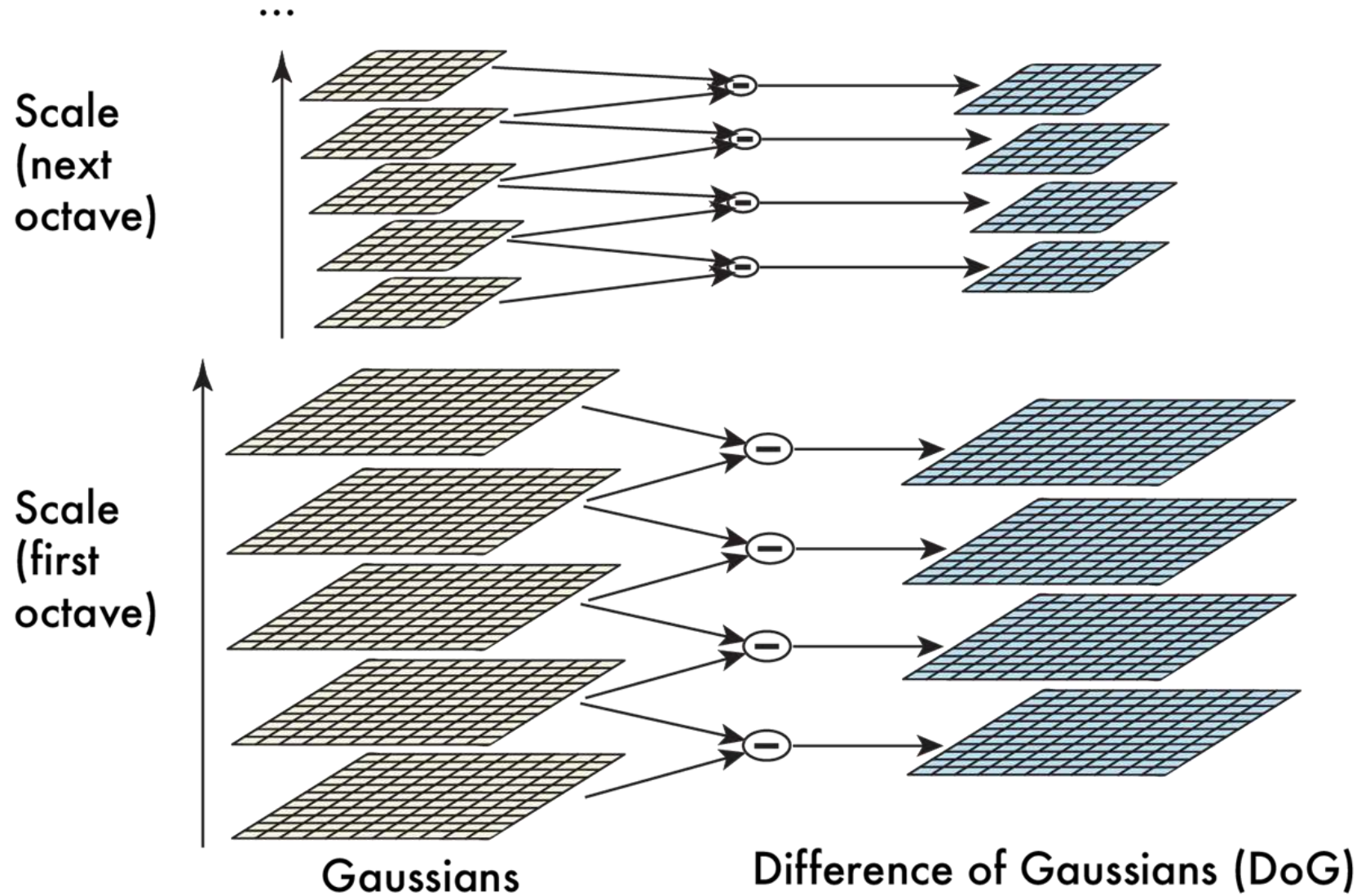
# Today's Agenda

- Basic feature descriptor and matching
- Histogram of Oriented Gradients
- SIFT
- Image transformations
- Estimate transformations

Want features invariant to scaling, rotation, etc.

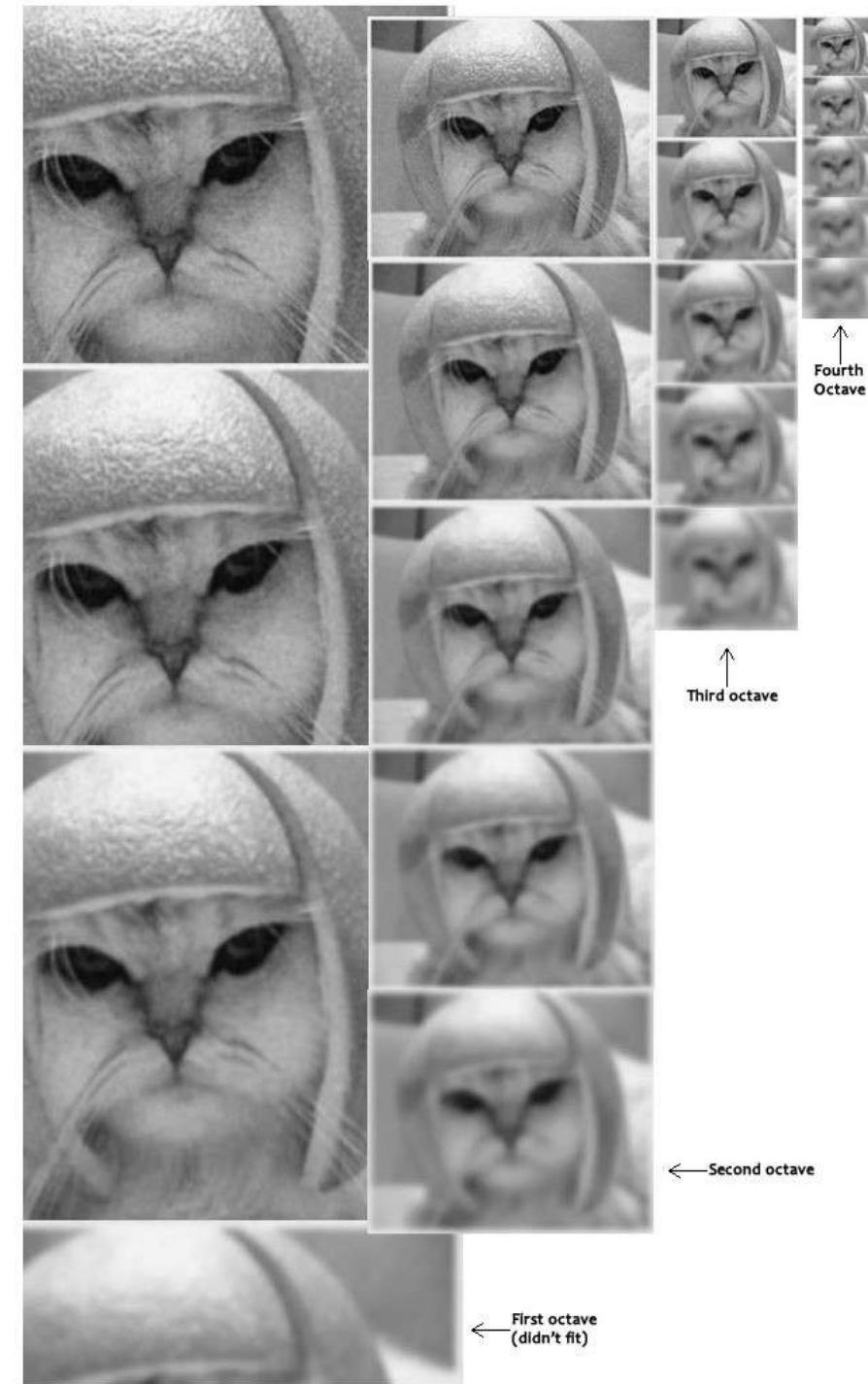
- Scale Invariant Feature Transform (SIFT)
  - Lowe et al. 2004, many images from that paper
- Get scale-invariant response map
- Find keypoints
- Extract rotation-invariant descriptors
  - Normalize based on orientation
  - Normalize based on lighting

## Extract DoG features at multiple scales



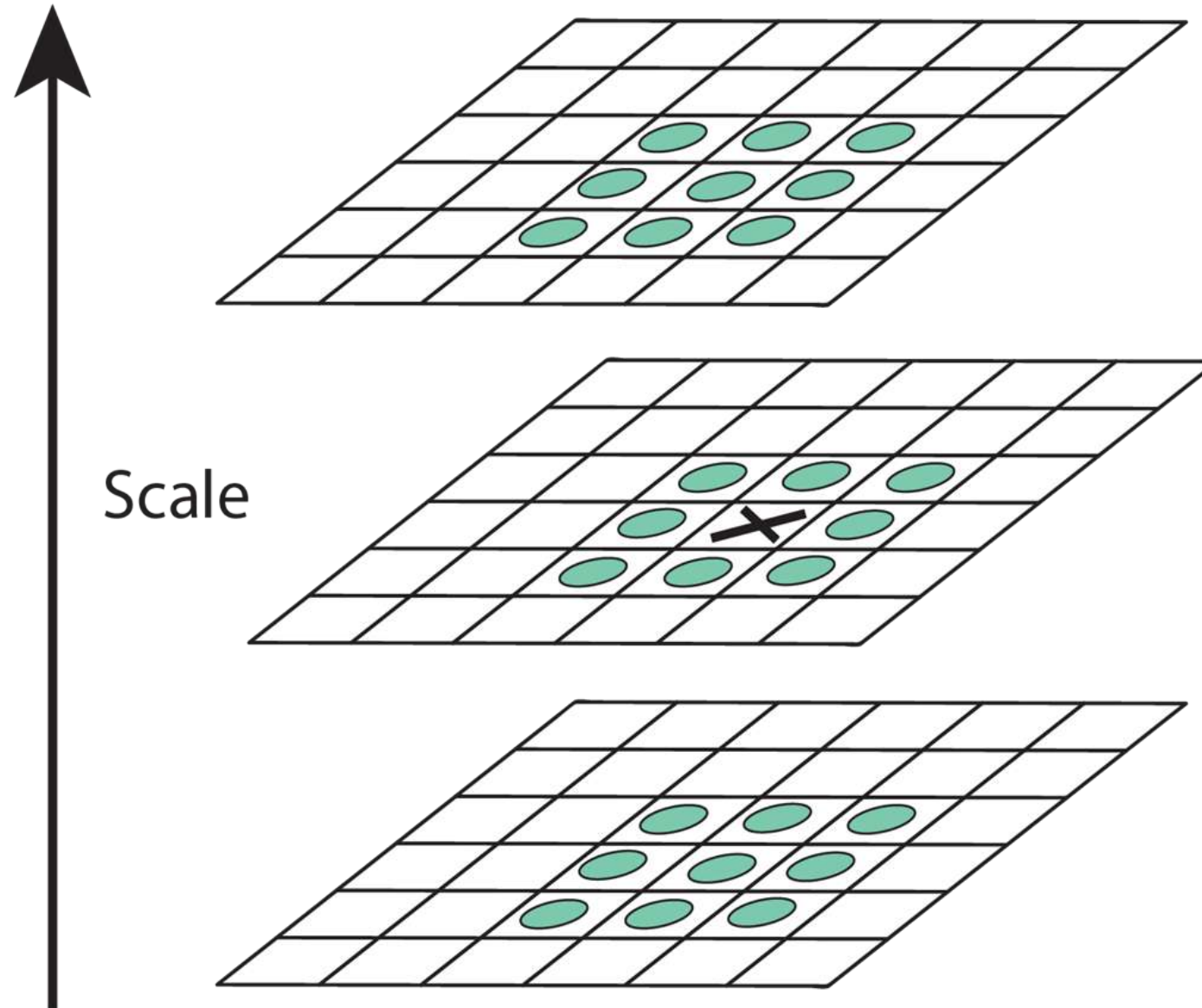
## Scale space

	scale →				
octave	0.707107	1.000000	1.414214	2.000000	2.828427
	1.414214	2.000000	2.828427	4.000000	5.656854
	2.828427	4.000000	5.656854	8.000000	11.313708
	5.656854	8.000000	11.313708	16.000000	22.627417





# Find local-maxima in location and scale



# Throw out weak responses and edges

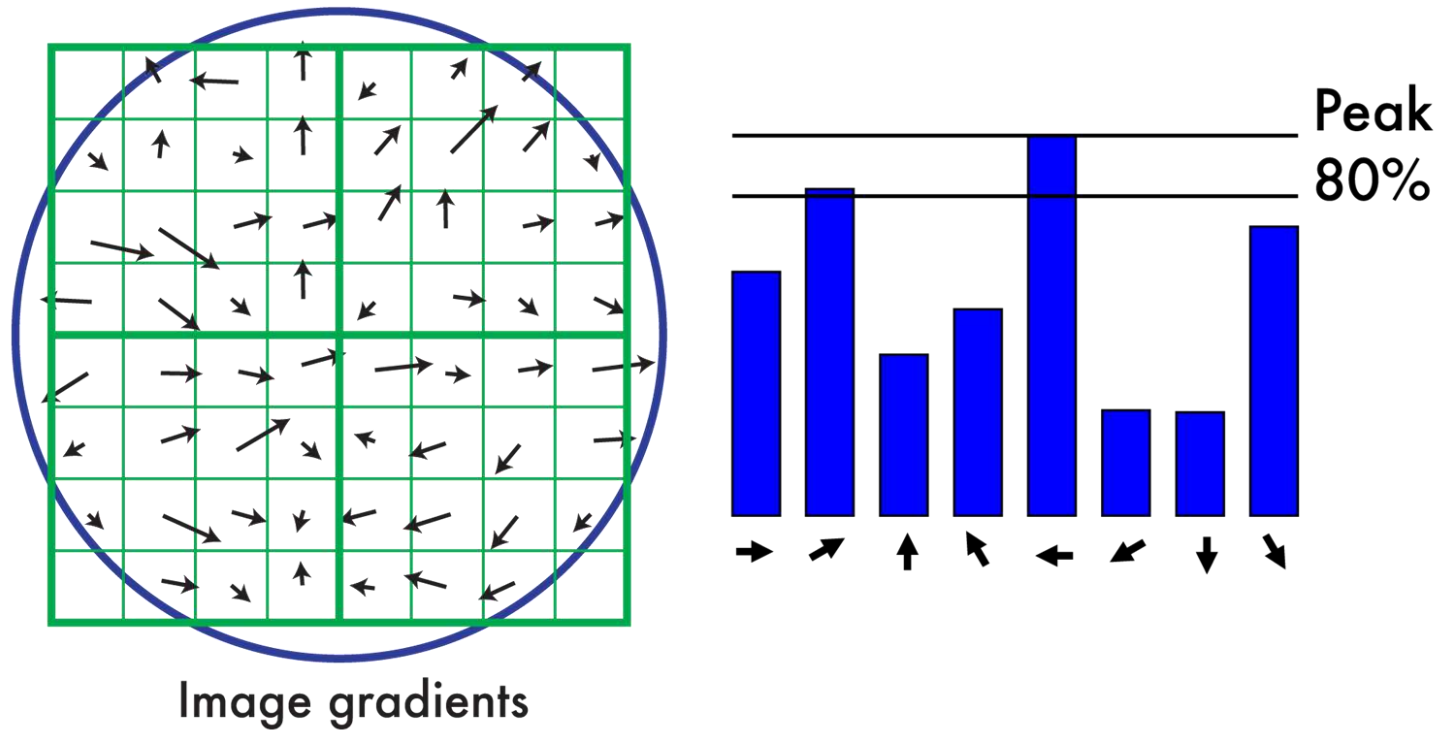
- Estimate gradients
  - Similar to before, look at nearby responses
  - Not whole image, only a few points! Faster!
  - Throw out weak responses
- Find cornery things
  - Same deal, structure matrix, use det and trace information
  - $r$  : ratio of larger to smaller eigenvalue

$$\frac{\text{Tr}(\mathbf{H})^2}{\text{Det}(\mathbf{H})} = \frac{(\alpha + \beta)^2}{\alpha\beta} = \frac{(r\beta + \beta)^2}{r\beta^2} = \frac{(r + 1)^2}{r};$$

$$\frac{\text{Tr}(\mathbf{H})^2}{\text{Det}(\mathbf{H})} < \frac{(r + 1)^2}{r}$$

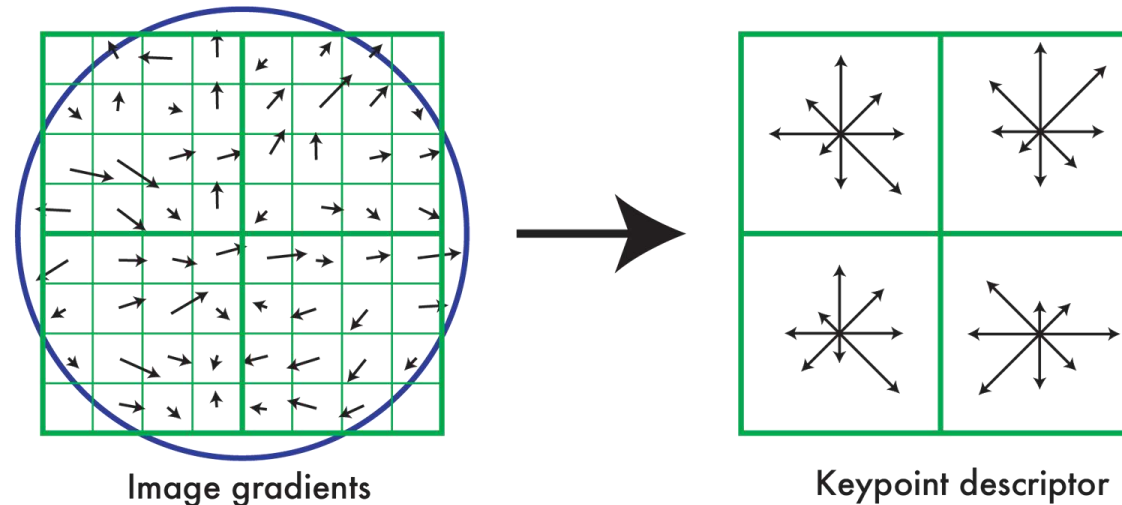
# Find main orientation of patches

- Look at weighted histogram of nearby gradients
  - Any gradient within 80% of peak gets its own descriptor
    - Multiple keypoints per pixel
  - Descriptors are normalized based on main orientation



## Keypoints are normalized gradient histograms

- Divide into subwindows (4x4)
- Bin gradients within subwindow, get histogram
  - Normalize to unit length
  - Clamp at maximum .2
  - Normalize again
  - Helps with lighting changes!



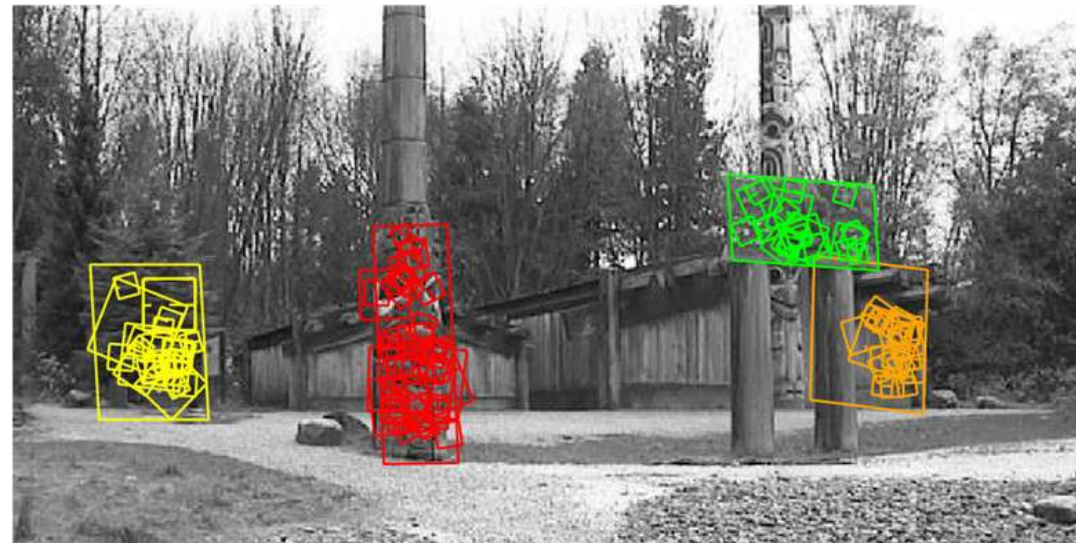


SIFT is great!

- Find good keypoints, describe them
- Finding objects, recognition, panoramas, etc.



## SIFT is great!



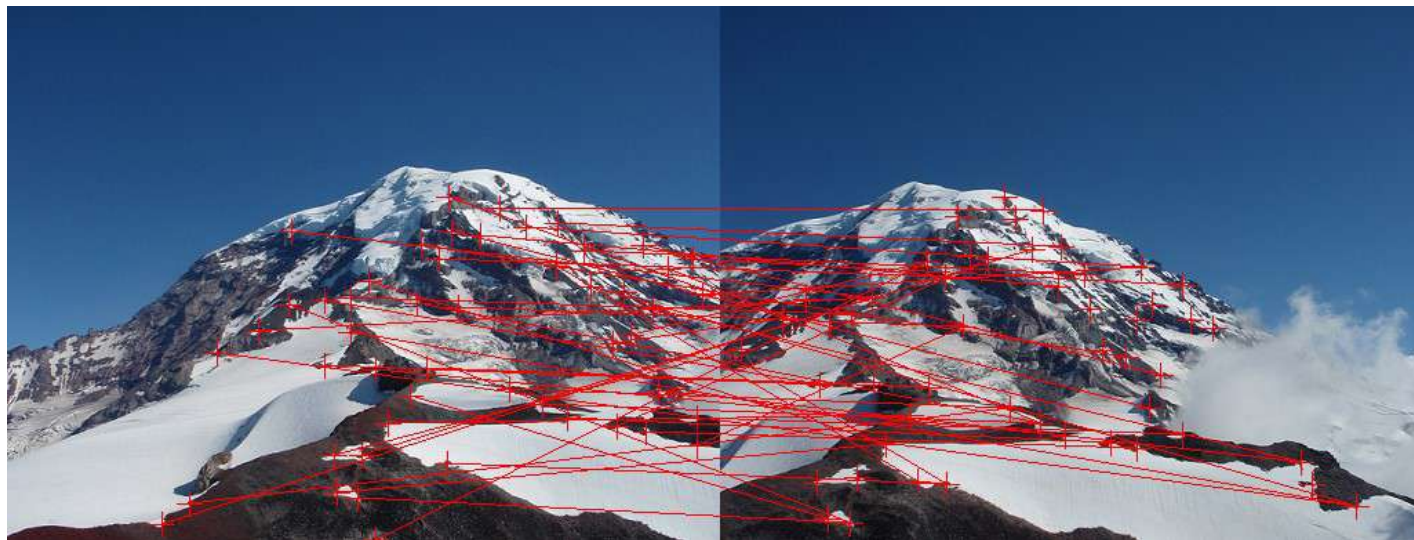
# Today's Agenda

- Basic feature descriptor and matching
- Histogram of Oriented Gradients
- SIFT
- Image transformations
- Estimate transformations



# Matching patches: descriptors!

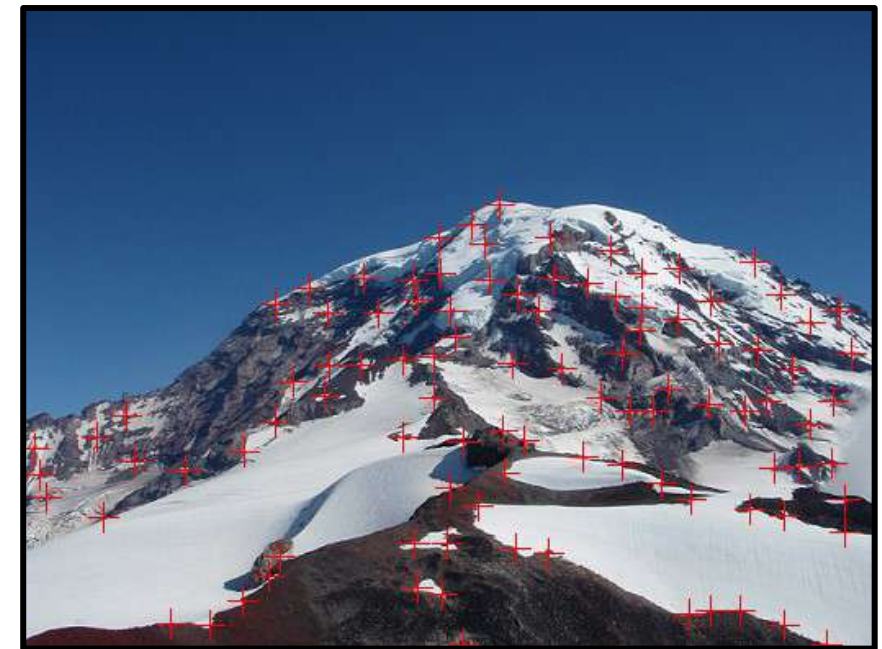
- Already have our patches that are likely “unique”-ish
- Loop over good patches in one image
  - Find best match in other image
- Do something with them?





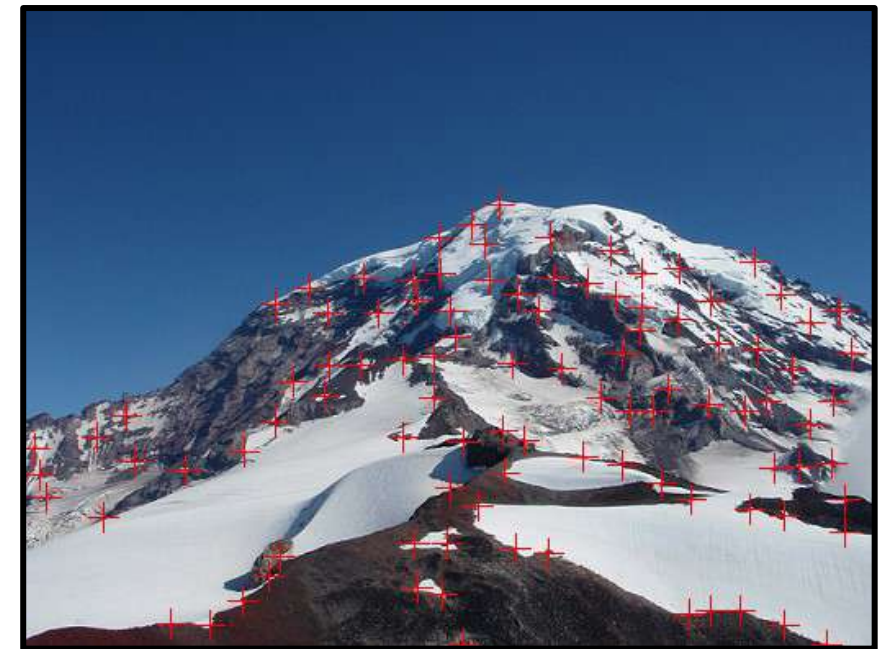
# Ok, we found corners, now what?

- Need to match image patches to each other
- Need to figure out transform between images



# Ok, we found corners, now what?

- Need to match image patches to each other
- Need to figure out transform between images
  - How can we transform images?
  - How do we solve for this transformation given matches?



# How can we transform images?

- Need to warp one image into the other
- Many different image transforms
  - Nested hierarchy of transformations

# How can we transform images?

- $\mathbf{x}$  is a point in our image where:
  - $\mathbf{x} = (x, y)$  or in matrix terms

$$\mathbf{x} = \begin{bmatrix} x \\ y \end{bmatrix}$$



# Say we want new coordinate system

- Map points from one image into another
- Often we can use matrix operations
- Given a point  $x$ , map to new point  $x'$  using  $M$

$$\mathbf{x}' = \mathbf{M} \mathbf{x}$$

# Scaling is just a matrix operation

- Map points from one image into another
- Often we can use matrix operations
- Given a point  $x$ , map to new point  $x'$  using  $M$

$$\mathbf{x}' = \mathbf{S} \mathbf{x}$$

$$\mathbf{x}' = \begin{bmatrix} \mathbf{S} & \mathbf{0} \\ \mathbf{0} & \mathbf{S} \end{bmatrix} \mathbf{x}$$

# Translation is harder...

- $x' = M x$ 
  - Want to move  $x'$  by  $dx$  and  $y'$  by  $dy$
  - How do we pick  $M$ ?
  - Can only add up multiples of  $x$  or  $y$
  - No easy way to add a constant!



# Translation: add another row

- $\bar{x}$  is  $x$  but with an added 1
- *Augmented vector*

$$\bar{x} = \begin{bmatrix} x \\ y \\ 1 \end{bmatrix}$$



# Translation: add another row

- $\bar{x}$  is  $x$  but with an added 1
- *Augmented vector*
- Now translation is easy

$$\bar{x} = \begin{bmatrix} x \\ y \\ 1 \end{bmatrix}$$

$$x' = [ I \quad t ] \bar{x}$$

# Reminder, I = Identity

Common to just use I as a generic, whatever size identity fits here.

$$\mathbf{I} = \begin{bmatrix} 1 & 0 & 0 & \dots & 0 \\ 0 & 1 & 0 & \dots & 0 \\ & & & \dots & \\ 0 & 0 & 0 & \dots & 1 \end{bmatrix} \quad \mathbf{I}_{2 \times 2} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \quad \mathbf{I}_{3 \times 3} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

# Translation: add another row

- $\bar{\mathbf{x}}$  is  $\mathbf{x}$  but with an added 1
- *Augmented vector*
- Now translation is easy

$$\bar{\mathbf{x}} = \begin{bmatrix} \mathbf{x} \\ y \\ 1 \end{bmatrix}$$

$$\mathbf{x}' = \begin{bmatrix} 1 & 0 & dx \\ 0 & 1 & dy \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ y \\ 1 \end{bmatrix}$$

$$\mathbf{x}' = [\mathbf{I} \ \mathbf{t}] \bar{\mathbf{x}}$$

# Translation: add another row

- $\bar{x}$  is  $x$  but with an added 1
- *Augmented vector*
- Now translation is easy
- $x' = 1*x + 0*y + dx*1$
- $y' = 0*x + 1*y + dy*1$

$$\bar{x} = \begin{bmatrix} x \\ y \\ 1 \end{bmatrix}$$

$$x' = \begin{bmatrix} 1 & 0 & dx \\ 0 & 1 & dy \end{bmatrix} \begin{bmatrix} x \\ y \\ 1 \end{bmatrix}$$

$$x' = [I \ t] \bar{x}$$



# Translation: add another row

- $\bar{x}$  is  $x$  but with an added 1
- *Augmented vector*
- Now translation is easy
- $x' = 1*x + 0*y + dx*1$
- $y' = 0*x + 1*y + dy*1$

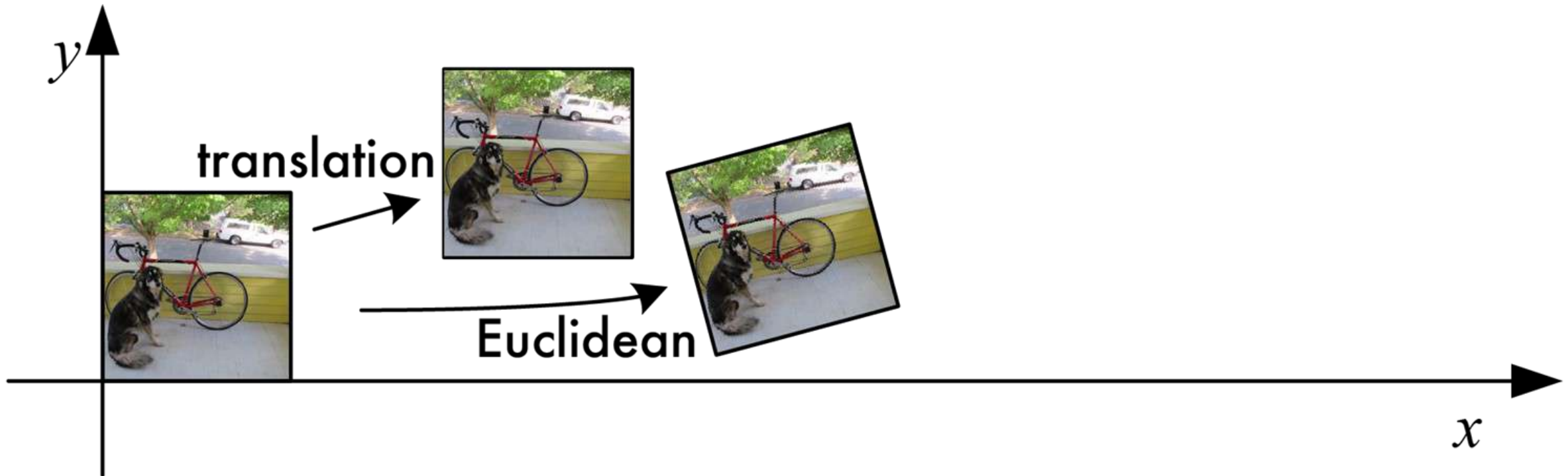
$$\bar{x} = \begin{bmatrix} x \\ y \\ 1 \end{bmatrix}$$

$$x' = \begin{bmatrix} 1 & 0 & dx \\ 0 & 1 & dy \end{bmatrix} \begin{bmatrix} x \\ y \\ 1 \end{bmatrix}$$

$$x' = [I \ t] \bar{x}$$

# Euclidean: rotation + translation

- Want to translate and rotate at same time
- Still just matrix operation



# Euclidean: rotation + translation

- Want to translate and rotate at same time
- Still just matrix operation

$$\mathbf{x}' = [ \mathbf{R} \quad \mathbf{t} ] \bar{\mathbf{x}}$$

# Euclidean: rotation + translation

- Want to translate and rotate at same time
- Still just matrix operation
- $\mathbf{R}$  is rotation matrix,  $\mathbf{t}$  is translation

$$\mathbf{x}' = [\mathbf{R} \quad \mathbf{t}] \bar{\mathbf{x}}$$

$$\mathbf{R} = \begin{bmatrix} \cos\theta & -\sin\theta \\ \sin\theta & \cos\theta \end{bmatrix}$$



# Euclidean: rotation + translation

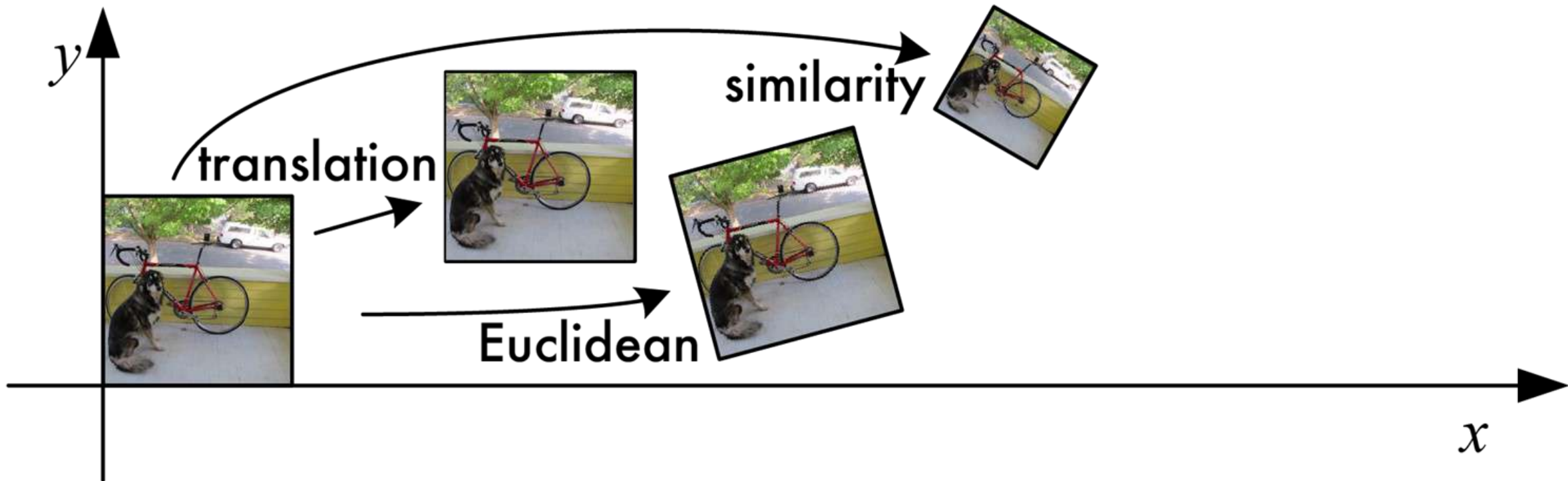
- Want to translate and rotate at same time
- Still just matrix operation
- $\mathbf{R}$  is rotation matrix,  $\mathbf{t}$  is translation

$$\mathbf{x}' = \begin{bmatrix} \cos\theta & -\sin\theta & dx \\ \sin\theta & \cos\theta & dy \\ & & 1 \end{bmatrix} \begin{bmatrix} x \\ y \\ 1 \end{bmatrix}$$

$$\mathbf{x}' = [\mathbf{R} \ \mathbf{t}] \bar{\mathbf{x}}$$

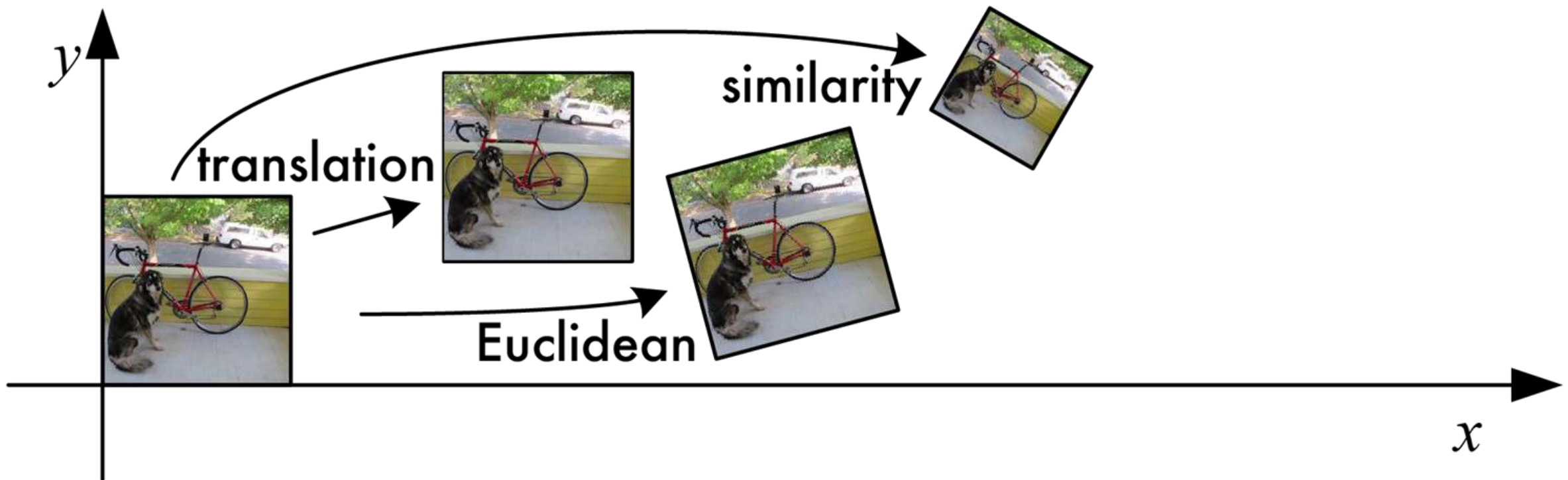
$$\mathbf{R} = \begin{bmatrix} \cos\theta & -\sin\theta \\ \sin\theta & \cos\theta \end{bmatrix}$$

# Similarity: scale, rotate, translate



# Similarity: scale, rotate, translate

$$\mathbf{x}' = [ \mathbf{sR} \quad \mathbf{t} ] \bar{\mathbf{x}}$$



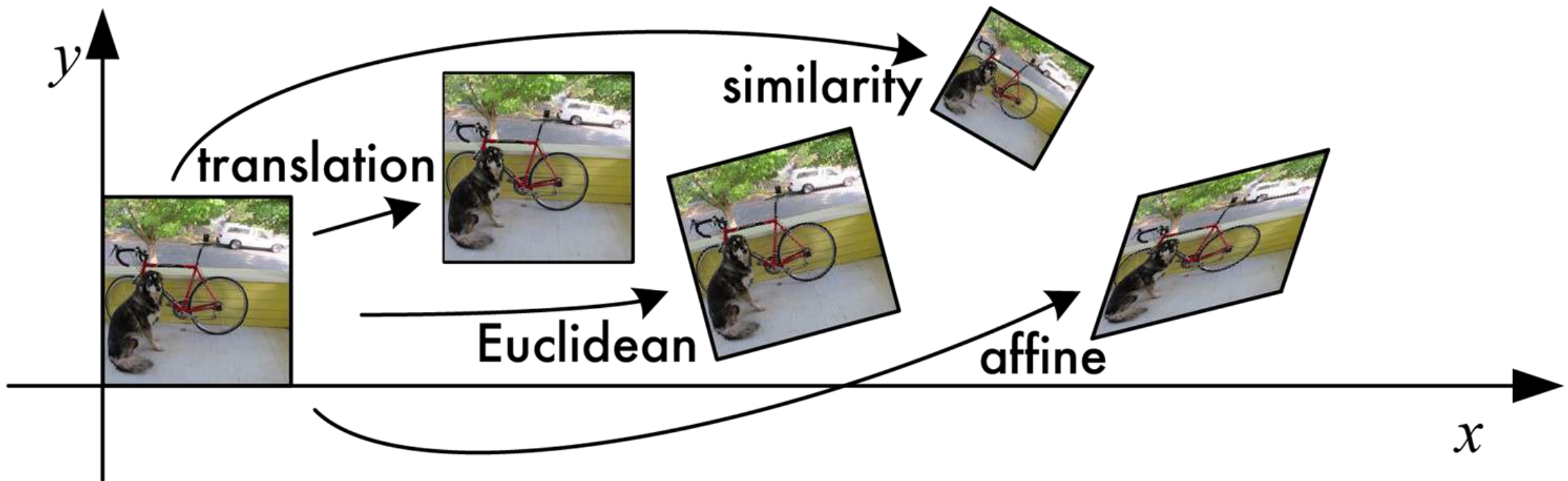
Similarity: scale, rotate, translate

$$\mathbf{x}' = [s\mathbf{R} \quad \mathbf{t}] \bar{\mathbf{x}}$$

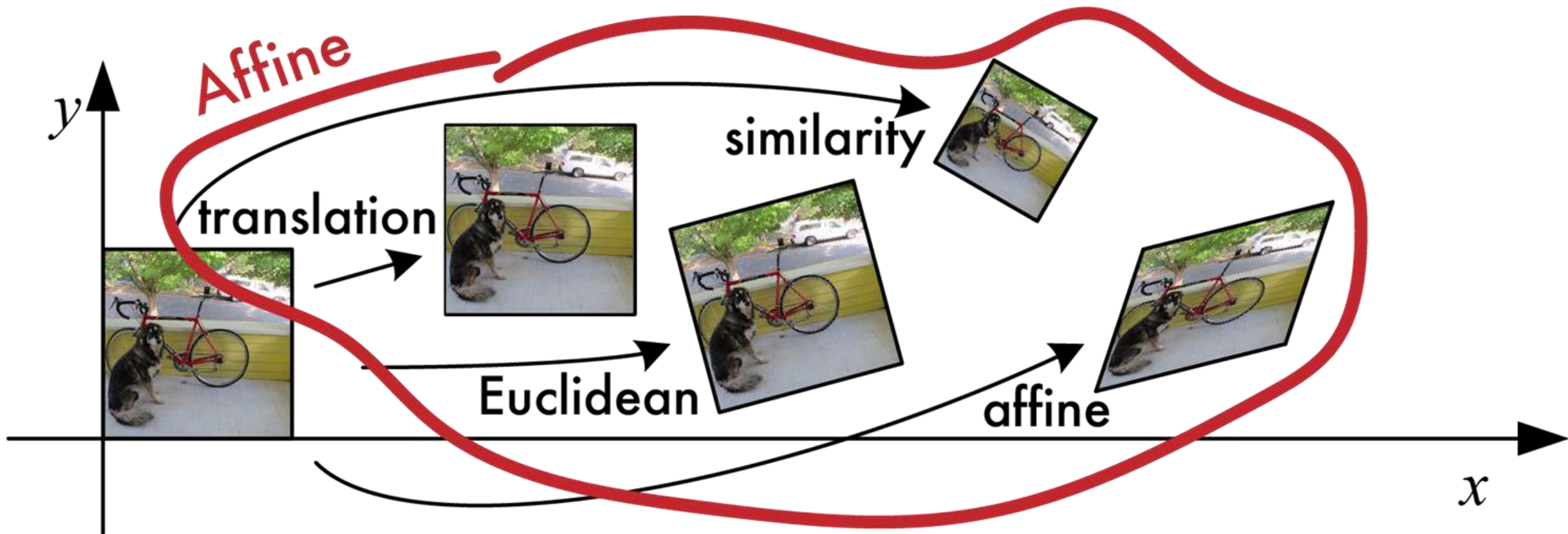
$$\mathbf{x}' = \begin{bmatrix} a & -b & dx \\ b & a & dy \end{bmatrix} \begin{bmatrix} x \\ y \\ 1 \end{bmatrix}$$



## Affine: scale, rotate, translate, shear



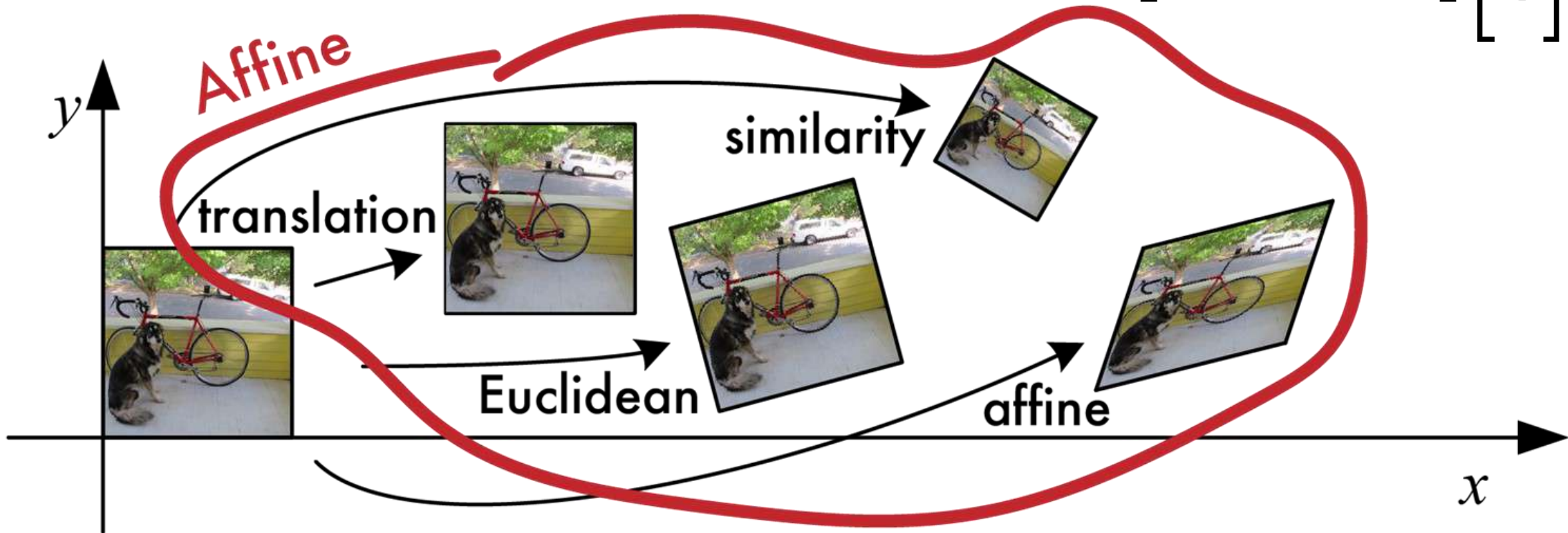
## Affine: scale, rotate, translate, shear



# Affine: scale, rotate, translate, shear

General case of 2x3 matrix

$$\mathbf{x}' = \begin{bmatrix} \mathbf{a}_{00} & \mathbf{a}_{01} & \mathbf{a}_{02} \\ \mathbf{a}_{10} & \mathbf{a}_{11} & \mathbf{a}_{12} \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \\ 1 \end{bmatrix}$$



# Combinations are still affine

Say you want to translate, then rotate, then translate back, then scale.

$$\mathbf{x}' = \mathbf{S} \mathbf{t} \mathbf{R} \mathbf{t} \bar{\mathbf{x}} = \mathbf{M} \bar{\mathbf{x}},$$

$$\text{If } \mathbf{M} = (\mathbf{S} \mathbf{t} \mathbf{R} \mathbf{t})$$

$\mathbf{M}$  is still affine transformation

Wait, but these are all 2x3, how to we multiply them together?





## Added row to transforms

$$\bar{\mathbf{x}}' = \begin{bmatrix} 1 & 0 & dx \\ 0 & 1 & dy \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x \\ y \\ 1 \end{bmatrix}$$

$$\bar{\mathbf{x}}' = \begin{bmatrix} \cos\theta & -\sin\theta & dx \\ \sin\theta & \cos\theta & dy \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x \\ y \\ 1 \end{bmatrix}$$

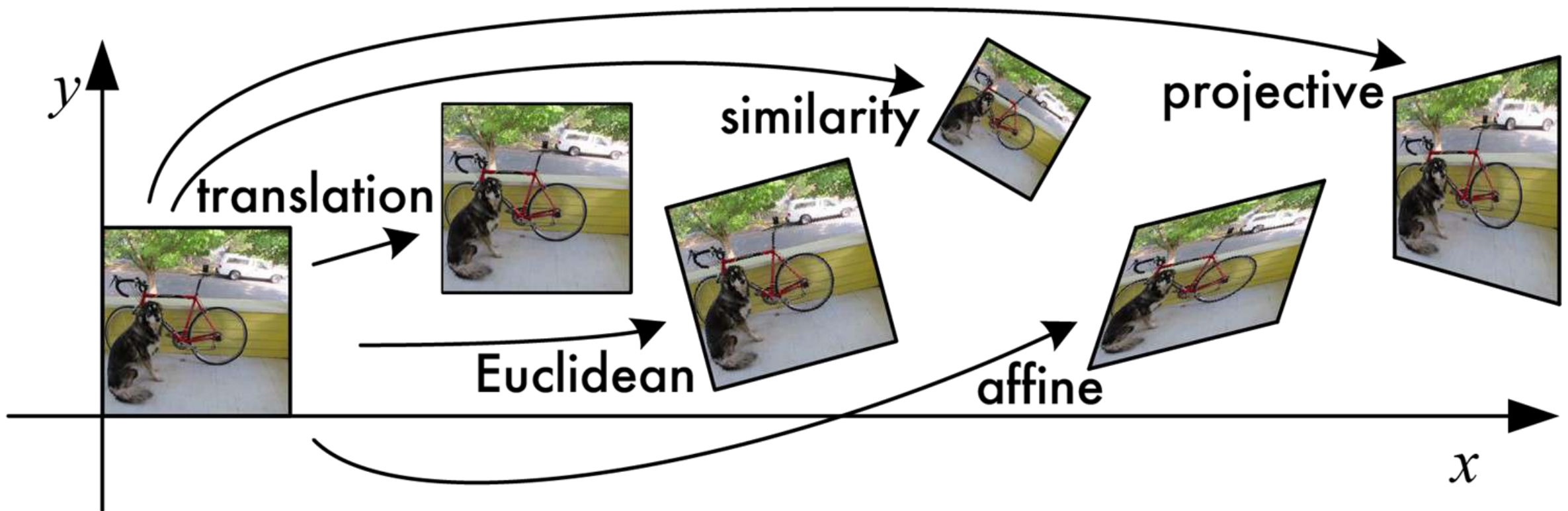
$$\bar{\mathbf{x}}' = \begin{bmatrix} a & -b & dx \\ b & a & dy \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x \\ y \\ 1 \end{bmatrix}$$

$$\bar{\mathbf{x}}' = \begin{bmatrix} a_{00} & a_{01} & a_{02} \\ a_{10} & a_{11} & a_{12} \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x \\ y \\ 1 \end{bmatrix}$$



## Projective transform

- Also known as homography
- Wait but affine was any 2x3 matrix...



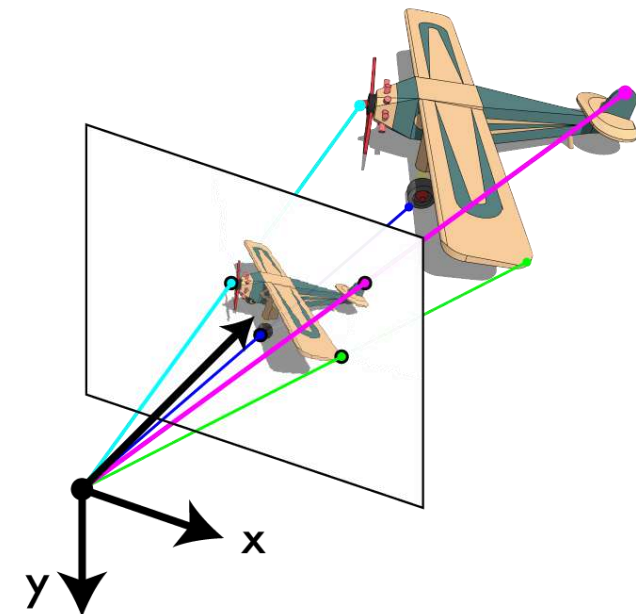
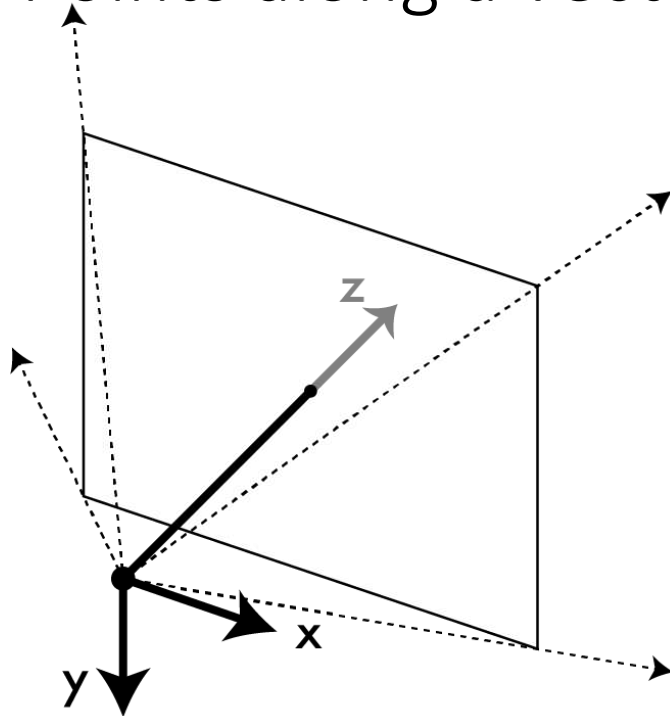
# Need some new coordinates!

- Homogeneous coordinate system
- Each point in 2d is actually a vector in 3d
- Equivalent up to scaling factor
- Have to normalize to get back to 2d

$$\tilde{\mathbf{x}} = \begin{bmatrix} \tilde{x} \\ \tilde{y} \\ \tilde{w} \end{bmatrix} \quad \bar{\mathbf{x}} = \tilde{\mathbf{x}} / \tilde{w}$$

# Why does this make sense?

- Remember our pinhole camera model
- Every point in 3d projects onto our viewing plane through our aperture
- Points along a vector are indistinguishable





# Projective transform

- Also known as homography
- Wait but affine was any 2x3 matrix...
- Homography is general 3x3 matrix
- Multiplication by scalar is equivalent

$$\tilde{\mathbf{x}}' = \tilde{\mathbf{H}} \tilde{\mathbf{x}}$$

# Projective transform

- Also known as homography
- Wait but affine was any 2x3 matrix...
- Homography is general 3x3 matrix
- Multiplication by scalar: equivalent projection

-  $3 \cdot H \sim H$

$$\tilde{\mathbf{x}}' = \begin{bmatrix} h_{00} & h_{01} & h_{02} \\ h_{10} & h_{11} & h_{12} \\ h_{20} & h_{21} & h_{22} \end{bmatrix} \begin{bmatrix} \tilde{x} \\ \tilde{y} \\ \tilde{w} \end{bmatrix}$$

$$\tilde{\mathbf{x}}' = \tilde{\mathbf{H}} \tilde{\mathbf{x}}$$

# Using homography to project point

- Multiply  $\tilde{\mathbf{x}}$  by  $\tilde{\mathbf{H}}$  to get  $\tilde{\mathbf{x}}'$
- Convert to  $\mathbf{x}'$  by dividing by  $w'$

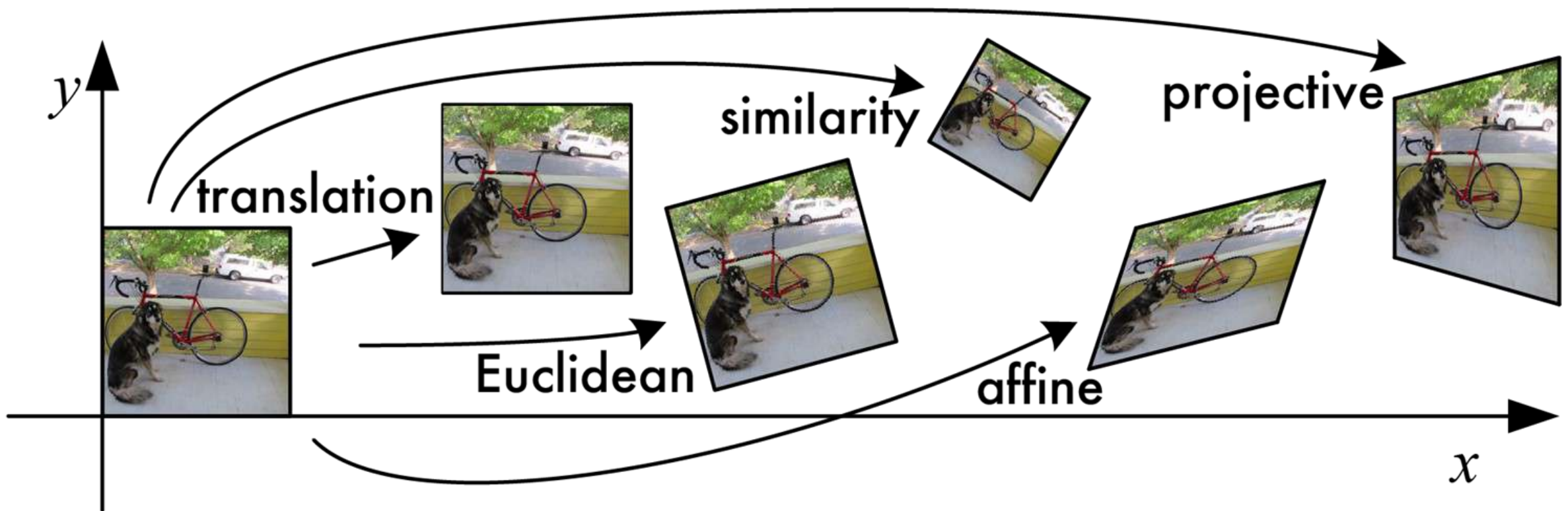
$$\tilde{\mathbf{x}}' = \tilde{\mathbf{H}} \tilde{\mathbf{x}}$$

$$\begin{bmatrix} \tilde{x}' \\ \tilde{y}' \\ \tilde{w}' \end{bmatrix} = \begin{bmatrix} h_{00} & h_{01} & h_{02} \\ h_{10} & h_{11} & h_{12} \\ h_{20} & h_{21} & h_{22} \end{bmatrix} \begin{bmatrix} \tilde{x} \\ \tilde{y} \\ \tilde{w} \end{bmatrix}$$

$$\bar{\mathbf{x}} = \tilde{\mathbf{x}} / \tilde{w}$$

## Lots to choose from

- What do each of them do?
- Which is right for panorama stitching?





# Today's Agenda

- Basic descriptor and matching
- Image transformations
- Estimate transformations

# How hard are they to recover?

$$\bar{\mathbf{x}}' = \begin{bmatrix} 1 & 0 & dx \\ 0 & 1 & dy \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x \\ y \\ 1 \end{bmatrix} \quad \bar{\mathbf{x}}' = \begin{bmatrix} \cos\theta & -\sin\theta & dx \\ \sin\theta & \cos\theta & dy \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x \\ y \\ 1 \end{bmatrix}$$

$$\bar{\mathbf{x}}' = \begin{bmatrix} a & -b & dx \\ b & a & dy \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x \\ y \\ 1 \end{bmatrix}$$

$$\bar{\mathbf{x}}' = \begin{bmatrix} a_{00} & a_{01} & a_{02} \\ a_{10} & a_{11} & a_{12} \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x \\ y \\ 1 \end{bmatrix}$$

$$\tilde{\mathbf{x}}' = \begin{bmatrix} h_{00} & h_{01} & h_{02} \\ h_{10} & h_{11} & h_{12} \\ h_{20} & h_{21} & h_{22} \end{bmatrix} \begin{bmatrix} \tilde{x} \\ \tilde{y} \\ \tilde{w} \end{bmatrix}$$



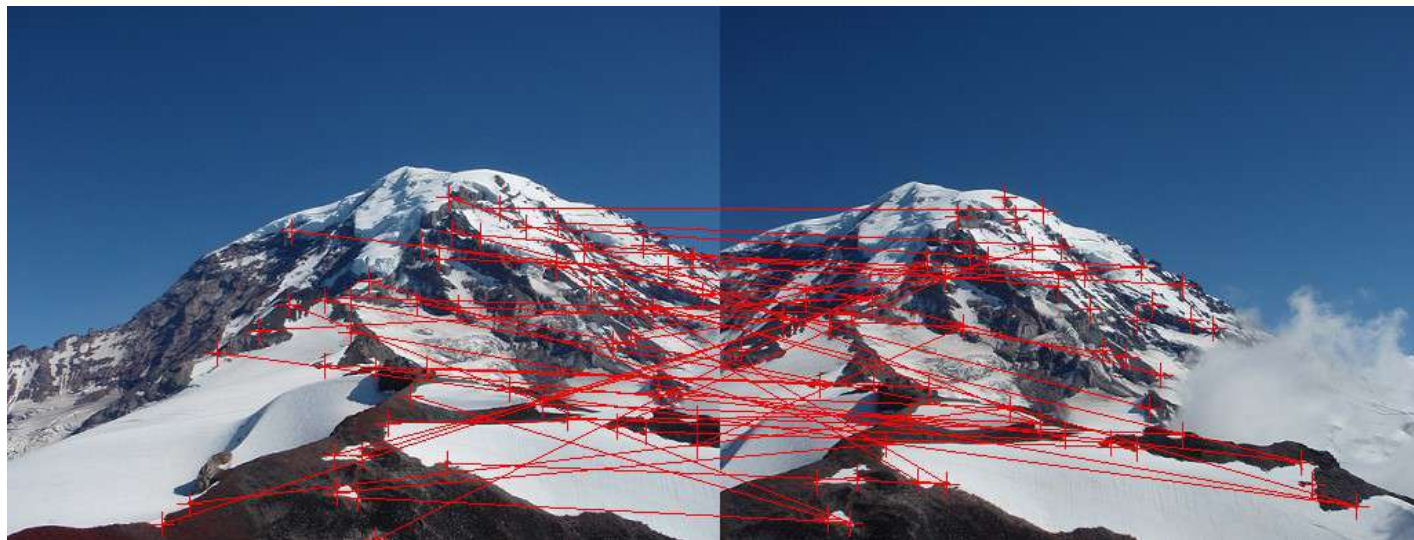
## Lots to choose from

Transformation	Matrix	# DoF	Preserves	Icon
translation	$\begin{bmatrix} \mathbf{I} &   & \mathbf{t} \end{bmatrix}_{2 \times 3}$	2	orientation	
rigid (Euclidean)	$\begin{bmatrix} \mathbf{R} &   & \mathbf{t} \end{bmatrix}_{2 \times 3}$	3	lengths	
similarity	$\begin{bmatrix} s\mathbf{R} &   & \mathbf{t} \end{bmatrix}_{2 \times 3}$	4	angles	
affine	$\begin{bmatrix} \mathbf{A} \end{bmatrix}_{2 \times 3}$	6	parallelism	
projective	$\begin{bmatrix} \tilde{\mathbf{H}} \end{bmatrix}_{3 \times 3}$	8	straight lines	



# Say we want affine transformation

- Have our matched points
- Want to estimate  $\mathbf{A}$  that maps from  $\mathbf{x}$  to  $\mathbf{x}'$
- $\mathbf{Ax} = \mathbf{x}'$





# Say we want affine transformation

- Have our matched points
- Want to estimate  $\mathbf{A}$  that maps from  $\mathbf{x}$  to  $\mathbf{x}'$
- $\mathbf{Ax} = \mathbf{x}'$
- How many degrees of freedom?

# Say we want affine transformation

- Have our matched points
- Want to estimate  $\mathbf{A}$  that maps from  $\mathbf{x}$  to  $\mathbf{x}'$
- $\mathbf{Ax} = \mathbf{x}'$
- How many degrees of freedom?
  - 6
- How many knowns do we get with one match?

$$\mathbf{x}' = \begin{bmatrix} \mathbf{a}_{00} & \mathbf{a}_{01} & \mathbf{a}_{02} \\ \mathbf{a}_{10} & \mathbf{a}_{11} & \mathbf{a}_{12} \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \\ 1 \end{bmatrix}$$

# Say we want affine transformation

- Have our matched points
- Want to estimate  $\mathbf{A}$  that maps from  $\mathbf{x}$  to  $\mathbf{x}'$
- $\mathbf{Ax} = \mathbf{x}'$
- How many degrees of freedom?
  - 6
- How many knowns do we get with one match?
  - 2
  - $n_x = a_{00} * m_x + a_{01} * m_y + a_{02} * 1$
  - $n_y = a_{10} * m_x + a_{11} * m_y + a_{12} * 1$

$$\mathbf{x}' = \begin{bmatrix} a_{00} & a_{01} & a_{02} \\ a_{10} & a_{11} & a_{12} \end{bmatrix} \begin{bmatrix} x \\ y \\ 1 \end{bmatrix}$$

# Say we want affine transformation

- How many knowns do we get with one match?
  - $n_x = a_{00} * m_x + a_{01} * m_y + a_{02} * 1$
  - $n_y = a_{10} * m_x + a_{11} * m_y + a_{12} * 1$
  - Solve linear system of equations  $M a = b$ 
    - $M^{-1} M a = M^{-1} b \Rightarrow a = M^{-1} b$
    - But  $M^{-1}$  does not exist in general - Why?
- Still works if overdetermined
  - Why???
  - Pseudoinverse – least squares solution
  - $M^T M a = M^T b$
  - $(M^T M)^{-1} (M^T M) a = (M^T M)^{-1} M^T b$
  - $\Rightarrow a = (M^T M)^{-1} M^T b$

$$\begin{matrix} & \mathbf{M} & & \mathbf{a} & & \mathbf{b} \\ & & & & & = \\ \left[ \begin{array}{cccccc} m_{x1} & m_{y1} & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & m_{x1} & m_{y1} & 1 \\ m_{x2} & m_{y2} & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & m_{x2} & m_{y2} & 1 \\ m_{x3} & m_{y3} & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & m_{x3} & m_{y3} & 1 \\ & & \dots & & & \\ & & \dots & & & \end{array} \right] & \left[ \begin{array}{c} a_{00} \\ a_{01} \\ a_{02} \\ a_{10} \\ a_{11} \\ a_{12} \end{array} \right] & & \left[ \begin{array}{c} n_{x1} \\ n_{y1} \\ n_{x2} \\ n_{y2} \\ n_{x3} \\ n_{y3} \\ \dots \\ \dots \end{array} \right]
 \end{matrix}$$



**MAI4CAREU**

Master programmes in Artificial  
Intelligence 4 Careers in Europe



# Thank you.





University of Cyprus – MSc Artificial Intelligence

# MAI644 – COMPUTER VISION

## Lecture 9: RANSAC, Panorama Stitching

**Melinos Averkiou**

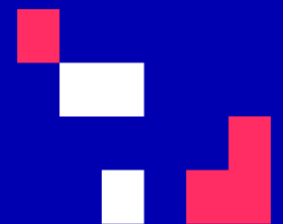
CYENS Centre of Excellence

University of Cyprus - Department of Computer Science

[m.averkiou@cyens.org.cy](mailto:m.averkiou@cyens.org.cy)



**CYENS**  
CENTRE OF EXCELLENCE



# Last time

- Basic feature descriptor and matching
- Histogram of Oriented Gradients
- SIFT
- Image transformations
- Estimate transformations

# Today's Agenda

- Linear least-squares
- RANSAC
- Panorama Stitching

**[material based on Joseph Redmon's course]**



# Today's Agenda

- Linear least-squares
- RANSAC
- Panorama Stitching

# Linear least squares

Want to solve overdetermined linear system:

$$- \quad M a = b$$

Want to minimize squared error:

$$|| b - M a ||^2 =$$

# Linear least squares

Want to solve overdetermined linear system:

$$- \quad M a = b$$

Want to minimize squared error:

$$|| b - M a ||^2 =$$

$$(b - M a)^T (b - M a)$$

# Linear least squares

Want to solve overdetermined linear system:

$$- \quad M a = b$$

Want to minimize squared error:

$$|| b - M a ||^2 =$$

$$(b - M a)^T (b - M a) =$$

$$b^T b - a^T M^T b - b^T M a + a^T M^T M a$$



# Linear least squares

Want to solve overdetermined linear system:

$$- \quad M a = b$$

Want to minimize squared error:

$$|| b - M a ||^2 =$$

$$(b - M a)^T (b - M a) =$$

$$b^T b - a^T M^T b - b^T M a + a^T M^T M a =$$

$$b^T b - 2a^T M^T b + a^T M^T M a$$

# Linear least squares

Want to minimize squared error:  $\|b - M a\|^2 =$

$$b^T b - 2a^T M^T b + a^T M^T M a$$

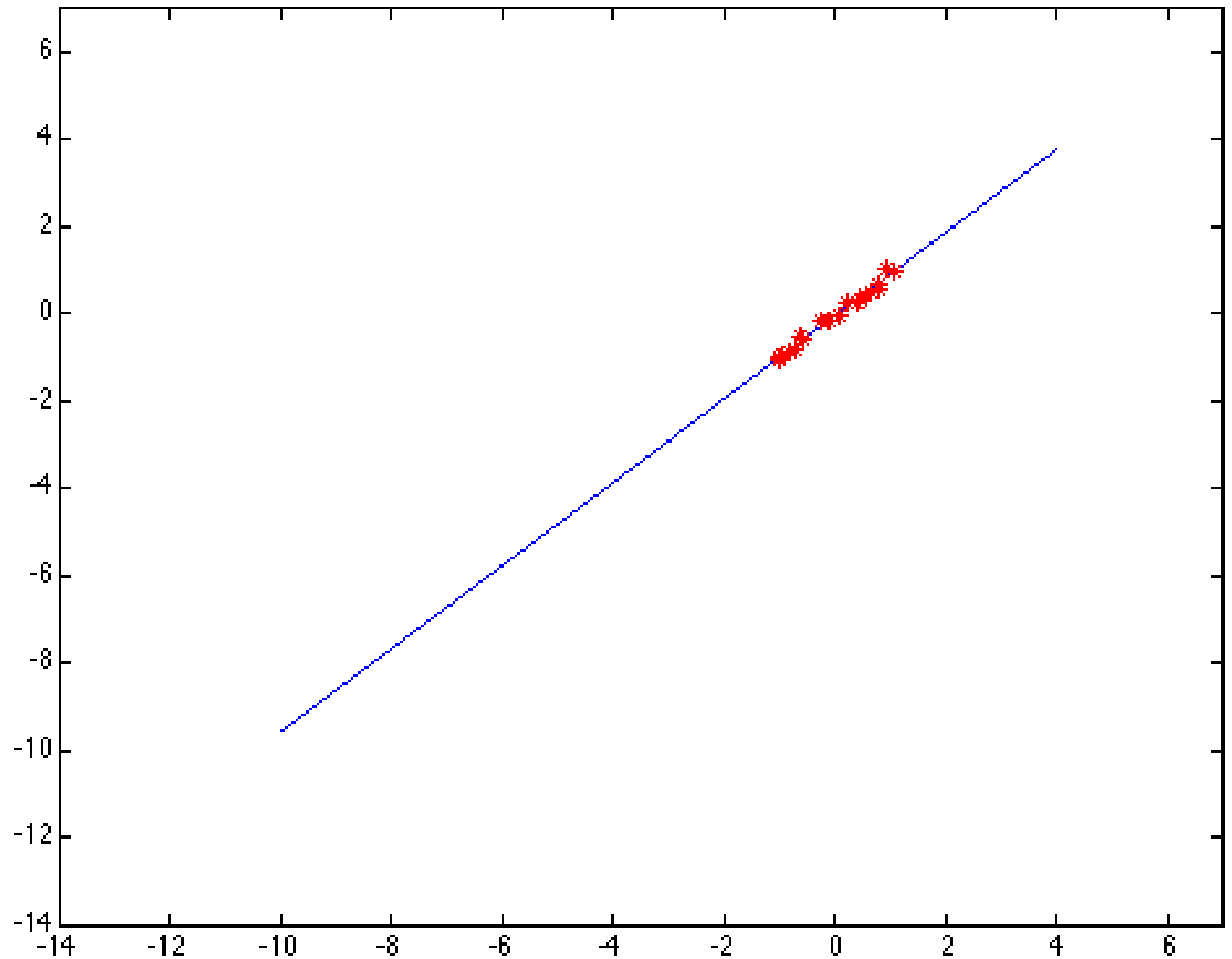
This is convex and minimized when gradient = 0. So we take the derivative wrt  $a$  and set = 0.

$$-M^T b + (M^T M) a = 0$$

$$(M^T M) a = M^T b$$

$$a = (M^T M)^{-1} M^T b$$

# So what does linear least squares do?

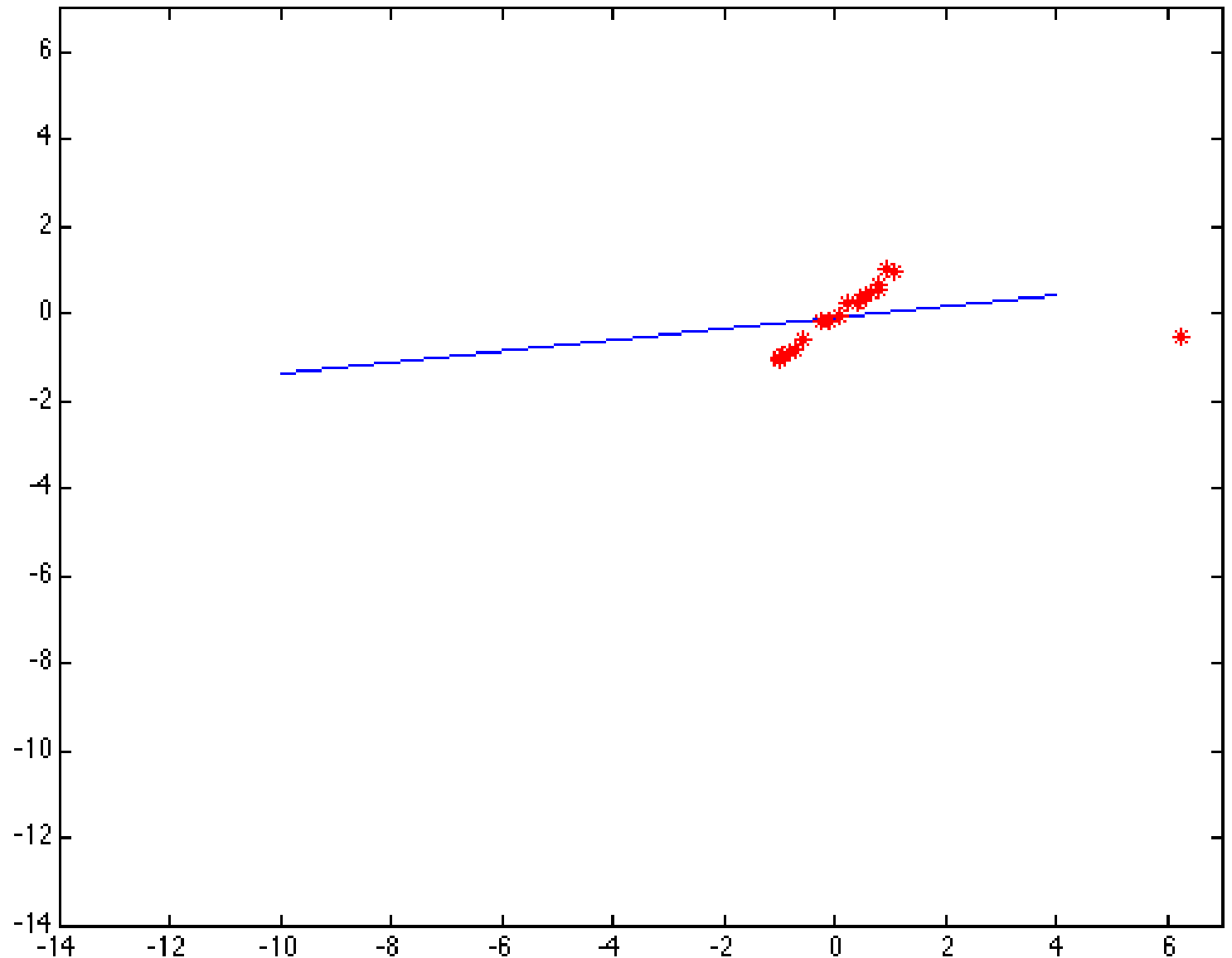


# So what does linear least squares do?

Error based on **squared** residual

Very scared of being wrong, even for just one point

Very bad at handling outliers in data





## Not a problem for us, our data is perfect...



Not really ...



# Today's Agenda

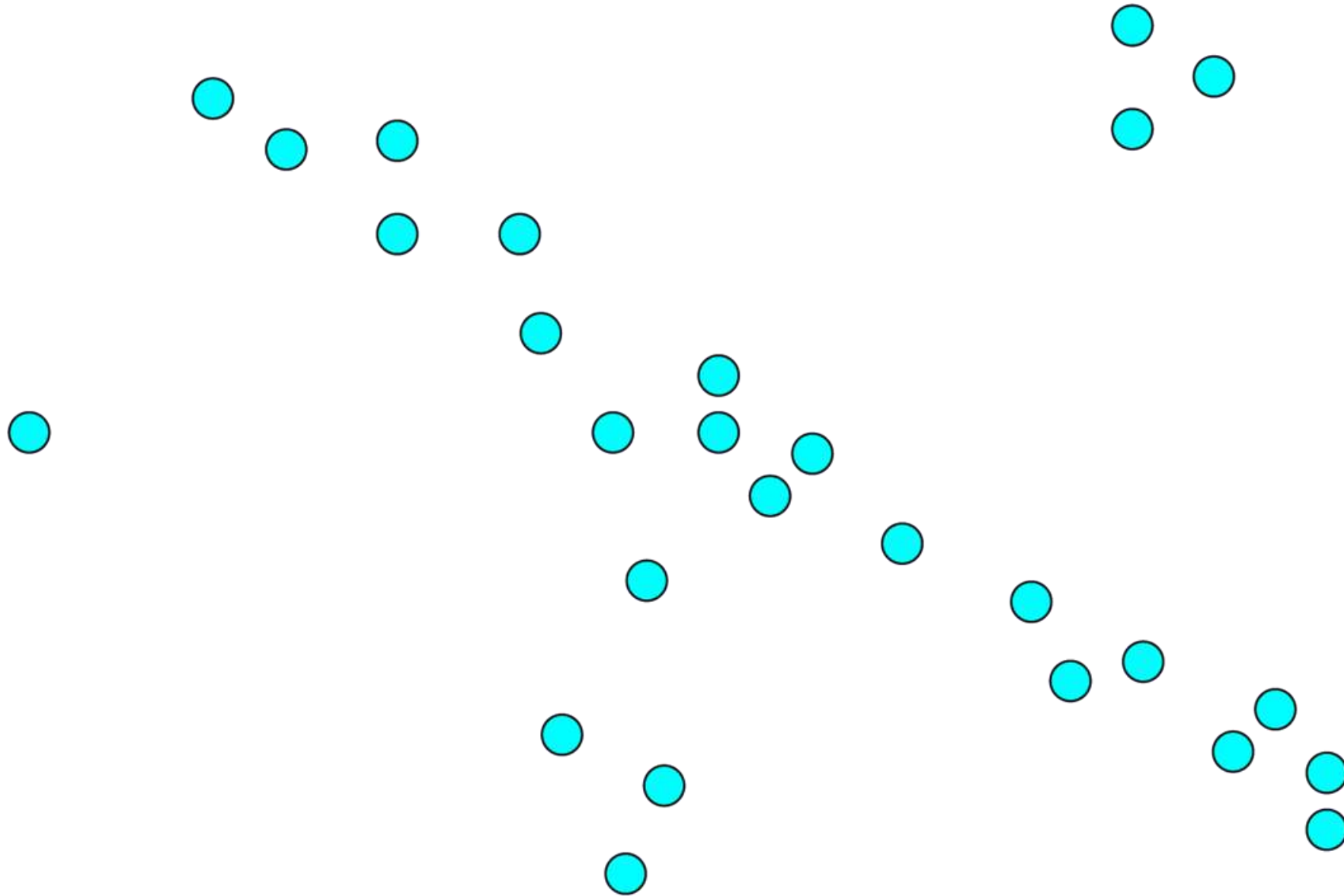
- Linear least-squares
- RANSAC
- Panorama Stitching



## RANSAC: RANdom SAMple Consensus

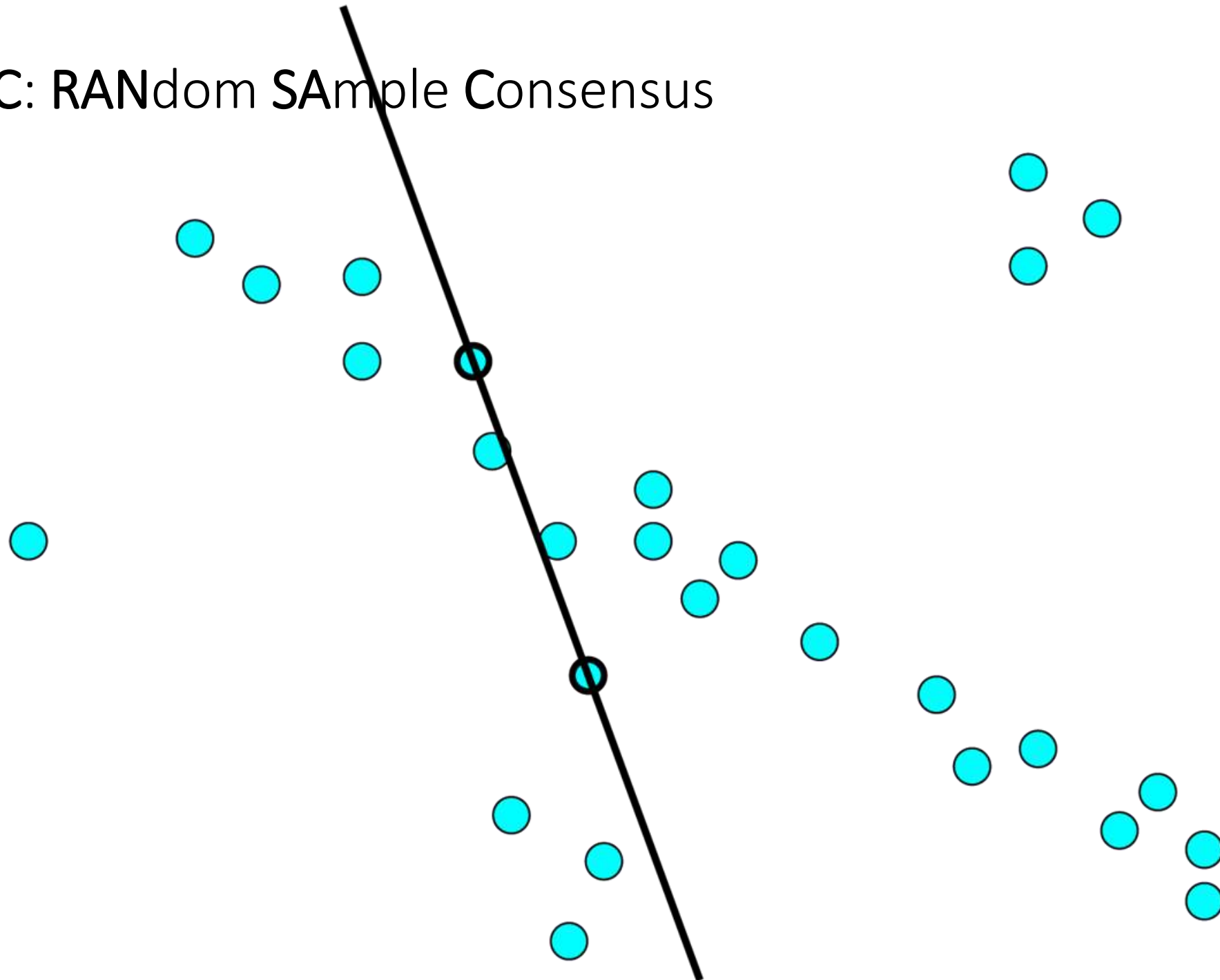
- How can we fit model to inliers but ignore outliers?
- Try a bunch of models, see which ones are best!
- Inliers will all agree on a model
- Outliers are basically random, will not agree

## RANSAC: RANdom SAMple Consensus



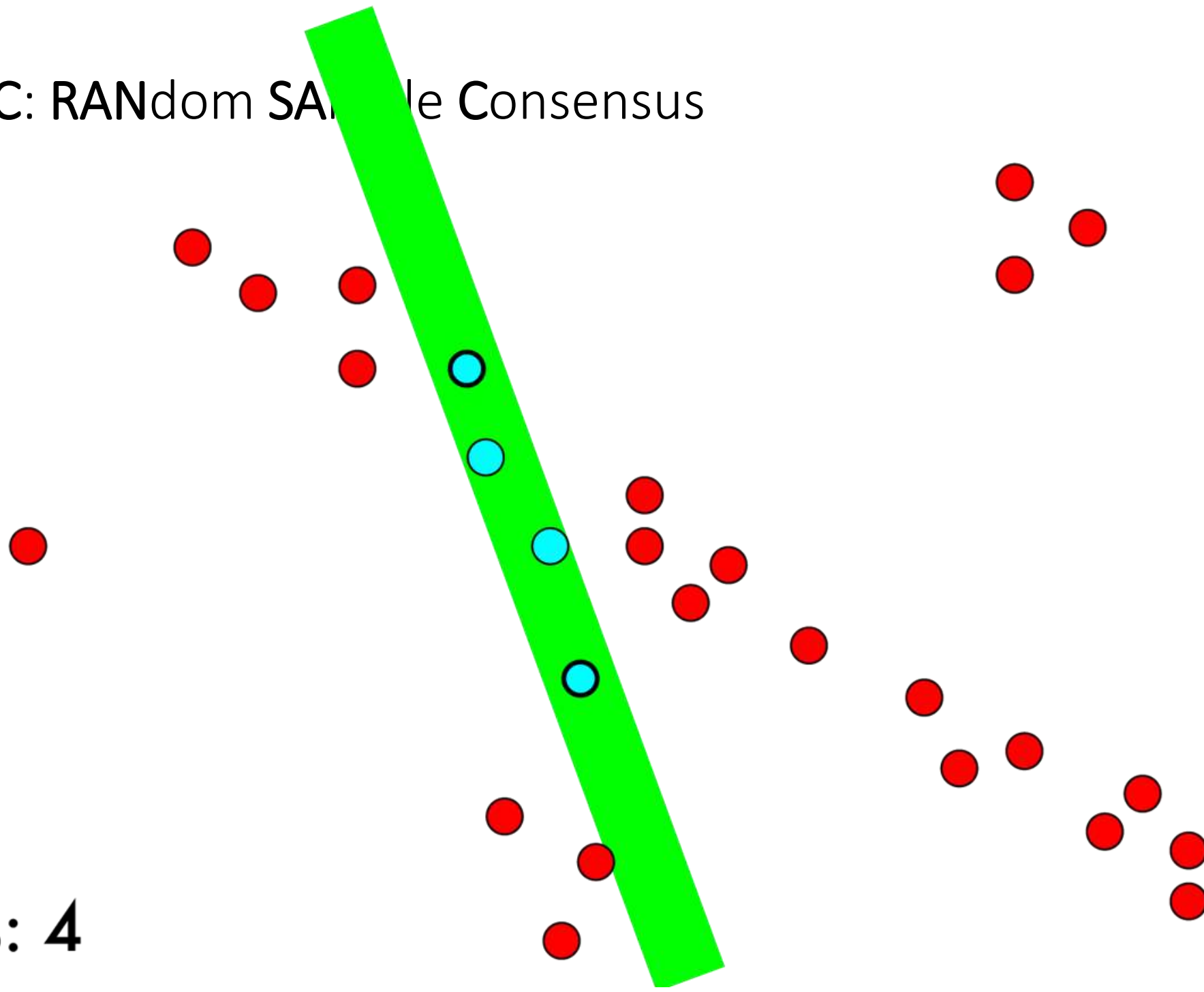


## RANSAC: RANdOm SAMple Consensus





## RANSAC: RANdOm SAMple Consensus

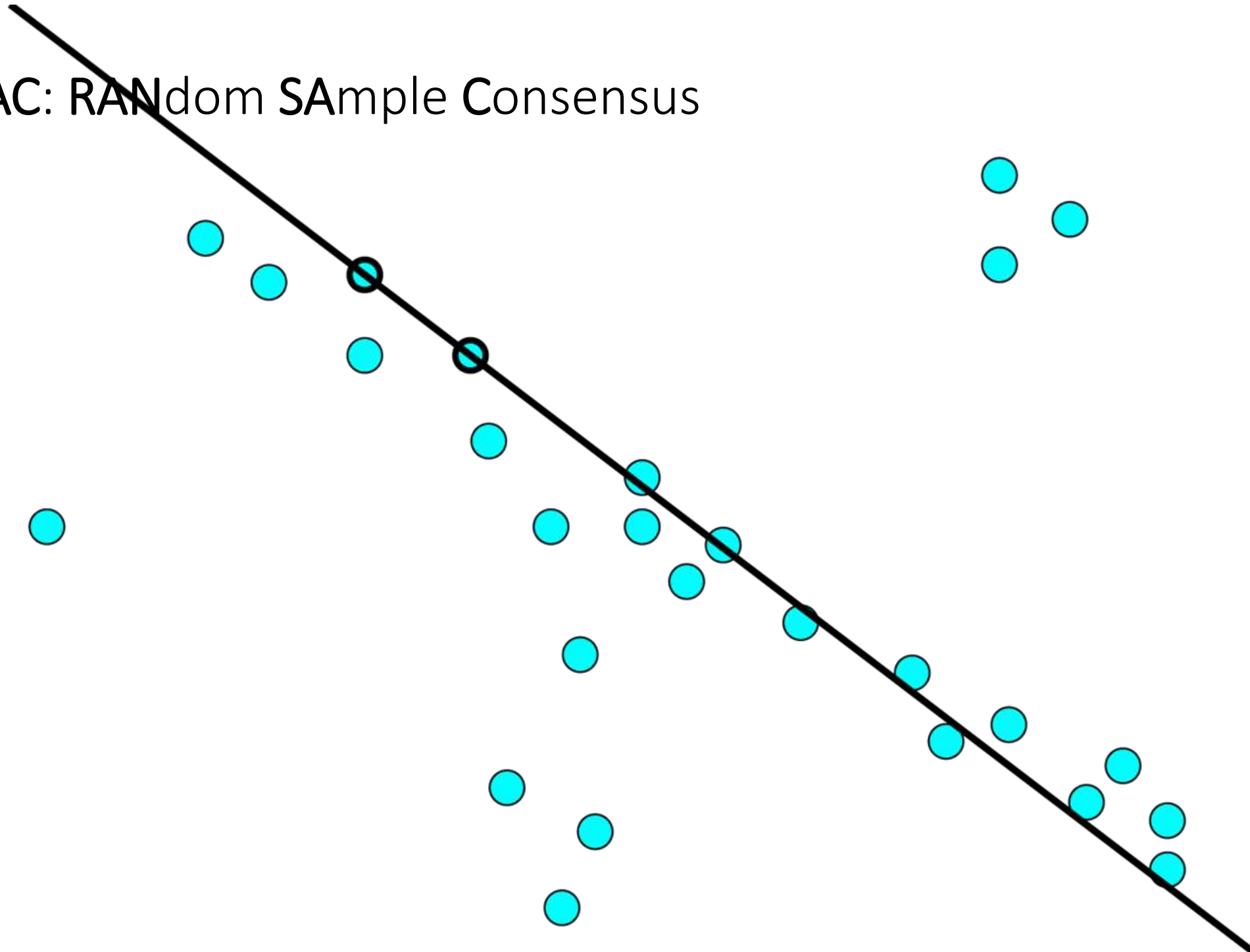


Inliers: 4

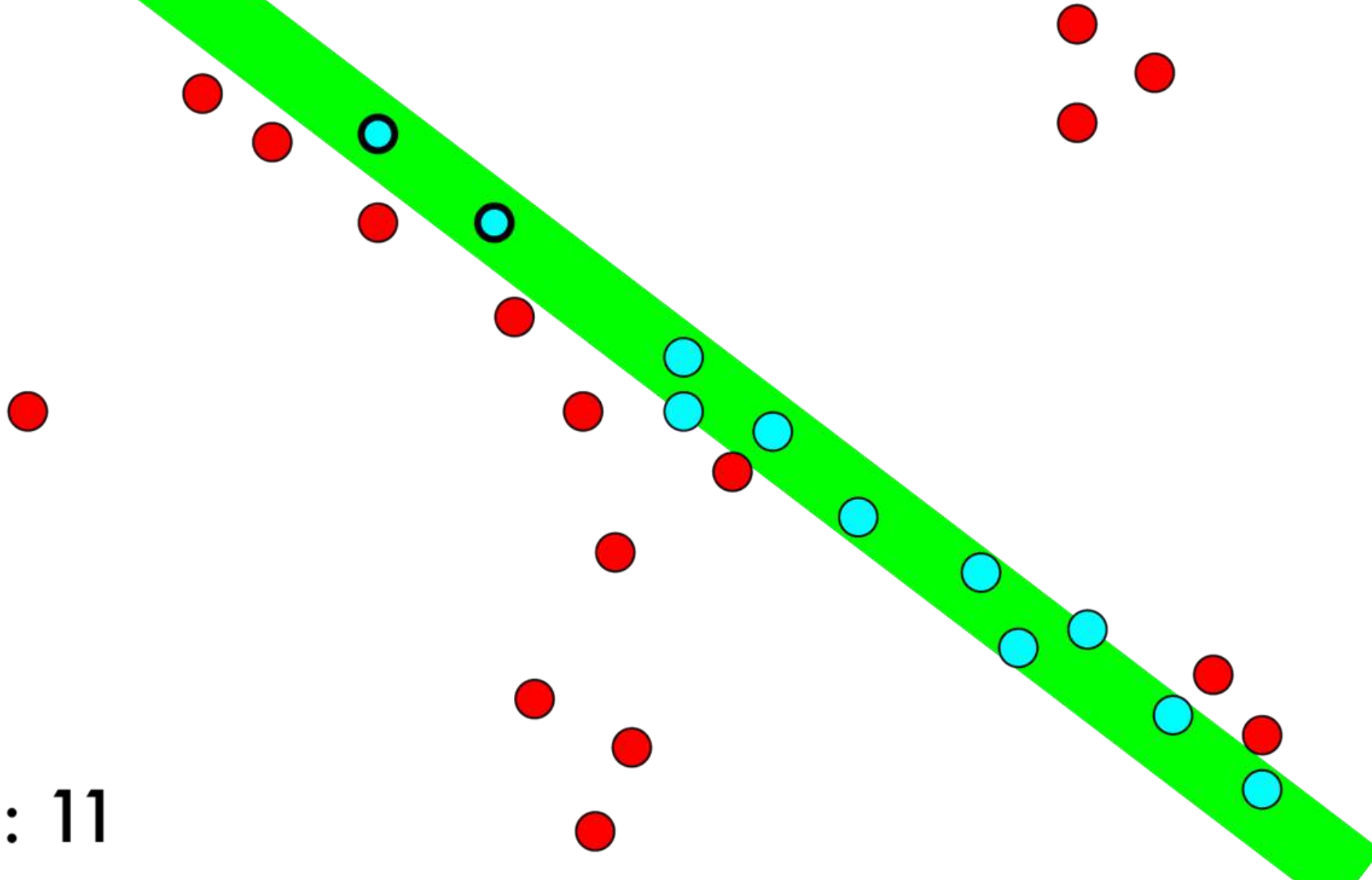




## RANSAC: RANdom SAMple Consensus



## RANSAC: Random Sample Consensus

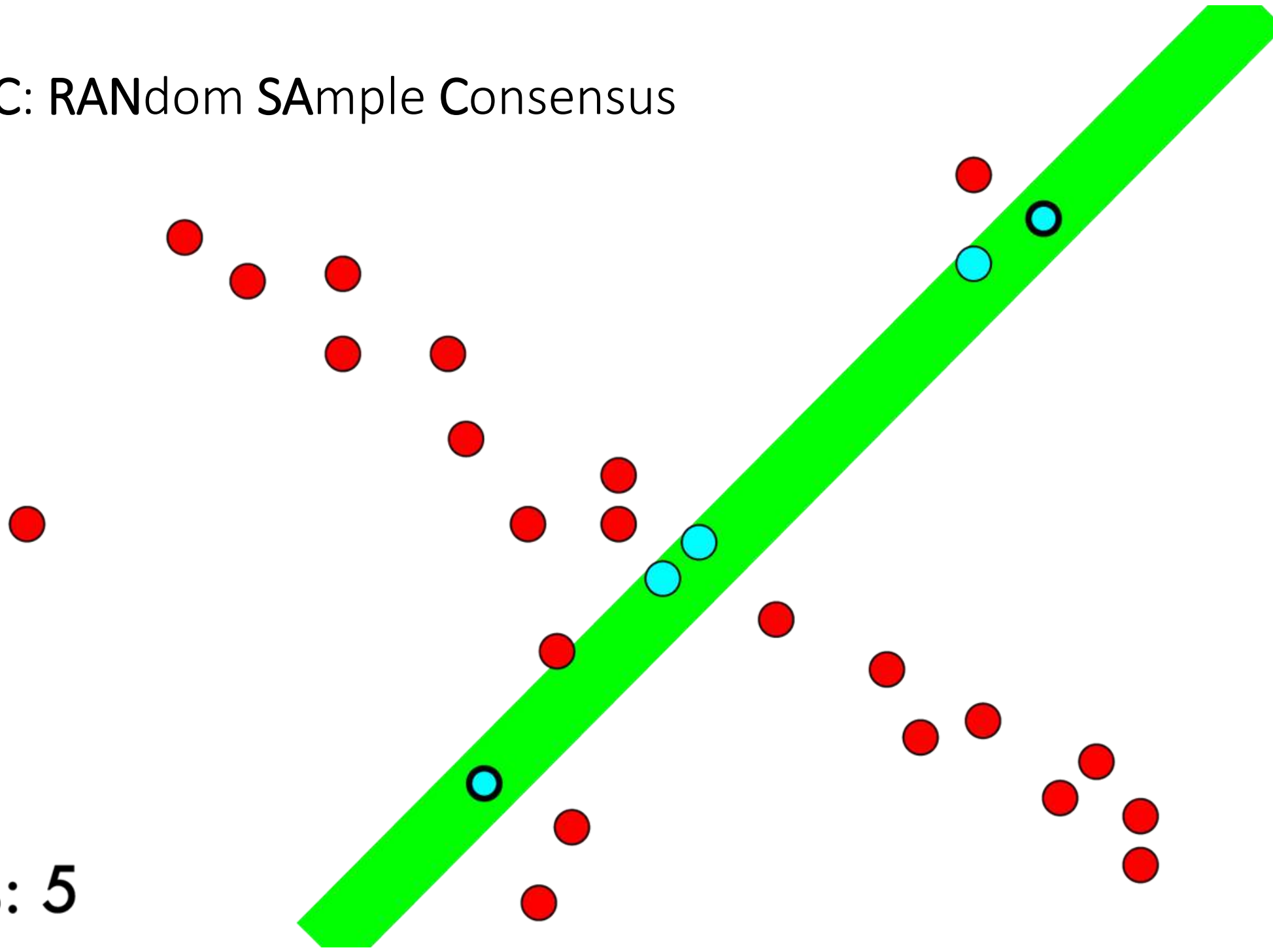


Inliers: 11





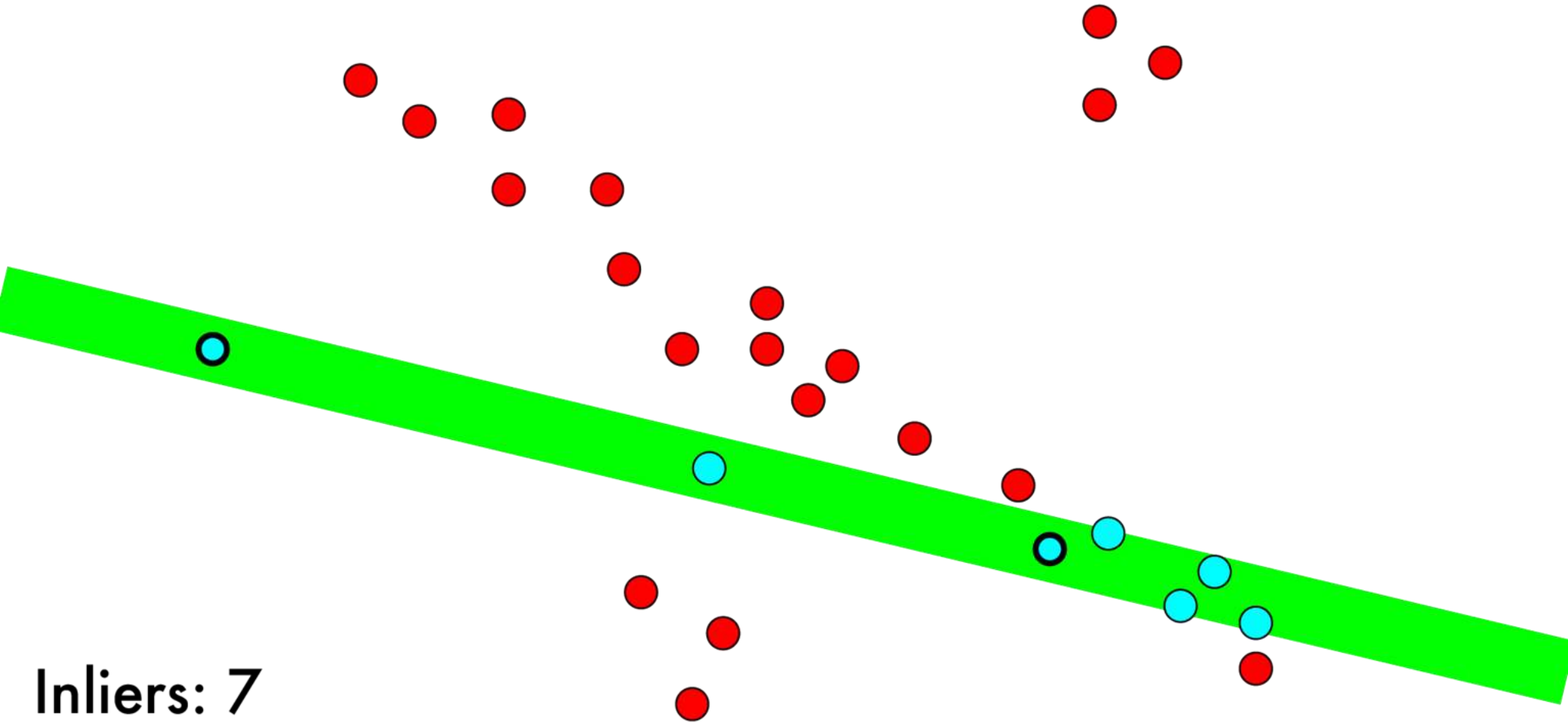
## RANSAC: RANdom SAMple Consensus



Inliers: 5

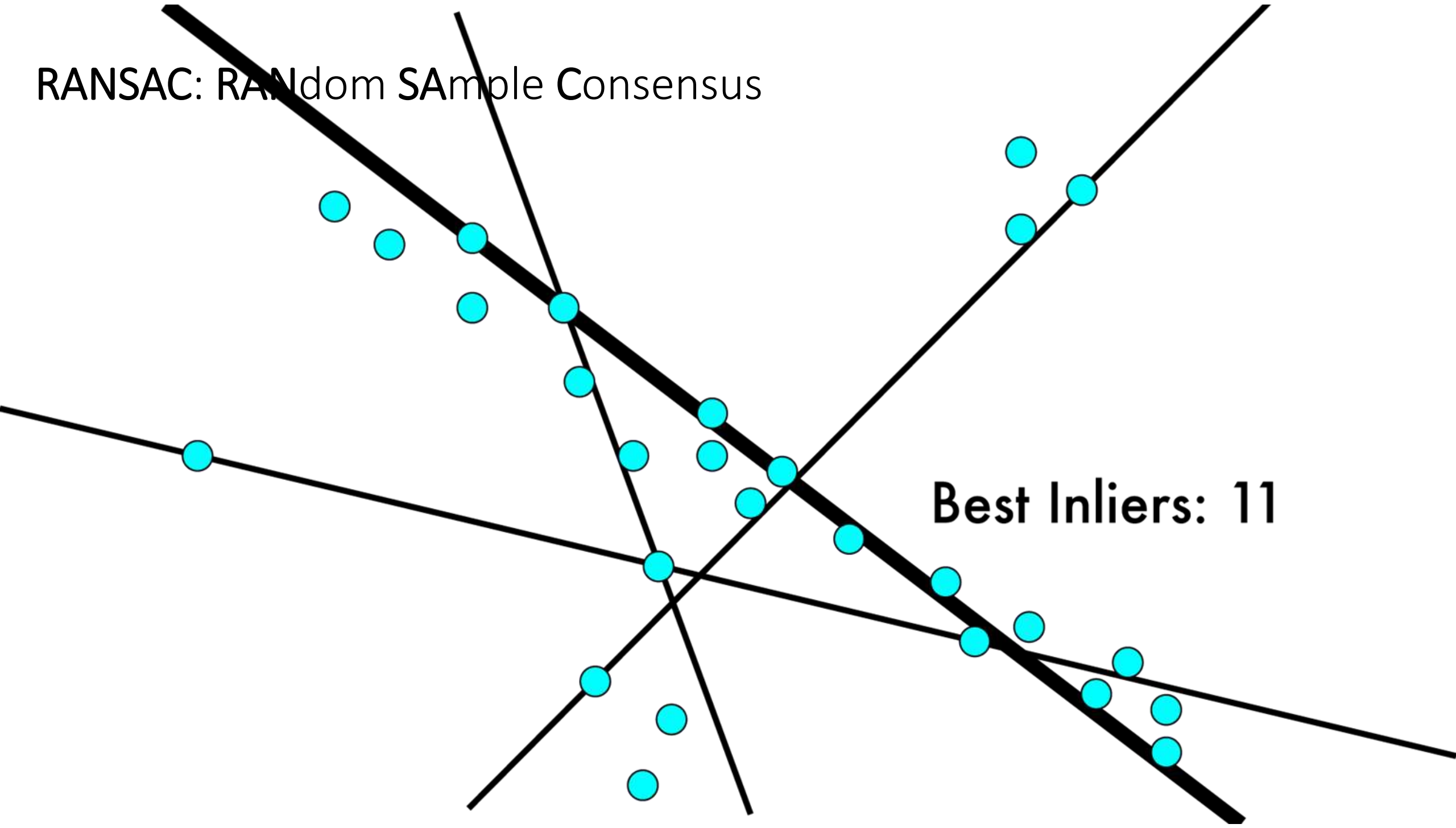


## RANSAC: RANdom SAMple Consensus



Inliers: 7

## RANSAC: RANDOM Sample Consensus



**Best Inliers: 11**

## RANSAC: RANdOm SAMple Consensus

- Parameters: data, model, n points to fit model, k iterations, t threshold, d “good” fit cutoff

```
bestmodel = None
```

```
bestfit = -INF
```

```
While i < k:
```

```
    sample = draw n random points from data
```

```
    Fit model to sample
```

```
    inliers = data within t of model
```

```
    if inliers > bestfit:
```

```
        Fit model to all inliers
```

```
        bestfit = fit
```

```
        bestmodel = model
```

```
    if inliers > d:
```

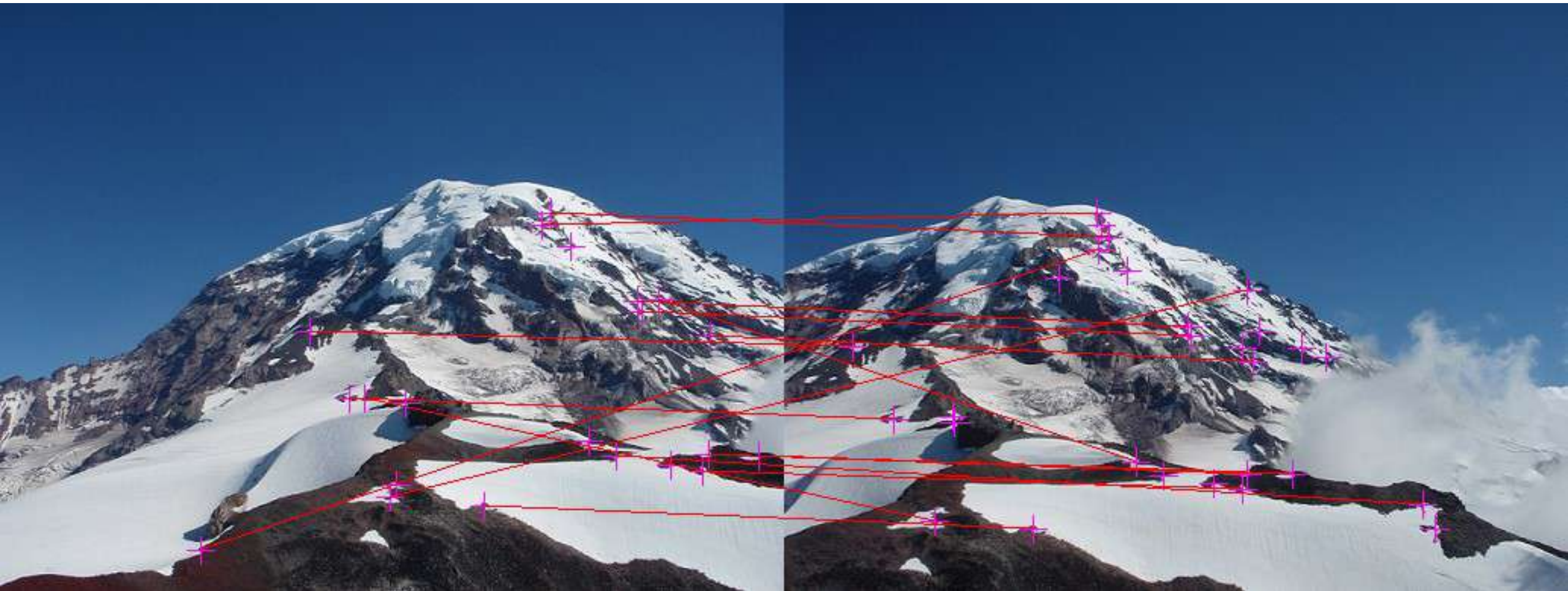
```
        return model
```

```
return bestmodel
```



# RANSAC: RANdom SAMple Consensus

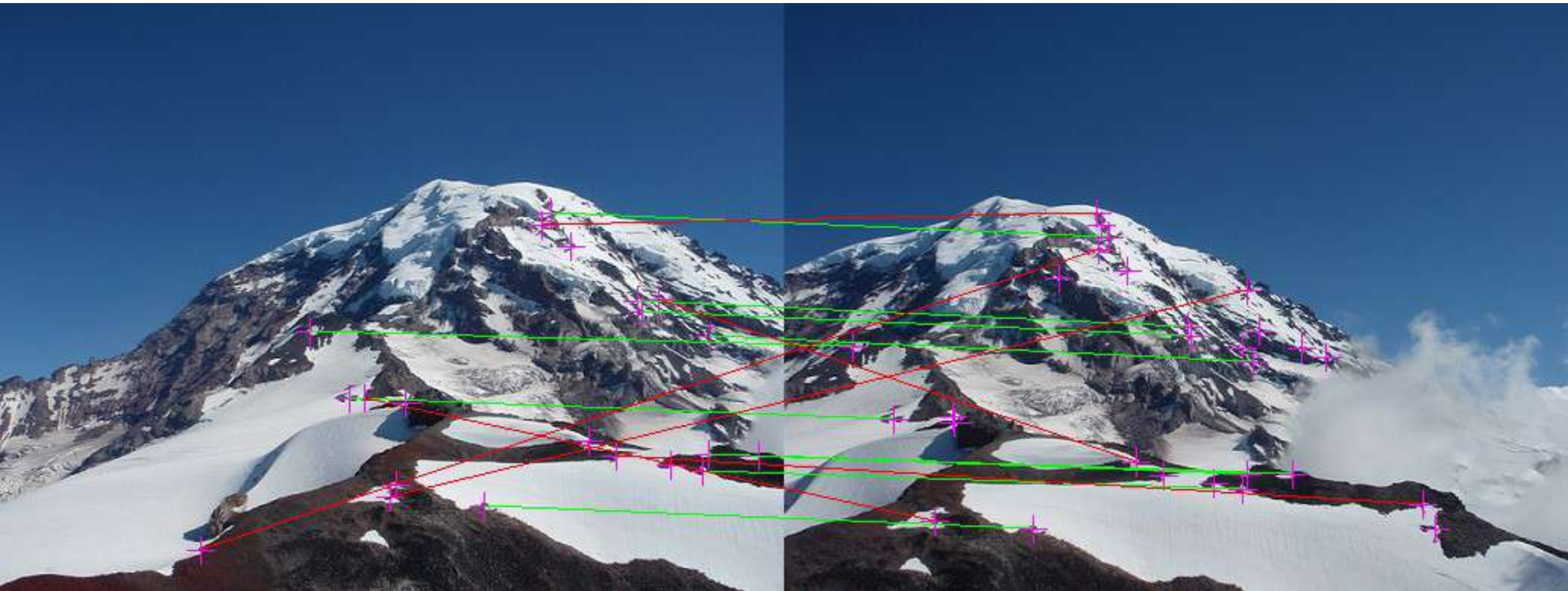
- Works well even with extreme noise.





# RANSAC: RANdom SAMple Consensus

- Works well even with extreme noise.



## RANSAC: RANdOm SAmpLe Consensus

- Parameters: data, model, n points to fit model, k iterations, t threshold, d “good” fit cutoff
- Lots of tunable parameters
- Want high probability of recovering “right” model
- t: often quite small, assume “good” inliers
- n: should be just enough to fit model, no extra
- k: can be very high
- d: should be  $\gg n$

## We can estimate affine..

- How many knowns do we get with one match?

- $n_x = a_{00} * m_x + a_{01} * m_y + a_{02} * 1$
- $n_y = a_{10} * m_x + a_{11} * m_y + a_{12} * 1$
- Solve linear system of equations  $\mathbf{M} \mathbf{a} = \mathbf{b}$ 
  - $M^{-1} M \mathbf{a} = M^{-1} \mathbf{b} \Rightarrow \mathbf{a} = M^{-1} \mathbf{b}$
  - But  $M^{-1}$  does not exist in general - Why?
- Still works if overdetermined
  - Why???
  - Pseudoinverse – least squares solution
  - $M^T M \mathbf{a} = M^T \mathbf{b}$
  - $(M^T M)^{-1} (M^T M) \mathbf{a} = (M^T M)^{-1} M^T \mathbf{b}$
  - $\Rightarrow \mathbf{a} = (M^T M)^{-1} M^T \mathbf{b}$

$$\begin{matrix} & \mathbf{M} & & \mathbf{a} & & \mathbf{b} \\ & & & & & \\ \left[ \begin{array}{cccccc} m_{x1} & m_{y1} & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & m_{x1} & m_{y1} & 1 \\ m_{x2} & m_{y2} & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & m_{x2} & m_{y2} & 1 \\ m_{x3} & m_{y3} & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & m_{x3} & m_{y3} & 1 \\ & \dots & & & & \\ & \dots & & & & \end{array} \right] & = & \left[ \begin{array}{c} a_{00} \\ a_{01} \\ a_{02} \\ a_{10} \\ a_{11} \\ a_{12} \end{array} \right] & & \left[ \begin{array}{c} n_{x1} \\ n_{y1} \\ n_{x2} \\ n_{y2} \\ n_{x3} \\ n_{y3} \\ \dots \\ \dots \end{array} \right]
 \end{matrix}$$



# We want projective (homography)

- What are our equations now?

- $n_x = (h_{00} * m_x + h_{01} * m_y + h_{02} * m_w) / (h_{20} * m_x + h_{21} * m_y + h_{22} * m_w)$
- $n_y = (h_{10} * m_x + h_{11} * m_y + h_{12} * m_w) / (h_{20} * m_x + h_{21} * m_y + h_{22} * m_w)$

$$\begin{bmatrix} \tilde{x}' \\ \tilde{y}' \\ \tilde{w}' \end{bmatrix} = \begin{bmatrix} h_{00} & h_{01} & h_{02} \\ h_{10} & h_{11} & h_{12} \\ h_{20} & h_{21} & h_{22} \end{bmatrix} \begin{bmatrix} \tilde{x} \\ \tilde{y} \\ \tilde{w} \end{bmatrix}$$

# We want projective (homography)

- What are our equations now?
  - $n_x = (h_{00} * m_x + h_{01} * m_y + h_{02} * m_w) / (h_{20} * m_x + h_{21} * m_y + h_{22} * m_w)$
  - $n_y = (h_{10} * m_x + h_{11} * m_y + h_{12} * m_w) / (h_{20} * m_x + h_{21} * m_y + h_{22} * m_w)$
- Assume  $h_{22}$  and  $m_w$  are 1, now 8 DOF
  - $n_x = (h_{00} * m_x + h_{01} * m_y + h_{02}) / (h_{20} * m_x + h_{21} * m_y + 1)$
  - $n_y = (h_{10} * m_x + h_{11} * m_y + h_{12}) / (h_{20} * m_x + h_{21} * m_y + 1)$

$$\begin{bmatrix} \tilde{x}' \\ \tilde{y}' \\ \tilde{w}' \end{bmatrix} = \begin{bmatrix} h_{00} & h_{01} & h_{02} \\ h_{10} & h_{11} & h_{12} \\ h_{20} & h_{21} & h_{22} \end{bmatrix} \begin{bmatrix} \tilde{x} \\ \tilde{y} \\ \tilde{w} \end{bmatrix}$$

# We want projective (homography)

- What are our equations now?
  - $n_x = (h_{00} * m_x + h_{01} * m_y + h_{02} * m_w) / (h_{20} * m_x + h_{21} * m_y + h_{22} * m_w)$
  - $n_y = (h_{10} * m_x + h_{11} * m_y + h_{12} * m_w) / (h_{20} * m_x + h_{21} * m_y + h_{22} * m_w)$
- Assume  $h_{22}$  and  $m_w$  are 1, now 8 DOF
  - $n_x = (h_{00} * m_x + h_{01} * m_y + h_{02}) / (h_{20} * m_x + h_{21} * m_y + 1)$
  - $n_y = (h_{10} * m_x + h_{11} * m_y + h_{12}) / (h_{20} * m_x + h_{21} * m_y + 1)$
- More algebra on  $n_x$ 
  - $n_x * (h_{20} * m_x + h_{21} * m_y + 1) = (h_{00} * m_x + h_{01} * m_y + h_{02})$
  - $n_x * h_{20} * m_x + n_x * h_{21} * m_y + n_x = h_{00} * m_x + h_{01} * m_y + h_{02}$
  - $n_x = h_{00} * m_x + h_{01} * m_y + h_{02} - n_x * h_{20} * m_x - n_x * h_{21} * m_y$
- Similar for  $n_y$









# Today's Agenda

- Linear least-squares
- RANSAC
- Panorama Stitching



# Panorama algorithm

Find corners in both images

Calculate descriptors

Match descriptors

RANSAC to find homography

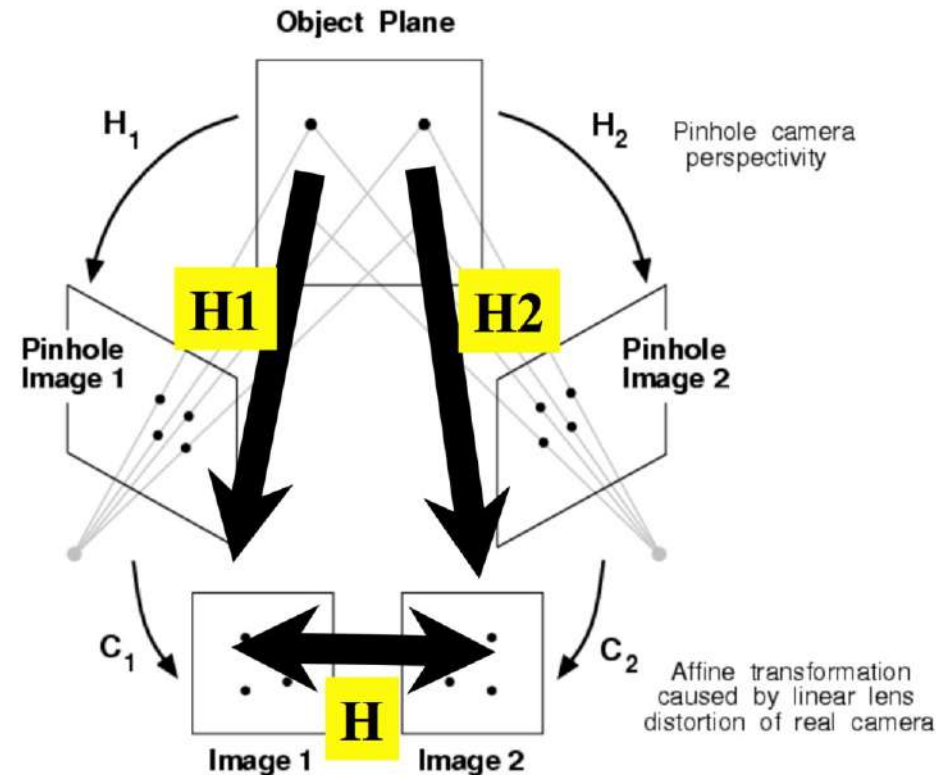
Stitch together images with homography





# Stitching panoramas

- We know homography is right choice under certain assumption:
  - Assume we are taking multiple images of planar object



## In practice



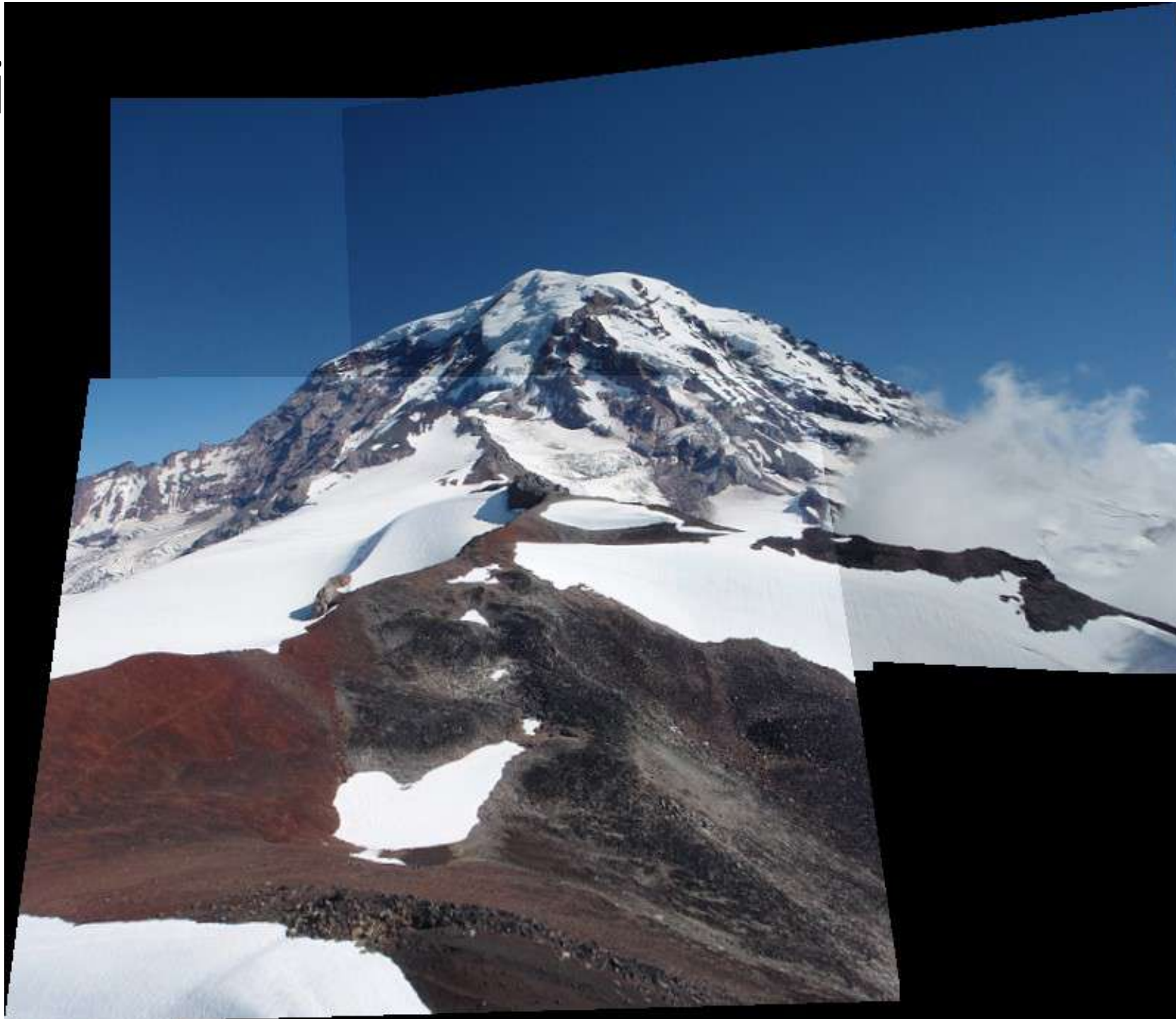
## In practice







## In practice





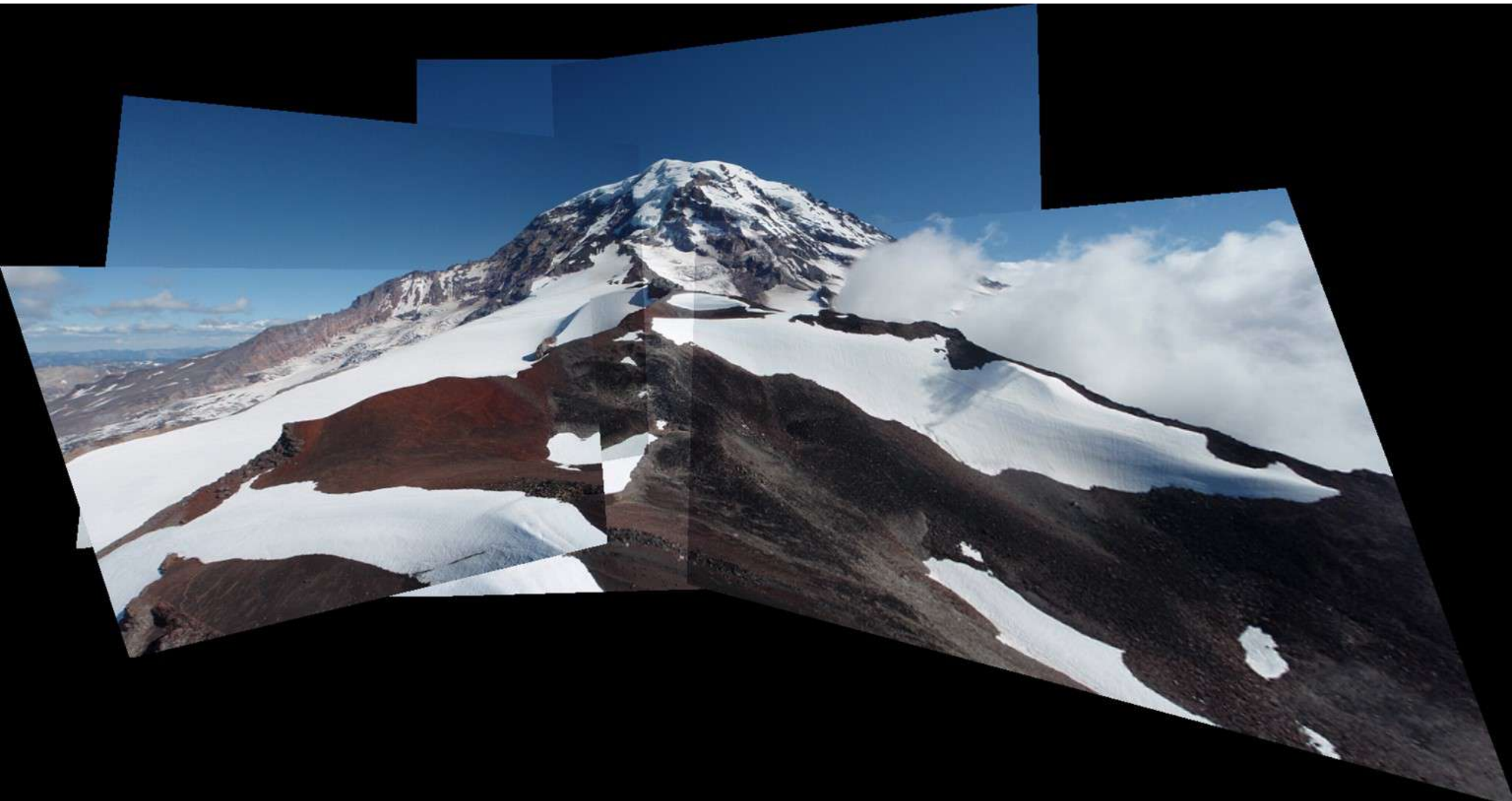


In p



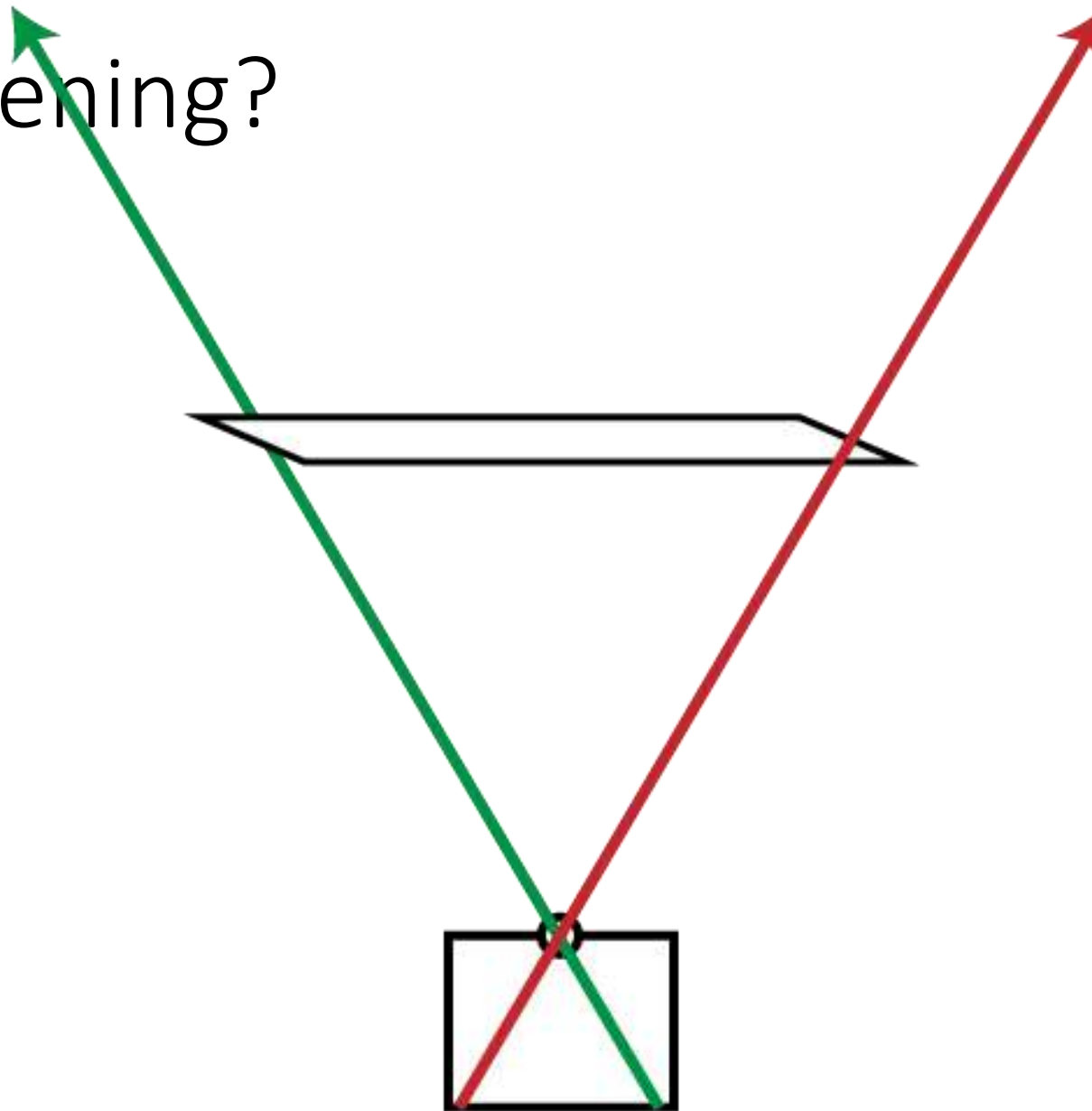








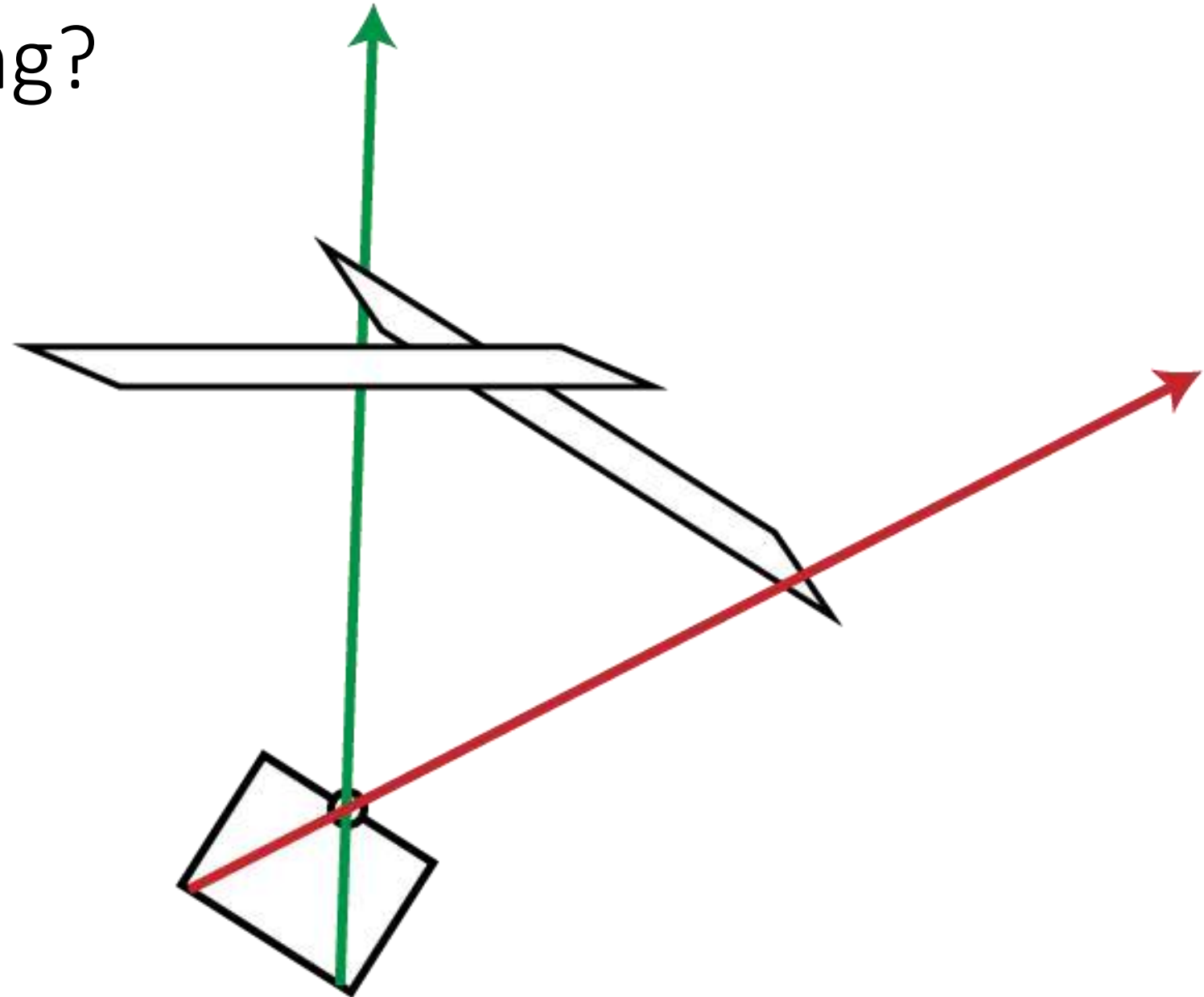
# What's happening?



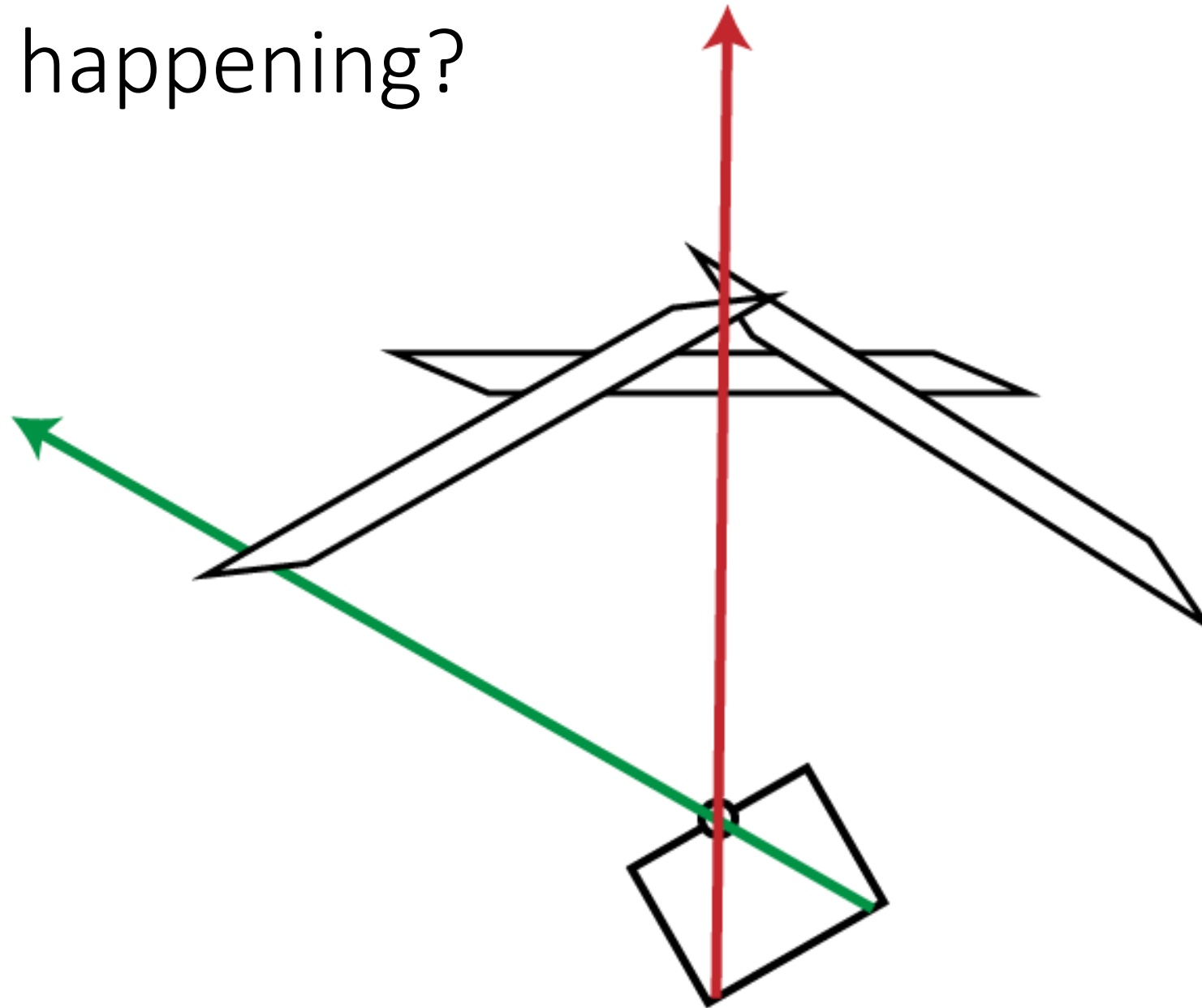




## What's happening?

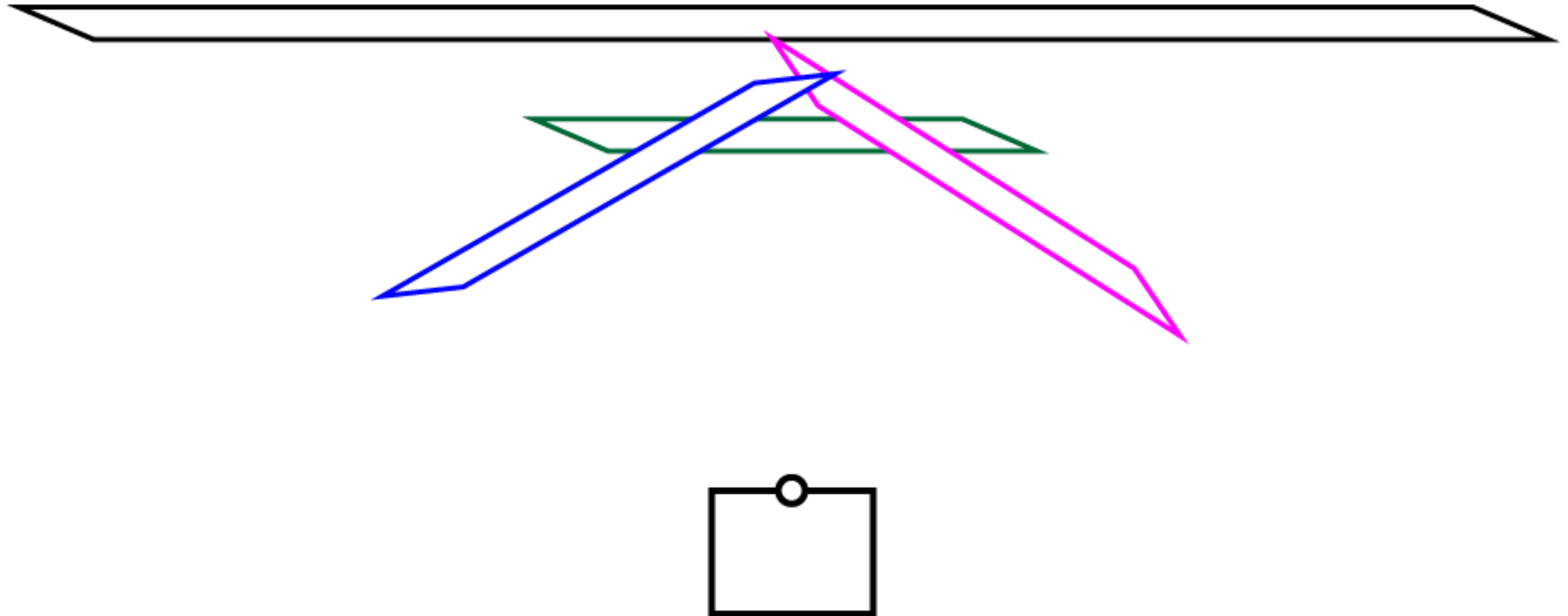


## What's happening?

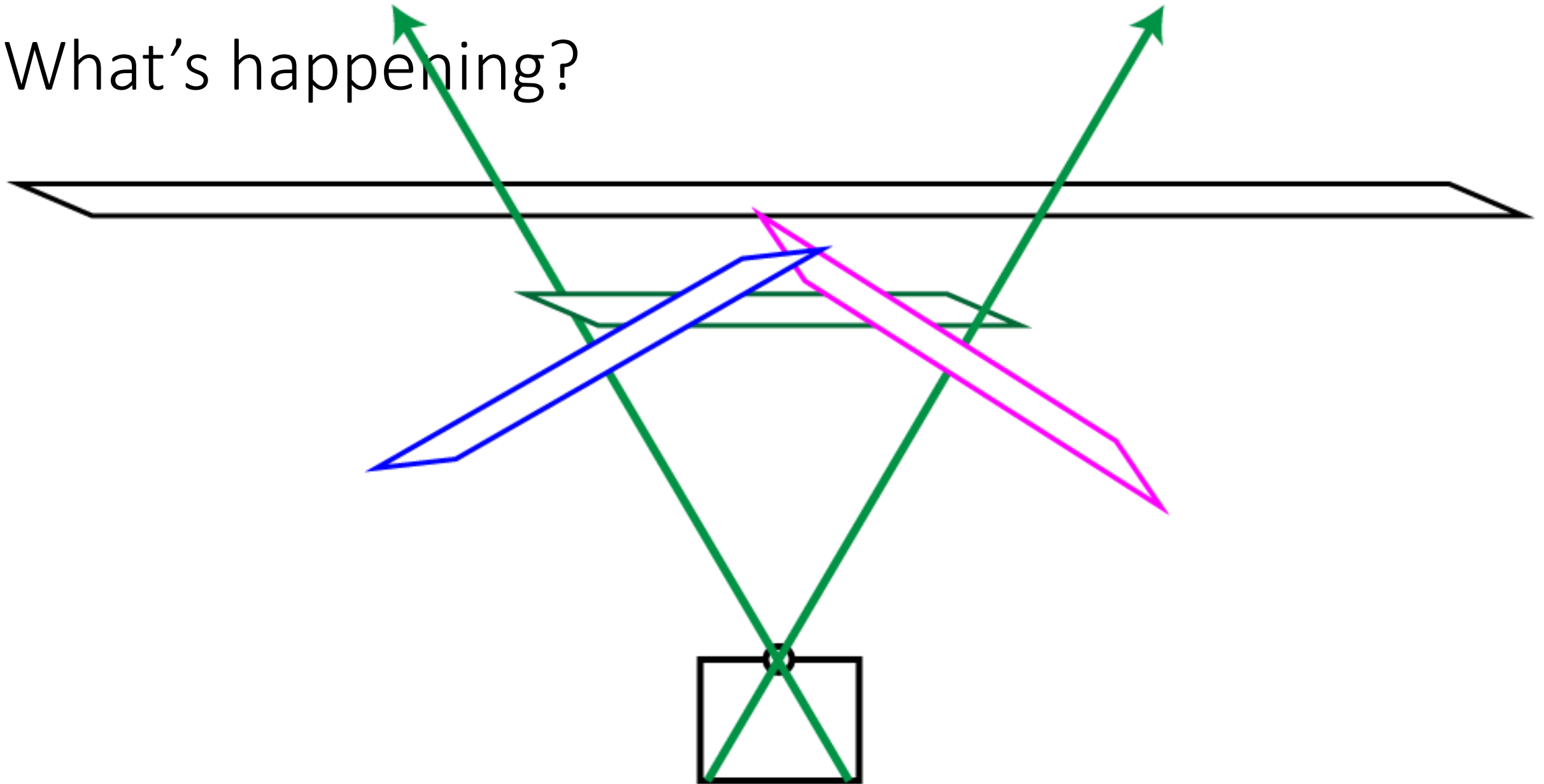




# What's happening?

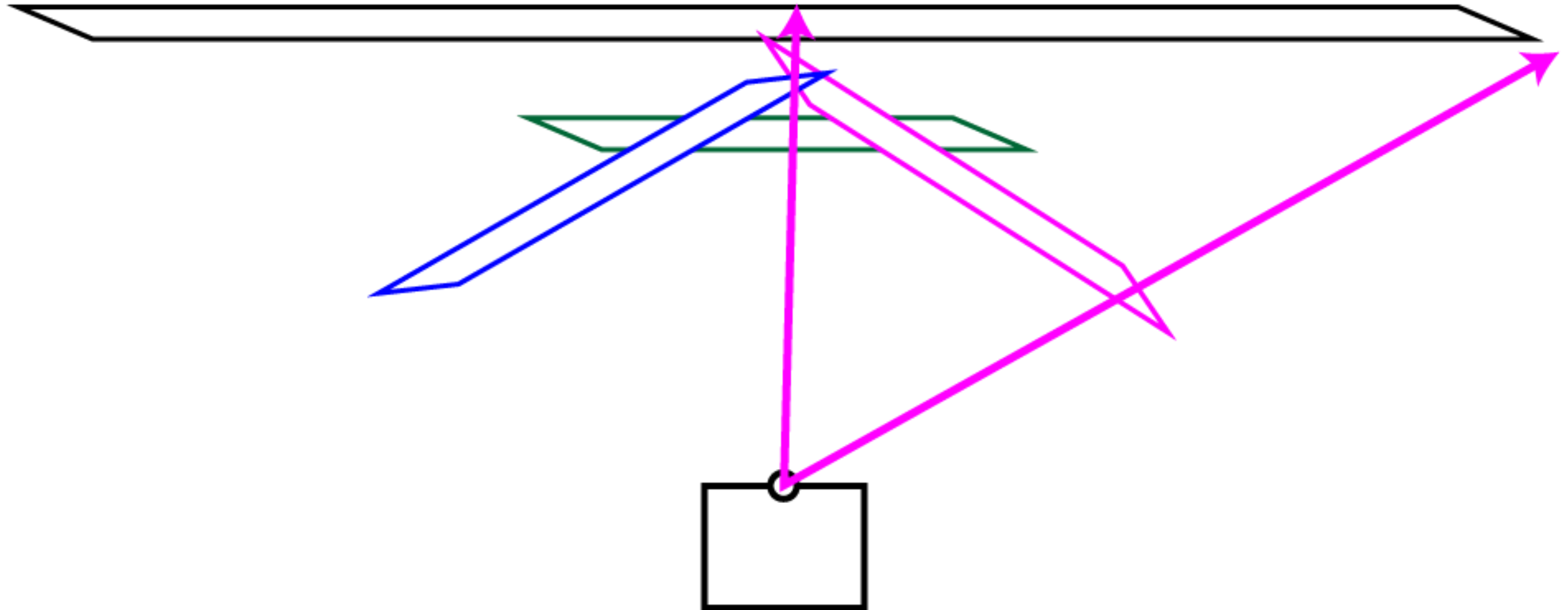


## What's happening?

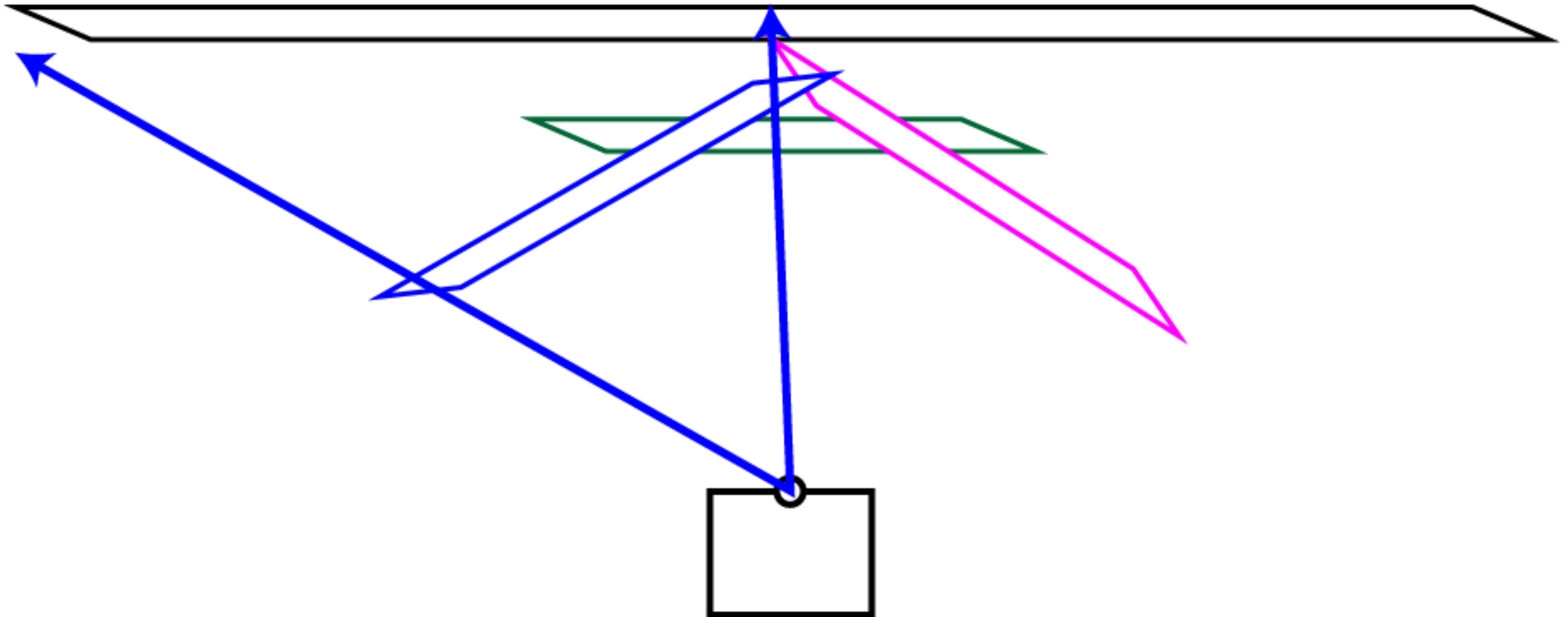




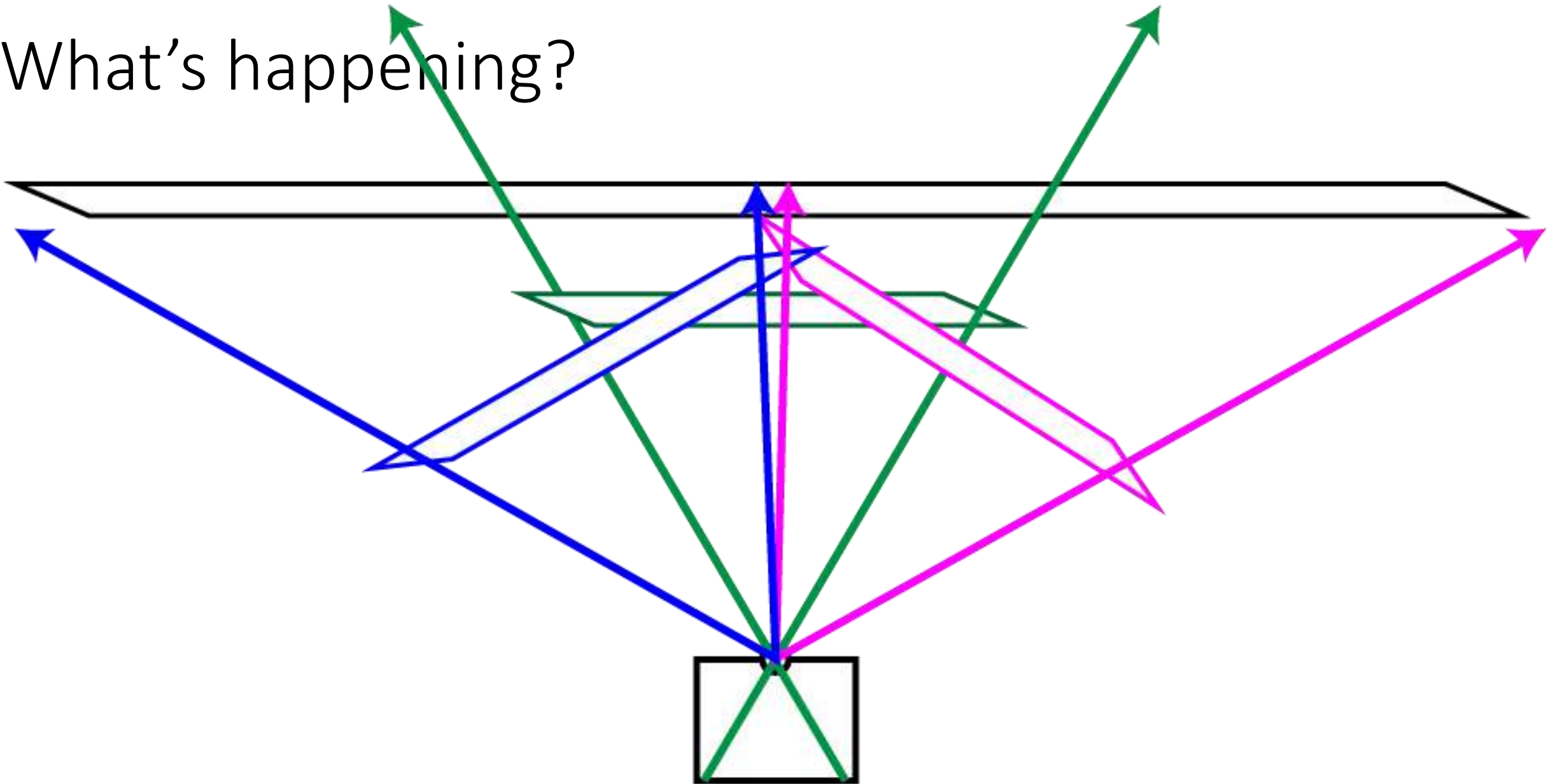
## What's happening?



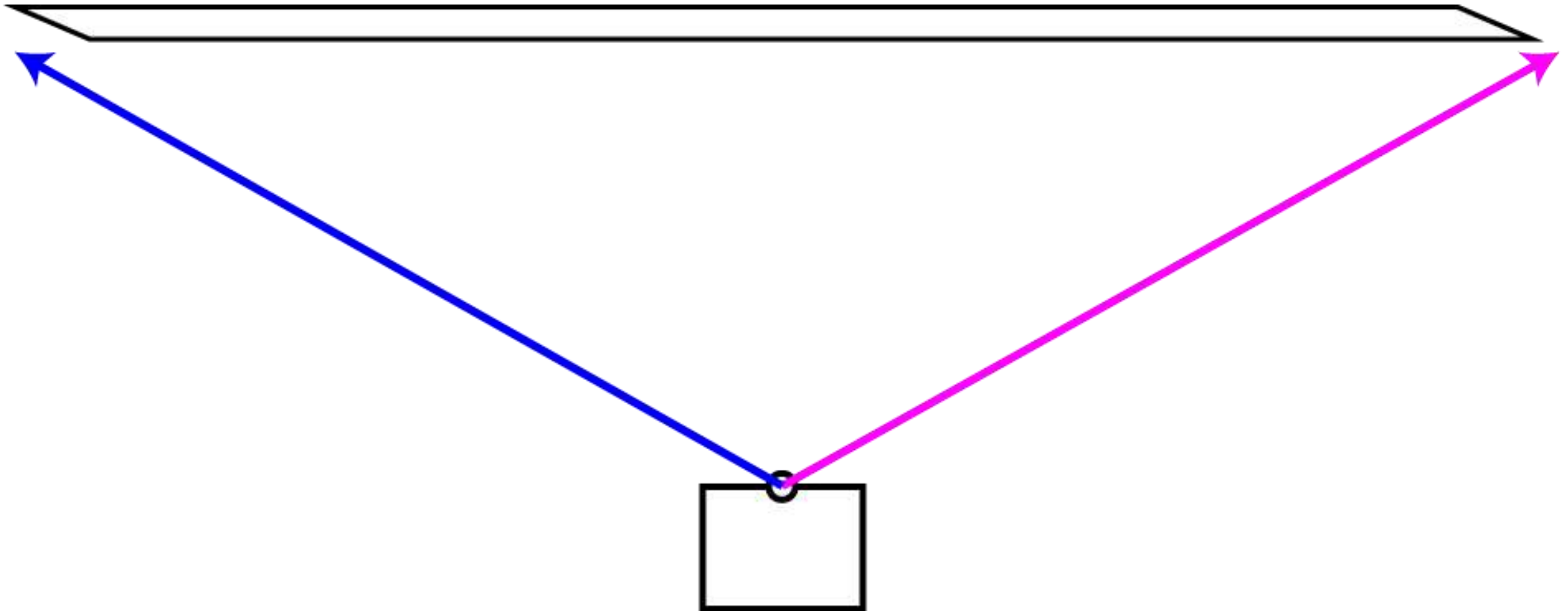
# What's happening?



## What's happening?



# What's happening?





## Very bad for big panoramas!







## Very bad for big panoramas!





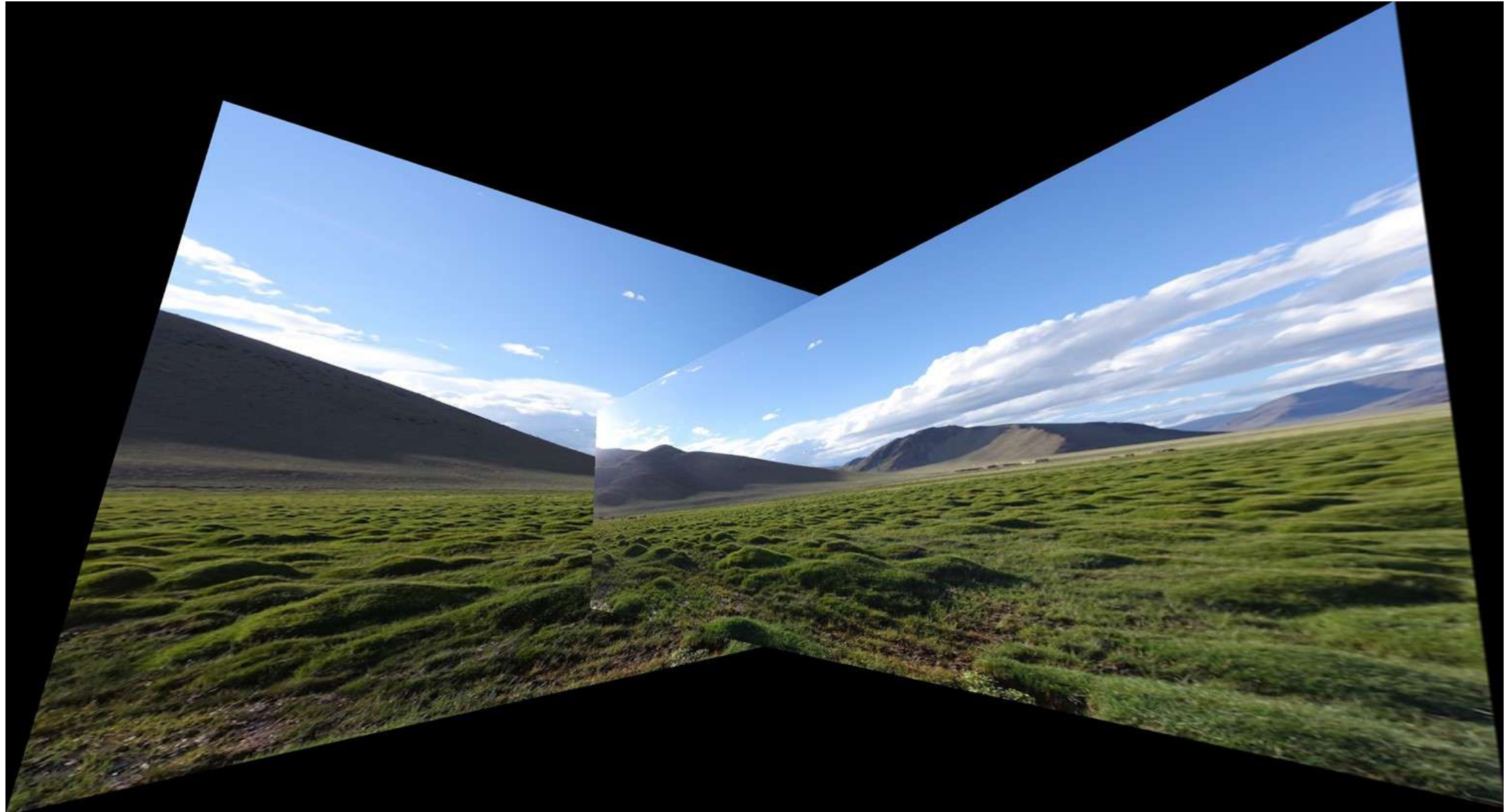
## Very bad for big panoramas!







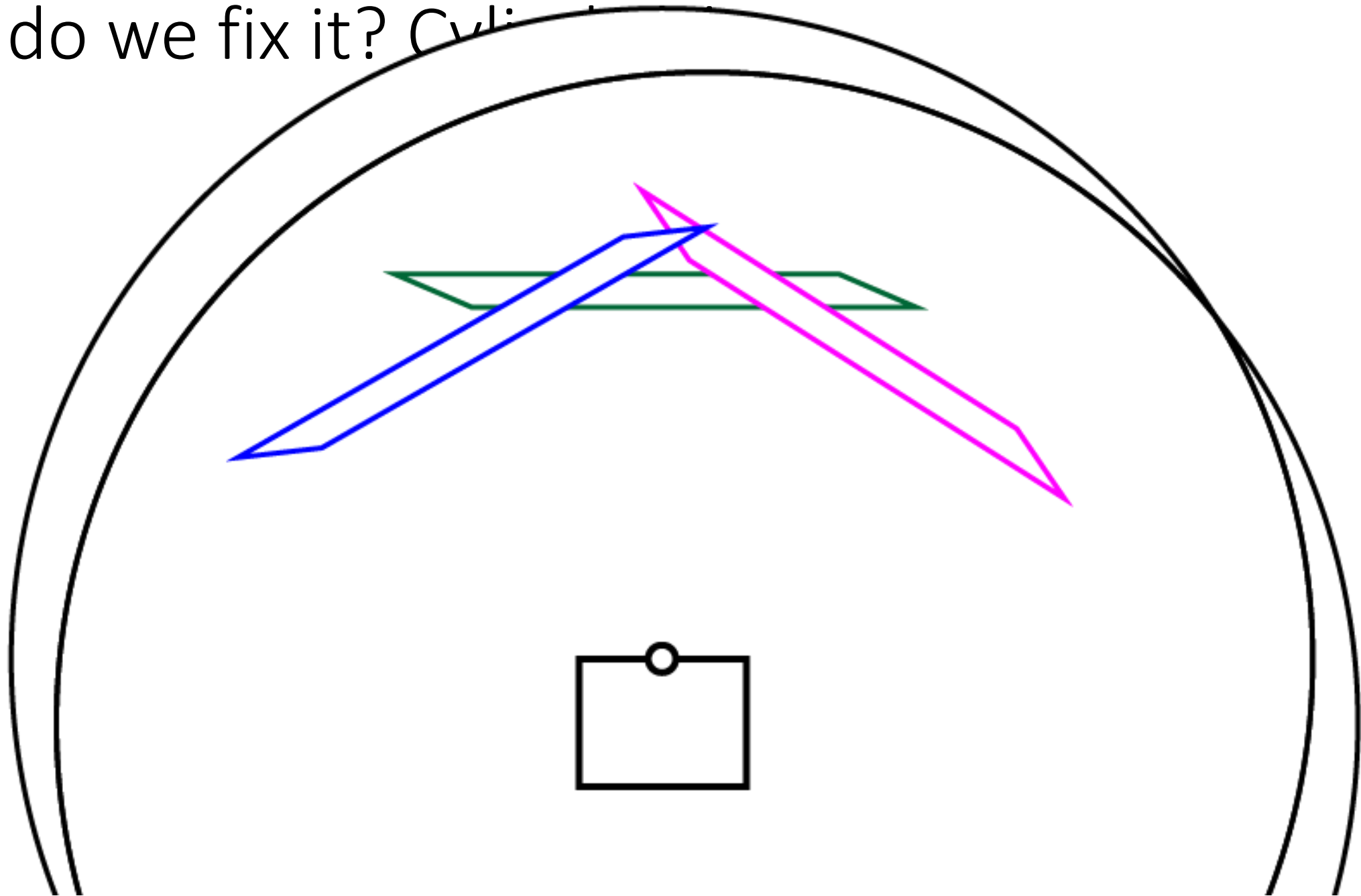
## Fails :-)



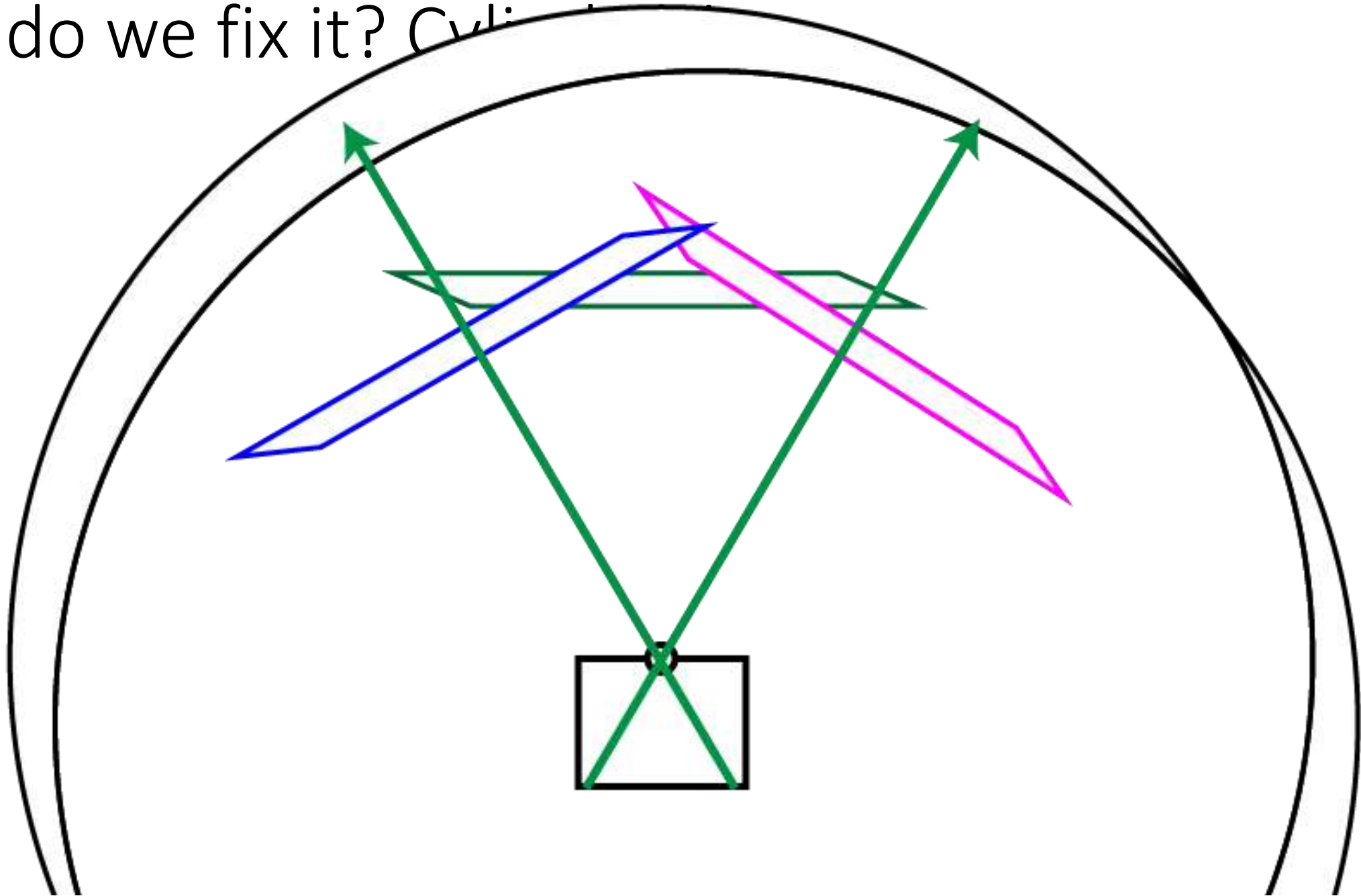


# How do we fix it? Cylinders!

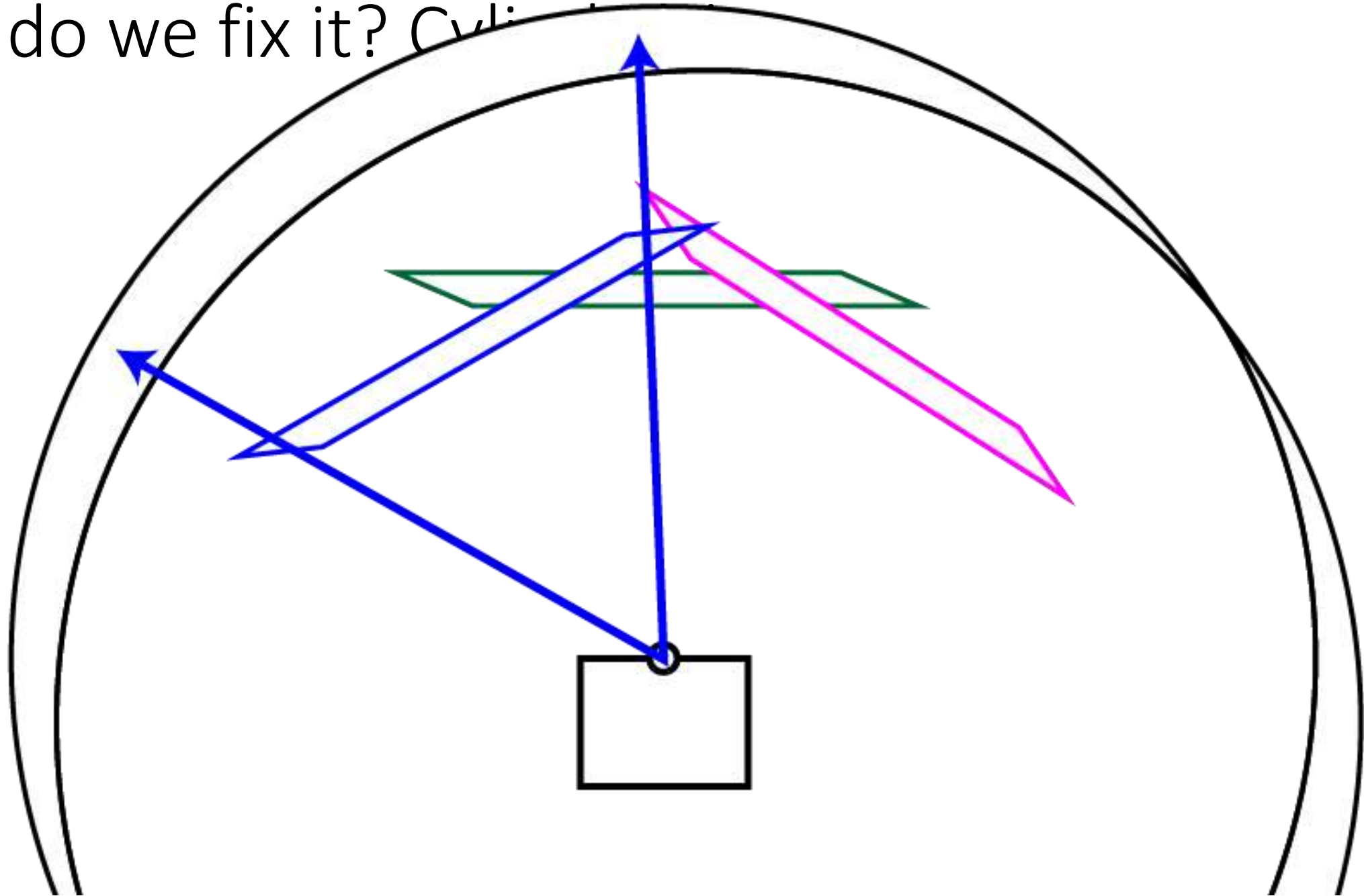
## How do we fix it? Cycle



# How do we fix it? Cycle

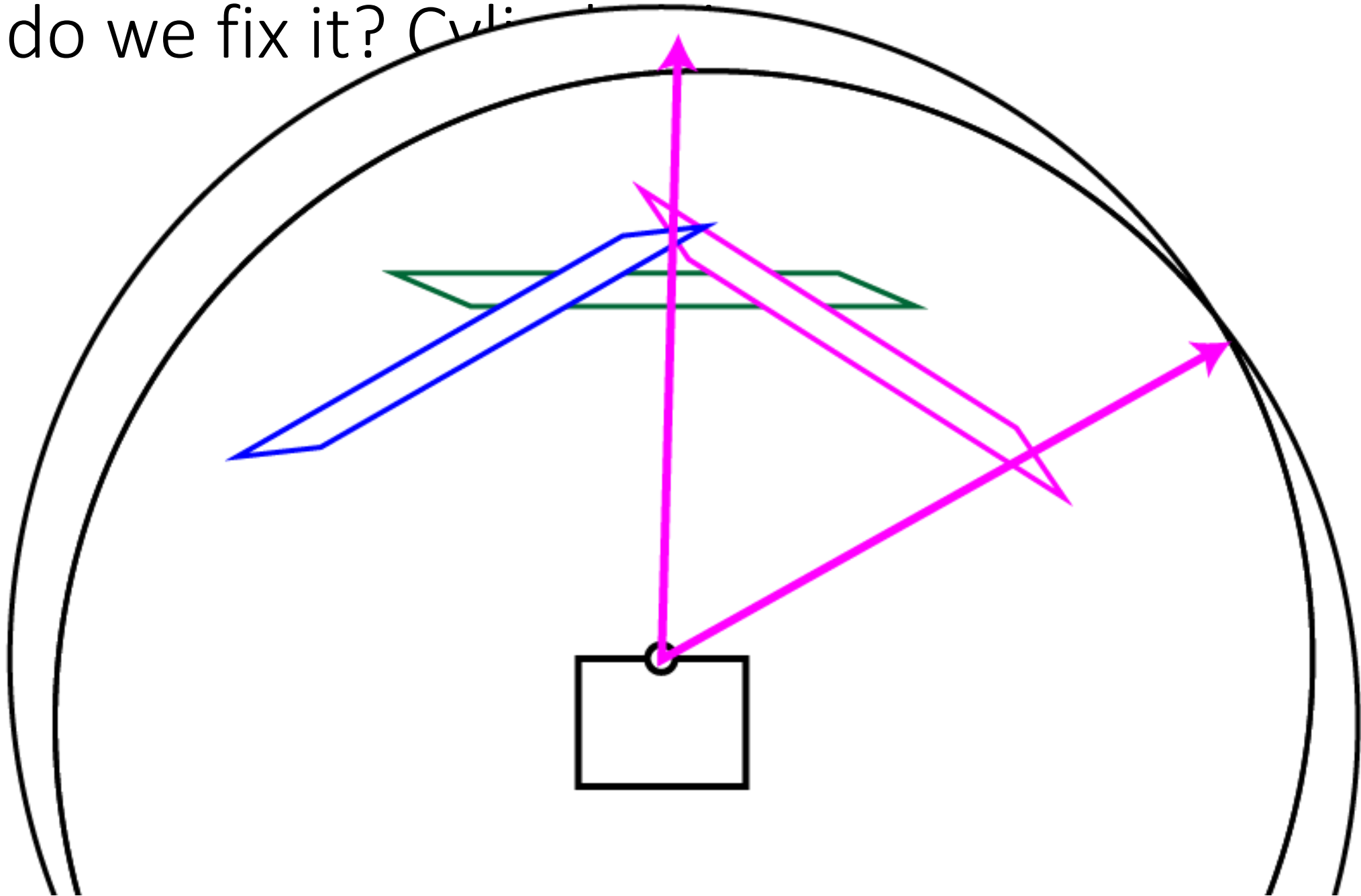


# How do we fix it? Cyclic

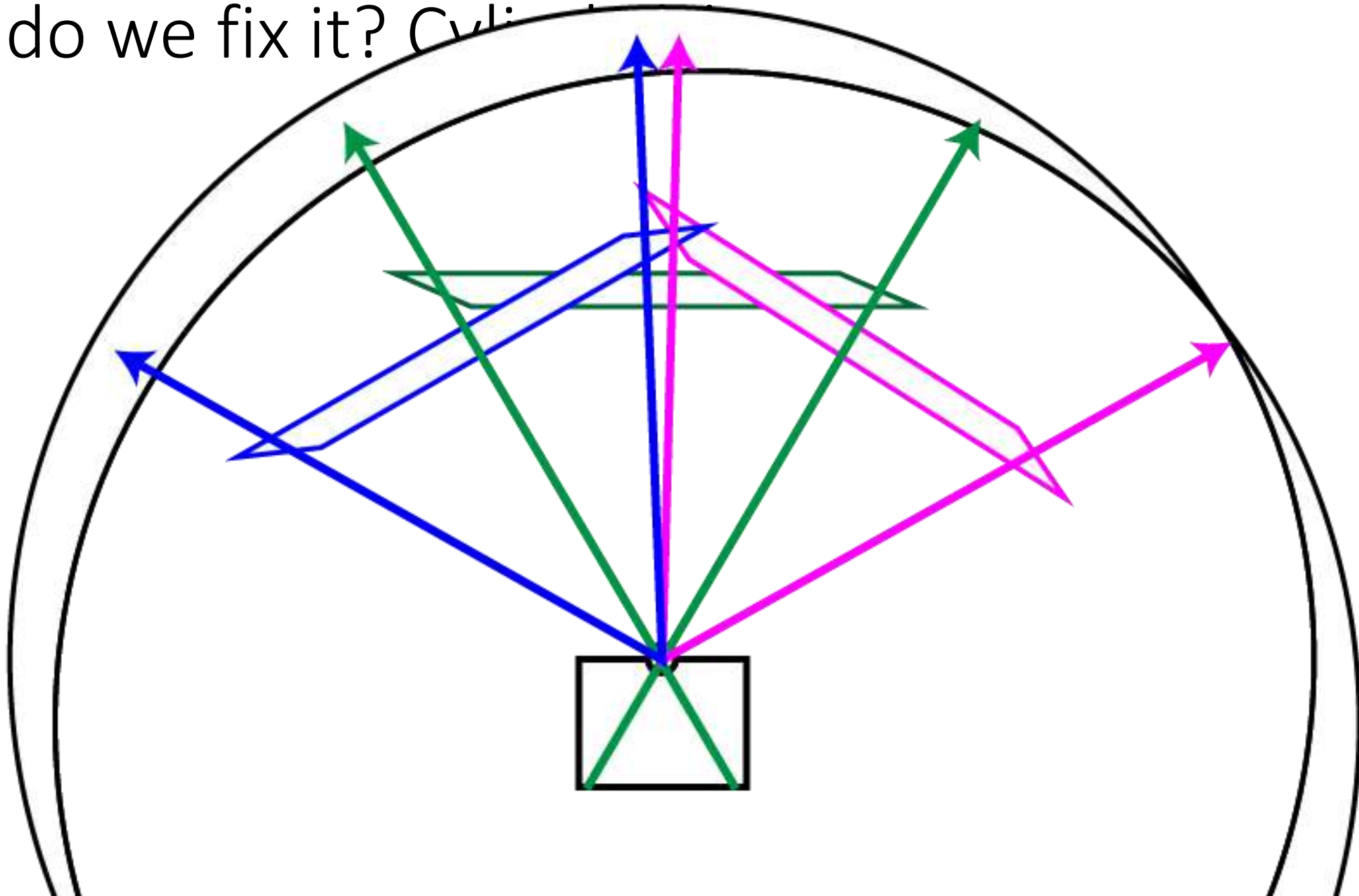




# How do we fix it? Cyl:



## How do we fix it? Cycle



# How do we fix it? Cylinders!

Calculate angle and height:

$$\theta = (x - x_c) / f$$

$$h = (y - y_c) / f$$

Find unit cylindrical coords:

$$X' = \sin(\theta)$$

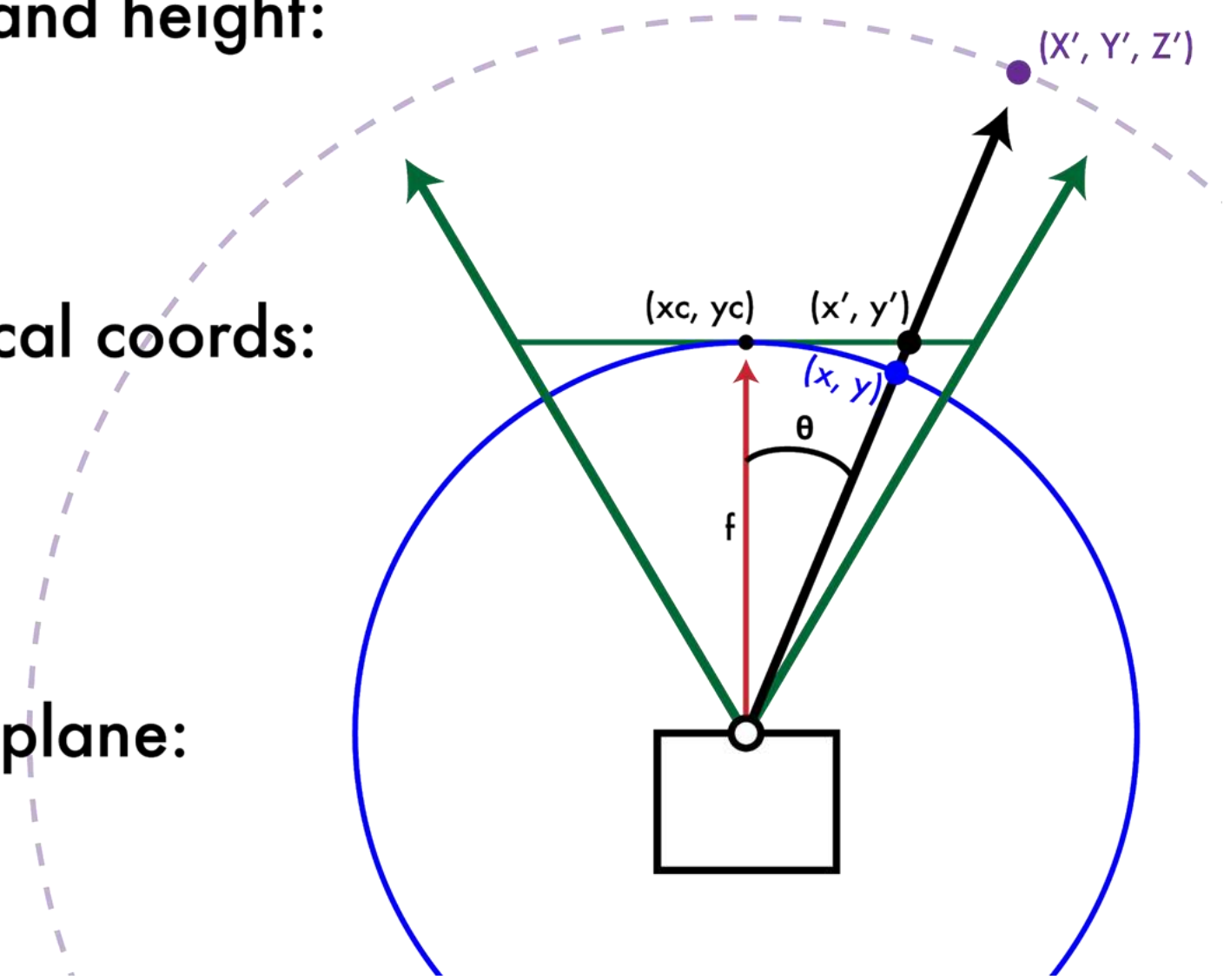
$$Y' = h$$

$$Z' = \cos(\theta)$$

Project to image plane:

$$x' = f X' / Z' + x_c$$

$$y' = f Y' / Z' + y_c$$





## Dependant on focal length!

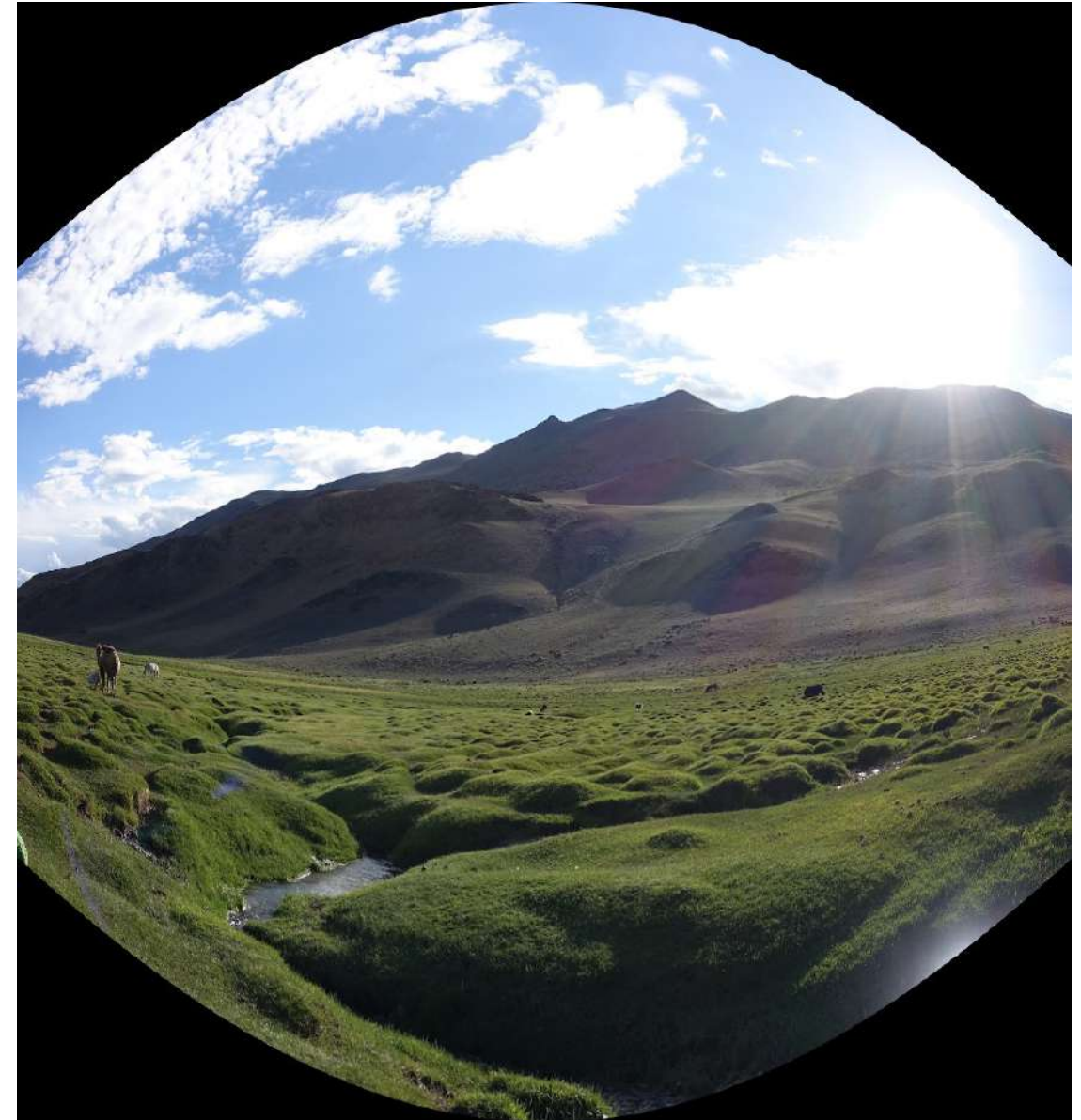




$f = 300$



$f = 500$



$$f = 1000$$





$$f = 1400$$







$$f = 10,000$$



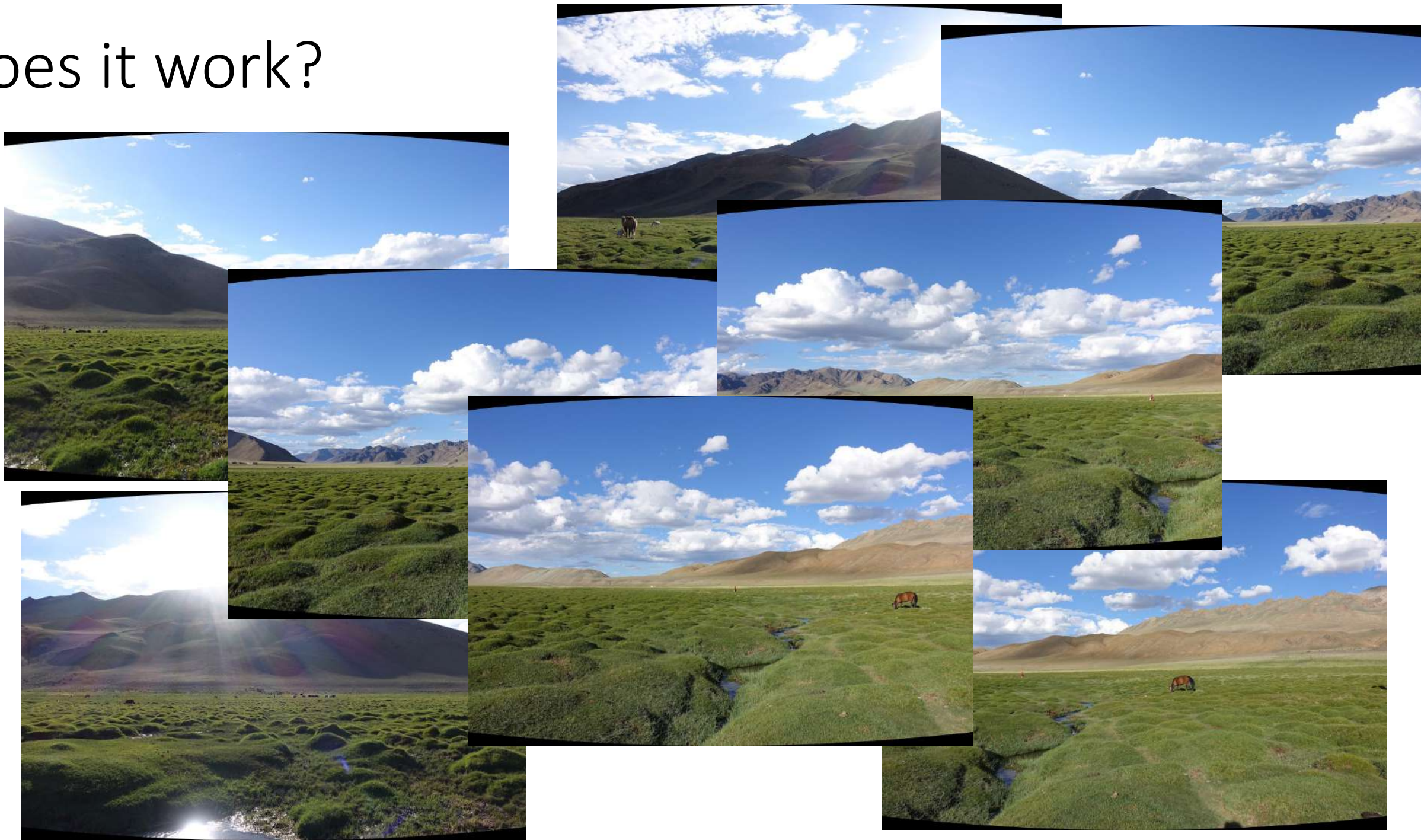


$$f = 10,000$$





## Does it work?







## Does it work?





## Does it work?



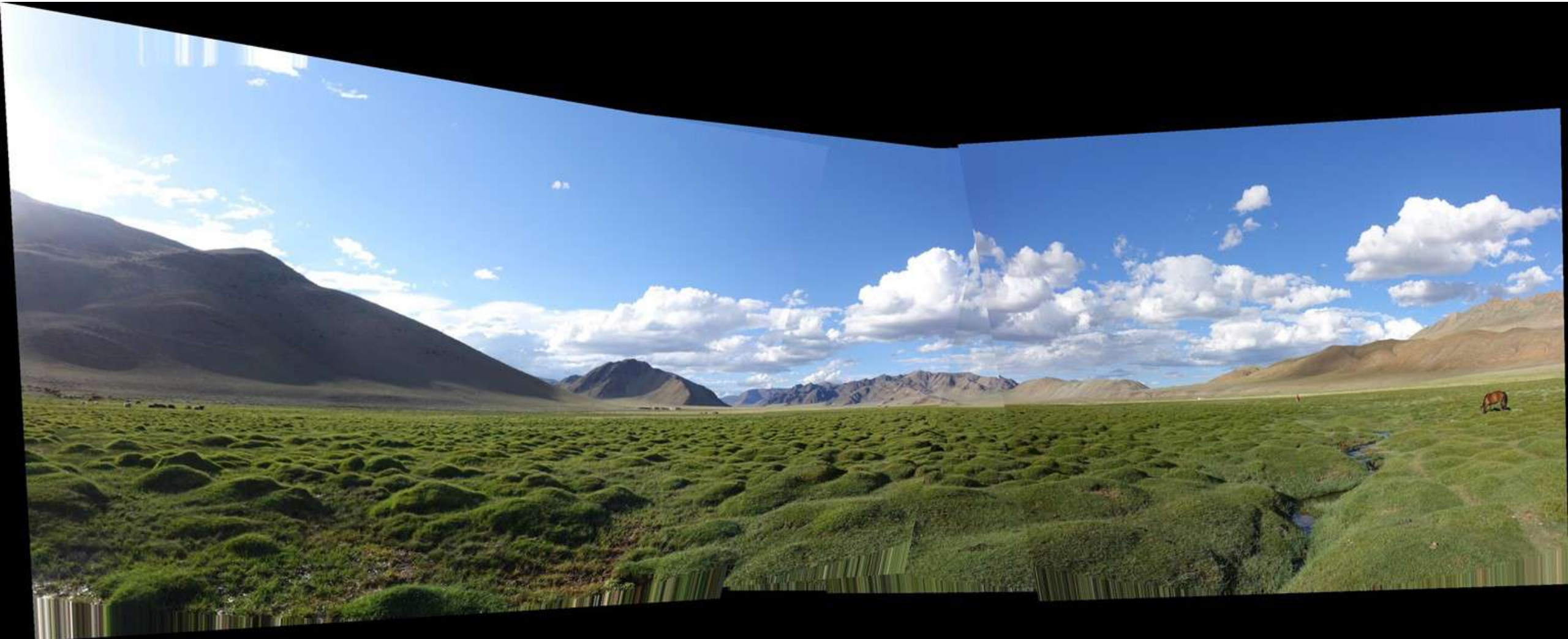


## Does it work?





## Does it work?







Yes! Assuming camera is level and rotating around its vertical axis



**MAI4CAREU**

Master programmes in Artificial  
Intelligence 4 Careers in Europe



# Thank you.







University of Cyprus – MSc Artificial Intelligence

# MAI644 – COMPUTER VISION

## Lecture 10: Visual Recognition – Segmentation

**Melinos Averkiou**

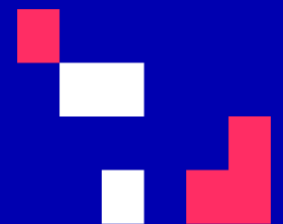
CYENS Centre of Excellence

University of Cyprus - Department of Computer Science

[m.averkiou@cyens.org.cy](mailto:m.averkiou@cyens.org.cy)



**CYENS**  
CENTRE OF EXCELLENCE



## Last time

- Linear least-squares
- RANSAC
- Panorama Stitching

# Today's Agenda

- Visual Recognition Tasks
- Introduction to segmentation and clustering
- Agglomerative clustering
- K-means clustering
- Mean-shift clustering
- Efficient Graph-based image segmentation

**Reading:** Forsyth Chapter 9

D. Comaniciu and P. Meer, [Mean Shift: A Robust Approach toward Feature Space Analysis](#), TPAMI 2002

[material based on Niebles-Krishna]

# Today's Agenda

- Visual Recognition Tasks
- Introduction to segmentation and clustering
- Agglomerative clustering
- K-means clustering
- Mean-shift clustering
- Efficient Graph-based image segmentation

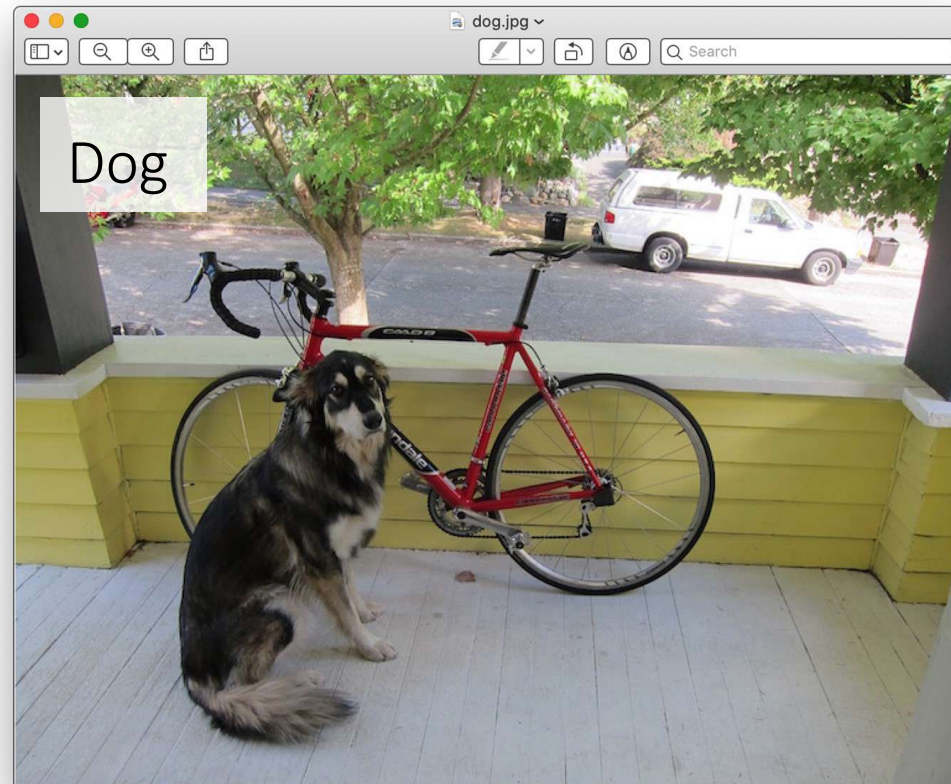
**Reading:** Forsyth Chapter 9

D. Comaniciu and P. Meer, [Mean Shift: A Robust Approach toward Feature Space Analysis](#), TPAMI 2002



# Classification

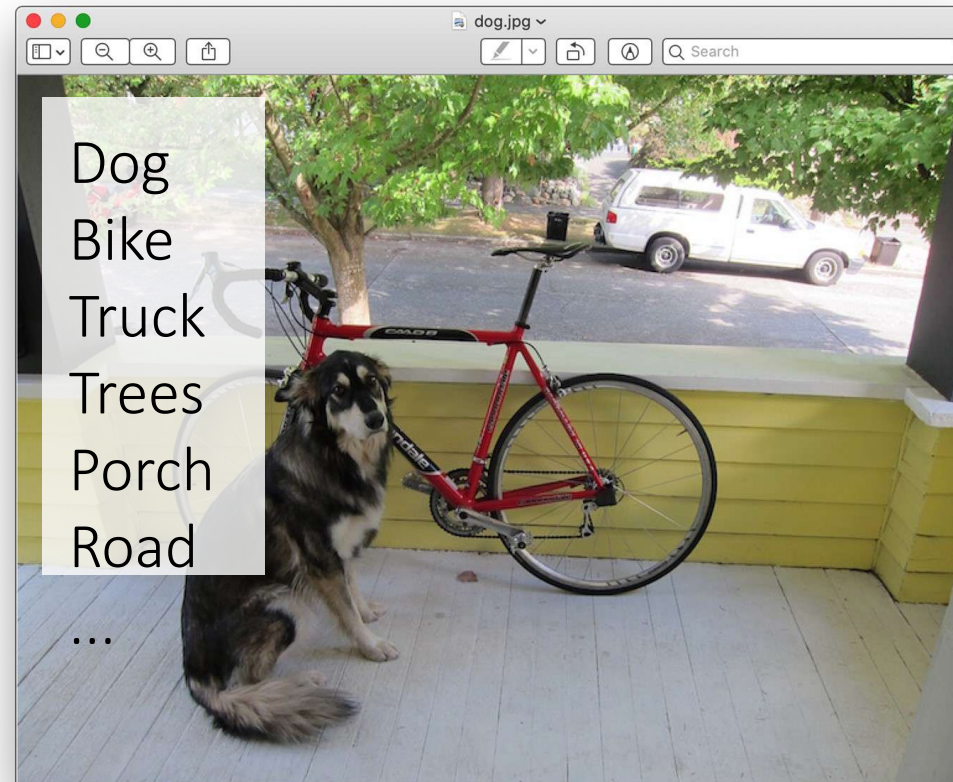
- What is in the image?



[Redmon]

# Tagging

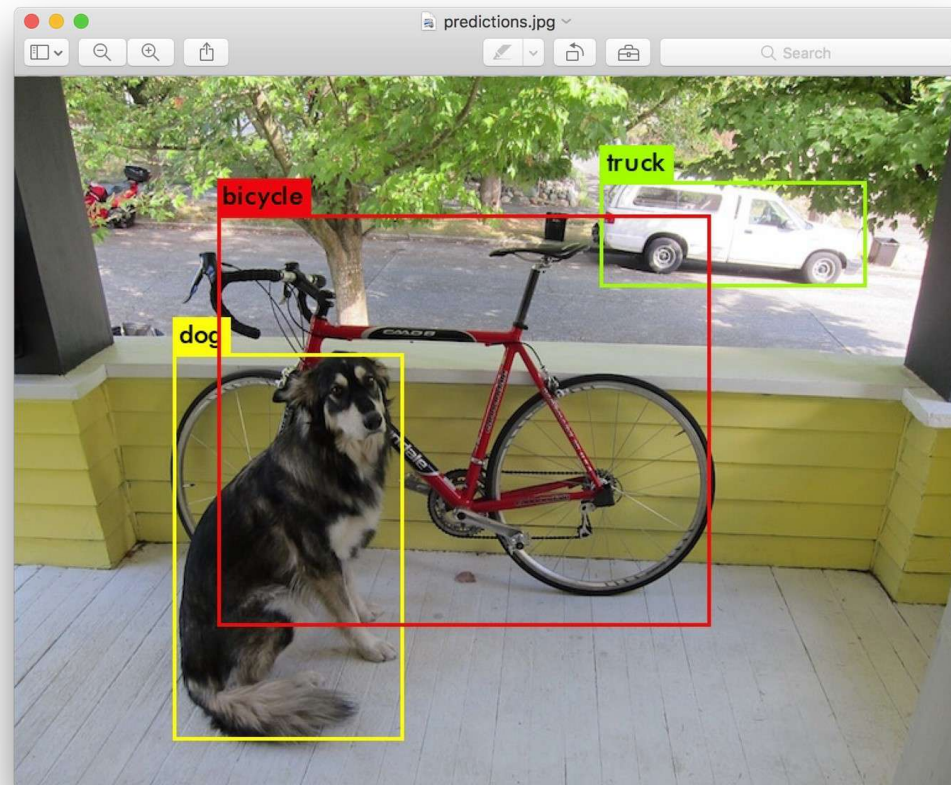
- What are ALL the things in the image?



[Redmon]

# Detection

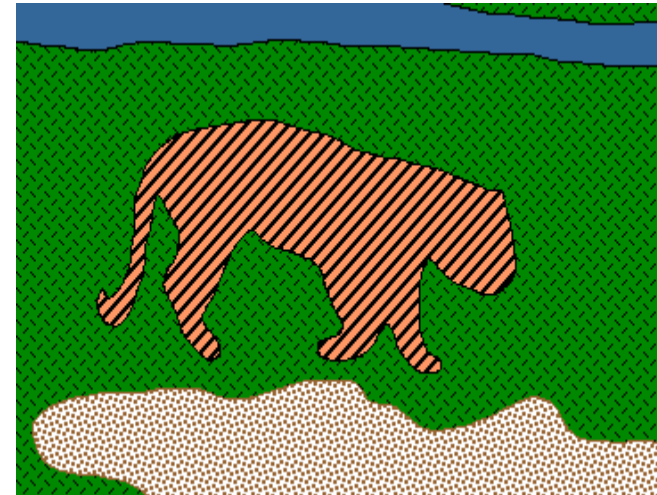
- What are ALL the things in the image?
- Where are they?



[Redmon]



# Segmentation





# Today's Agenda

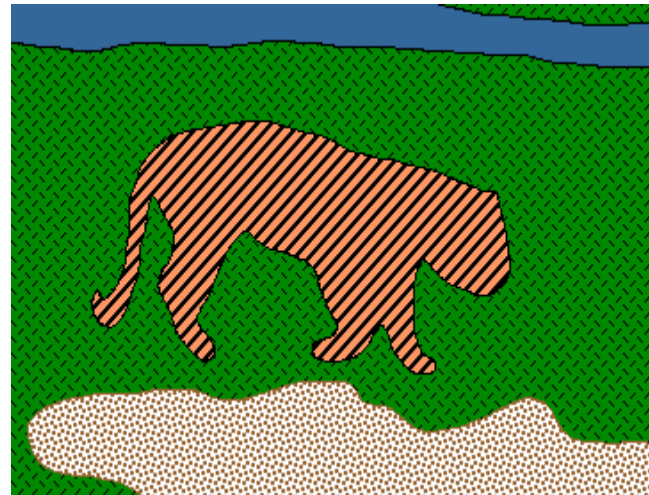
- Visual Recognition Tasks
- Introduction to segmentation and clustering
- Agglomerative clustering
- K-means clustering
- Mean-shift clustering
- Efficient Graph-based image segmentation

**Reading:** Forsyth Chapter 9

D. Comaniciu and P. Meer, [Mean Shift: A Robust Approach toward Feature Space Analysis](#), TPAMI 2002

# Image Segmentation

Goal: identify groups of pixels that go together



Slide credit: Steve Seitz, Kristen Grauman

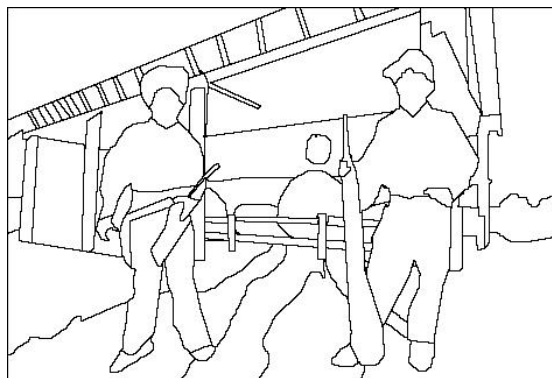
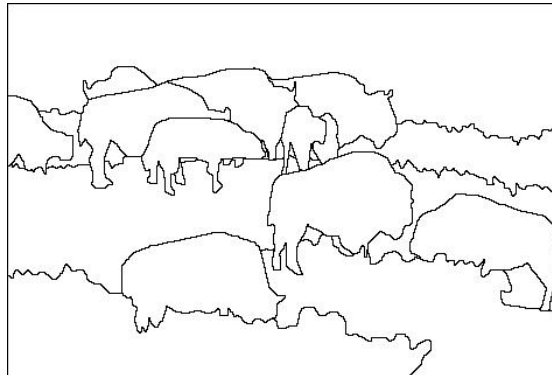
# The Goals of Segmentation

- Separate image into coherent “objects”

Image



Human segmentation



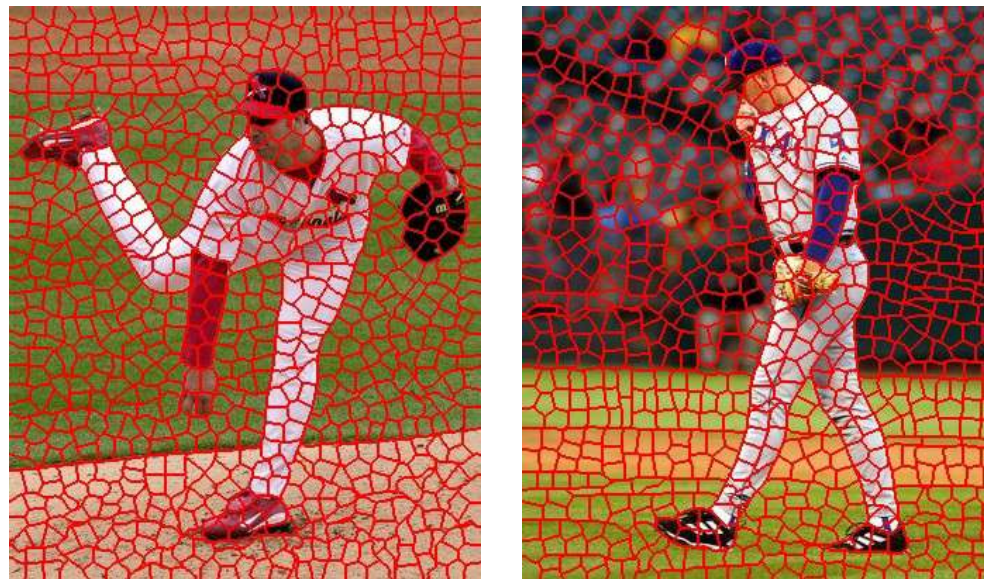
Slide credit: Svetlana Lazebnik



# The Goals of Segmentation

- Separate image into coherent “objects”
- Group together similar-looking pixels for efficiency of further processing

“superpixels”



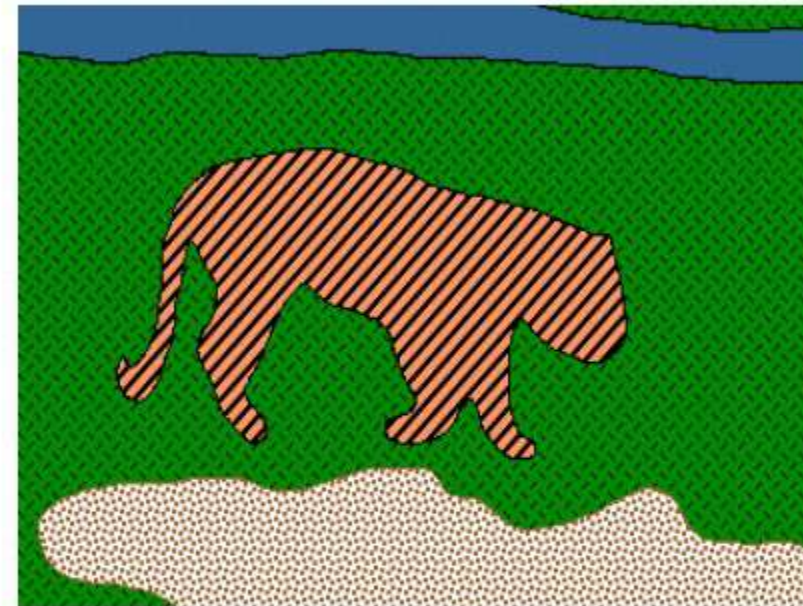
X. Ren and J. Malik. [Learning a classification model for segmentation](#). ICCV 2003.

Slide credit: Svetlana Lazebnik



# Types of Segmentation

- Semantic segmentation: Assign labels



Tiger  
Water

Grass  
Dirt

# Types of Segmentation

- Semantic segmentation: Assign labels



# Types of Segmentation

- Instance segmentation: Assign labels per object



<http://www.youtube.com/watch?v=OOT3UIXZztE>



# Types of Segmentation

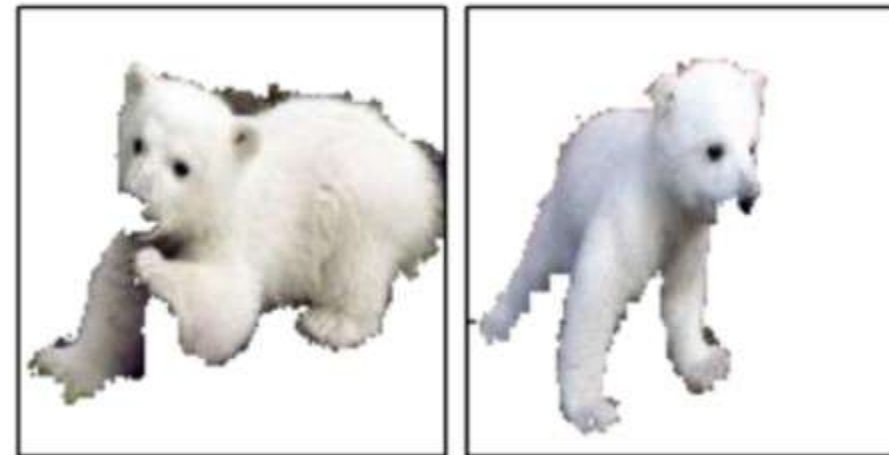
- Foreground / background segmentation





# Types of Segmentation

- Co-segmentation: Segment common object in multiple images



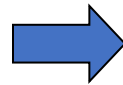
## Application: as a result



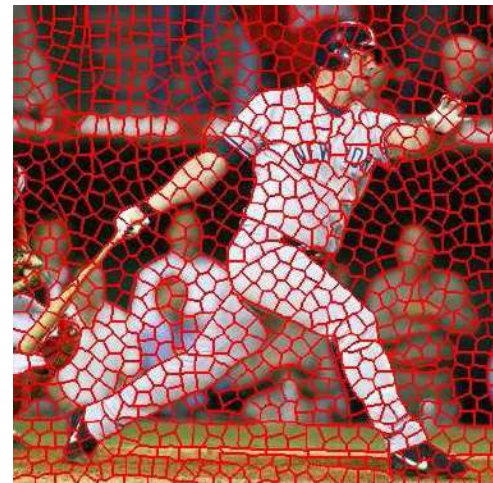
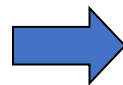
GrabCut: Rother et al. 2004



## Application: for efficiency – e.g., speed up recognition



[Felzenszwalb and Huttenlocher 2004]



[Hoiem et al. 2005, Mori 2005]

[Shi and Malik 2001]

Slide: Derek Hoiem

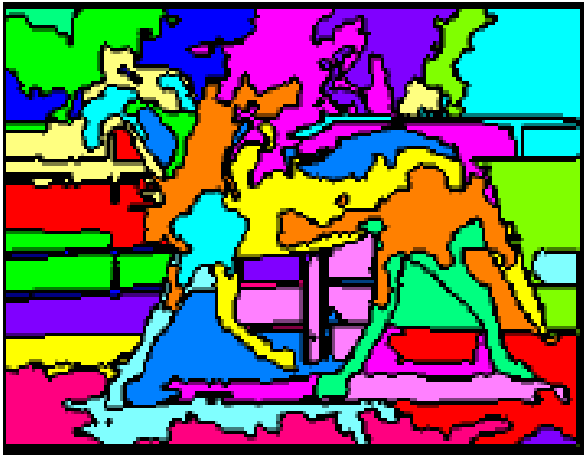
## Application: better classification



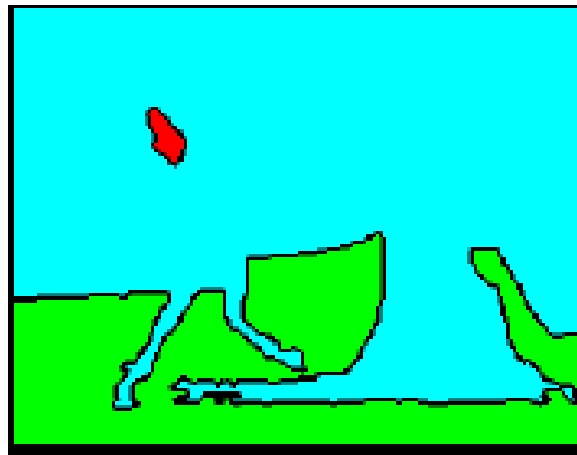
Angelova and Zhu, 2013



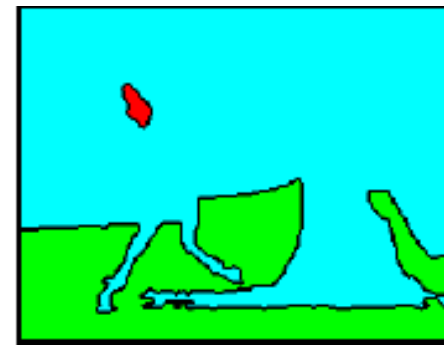
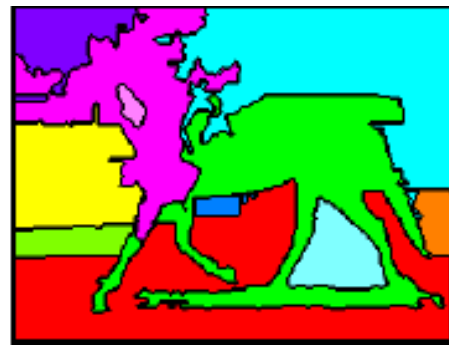
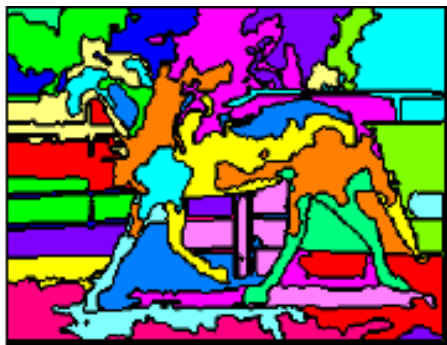
## Over/under segmentation



Oversegmentation



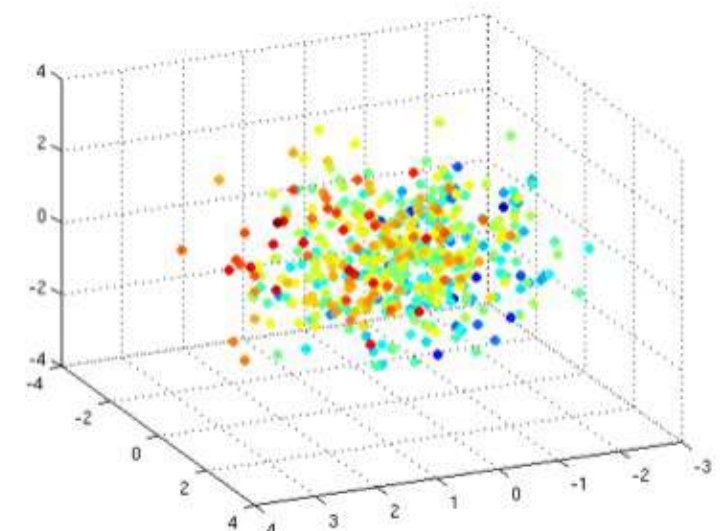
Undersegmentation



Multiple Segmentations

# One way to think about segmentation is Clustering

- Pixels are points in a (high-dimensional) feature space, e.g.
  - color: 3D
  - color + location: 5D
- Cluster pixels into segments



# One way to think about segmentation is Clustering

Clustering: group together similar data points and represent them as a single *entity*

Clustering is an unsupervised learning method

Key Challenges:

- 1) What makes two points/images/patches similar?
- 2) How do we compute an overall grouping from pairwise similarities?

Slide: Derek Hoiem

# Distance vs Similarity Measures

Let  $x$  and  $x'$  be two objects from the dataset.

The distance or similarity between  $x$  and  $x'$  is a real number,  $\text{dist}(x, x')$  or  $\text{sim}(x, x')$

- The Euclidian distance is defined as

$$\text{dist}(x, x') = \sqrt{\sum (x_i - x'_i)^2}$$

- In contrast, cosine similarity measure would be

$$\begin{aligned} \text{sim}(x, x') &= \cos(\theta) \\ &= \frac{x^\top x'}{\|x\| \cdot \|x'\|} \\ &= \frac{x^\top x'}{\sqrt{x^\top x} \sqrt{x'^\top x'}} \end{aligned}$$



# Desirable Properties of a Clustering Algorithm

1. Scalability in terms of both time and space
2. Ability to deal with different data types
3. Minimal requirements for domain knowledge to determine algorithm parameters
  - Don't need to know how many objects there are or what those object categories will be.
4. Interpretability and usability are optional
  - Incorporation of user-specified constraints

# General ideas

- Bottom-up clustering
  - pixels belong together because they are locally coherent
- Top-down clustering
  - pixels belong together because they lie on the same visual entity (object)

# Clustering algorithms

- Agglomerative clustering
- K-means
- Mean-shift clustering

# Today's Agenda

- Visual Recognition Tasks
- Introduction to segmentation and clustering
- Agglomerative clustering
- K-means clustering
- Mean-shift clustering
- Efficient Graph-based image segmentation

**Reading:** Forsyth Chapter 9

D. Comaniciu and P. Meer, [Mean Shift: A Robust Approach toward Feature Space Analysis](#), TPAMI 2002

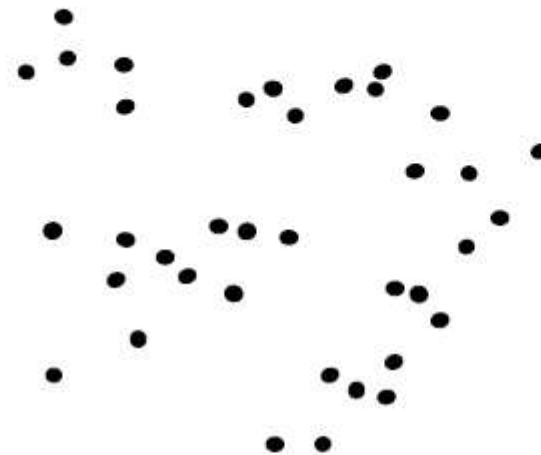


# Agglomerative Hierarchical Clustering - Algorithm

1. Initially each item  $x_1, \dots, x_n$  is in its own cluster  $C_1, \dots, C_n$ .
2. Repeat until there is only one cluster left:
3. Merge the nearest clusters, say  $C_i$  and  $C_j$ .



# Agglomerative Hierarchical Clustering - Algorithm

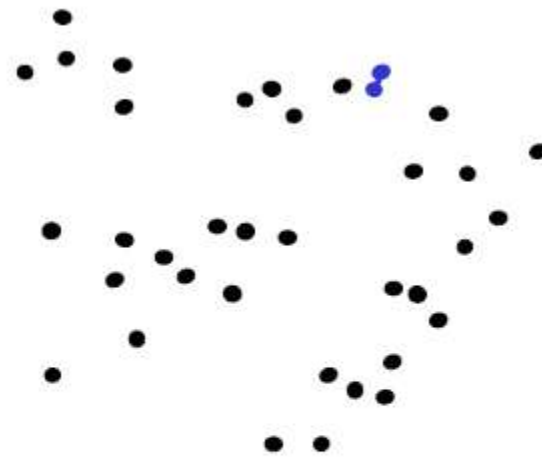


1. Say "Every point is its own cluster"

Slide credit: Andrew Moore



# Agglomerative Hierarchical Clustering - Algorithm

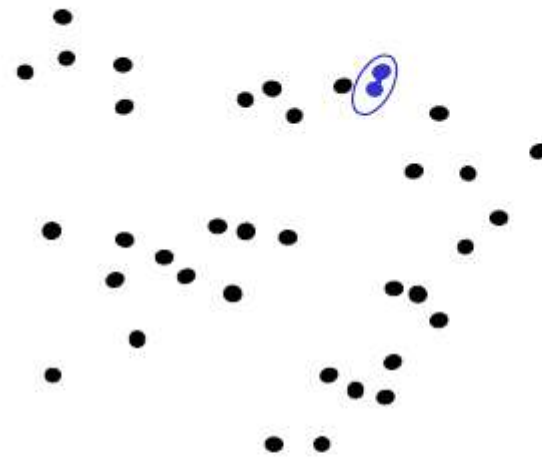


1. Say "Every point is its own cluster"
2. Find "most similar" pair of clusters



Slide credit: Andrew Moore

# Agglomerative Hierarchical Clustering - Algorithm



1. Say "Every point is its own cluster"
2. Find "most similar" pair of clusters
3. Merge it into a parent cluster

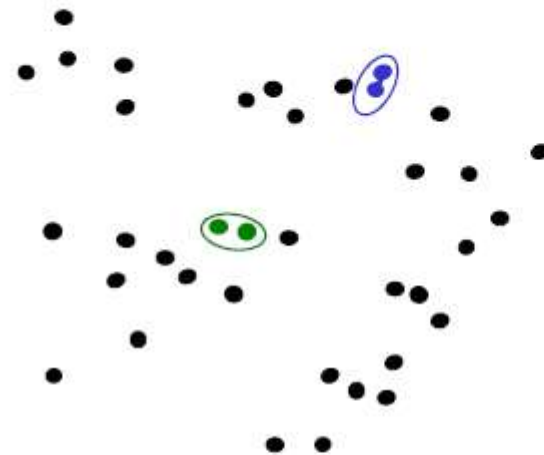


Slide credit: Andrew Moore





# Agglomerative Hierarchical Clustering - Algorithm



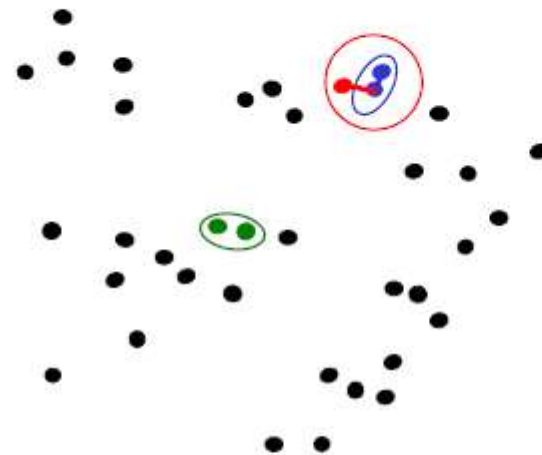
1. Say "Every point is its own cluster"
2. Find "most similar" pair of clusters
3. Merge it into a parent cluster
4. Repeat



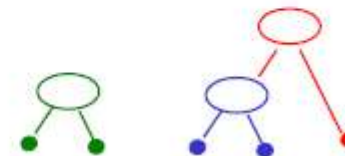
Slide credit: Andrew Moore



# Agglomerative Hierarchical Clustering - Algorithm

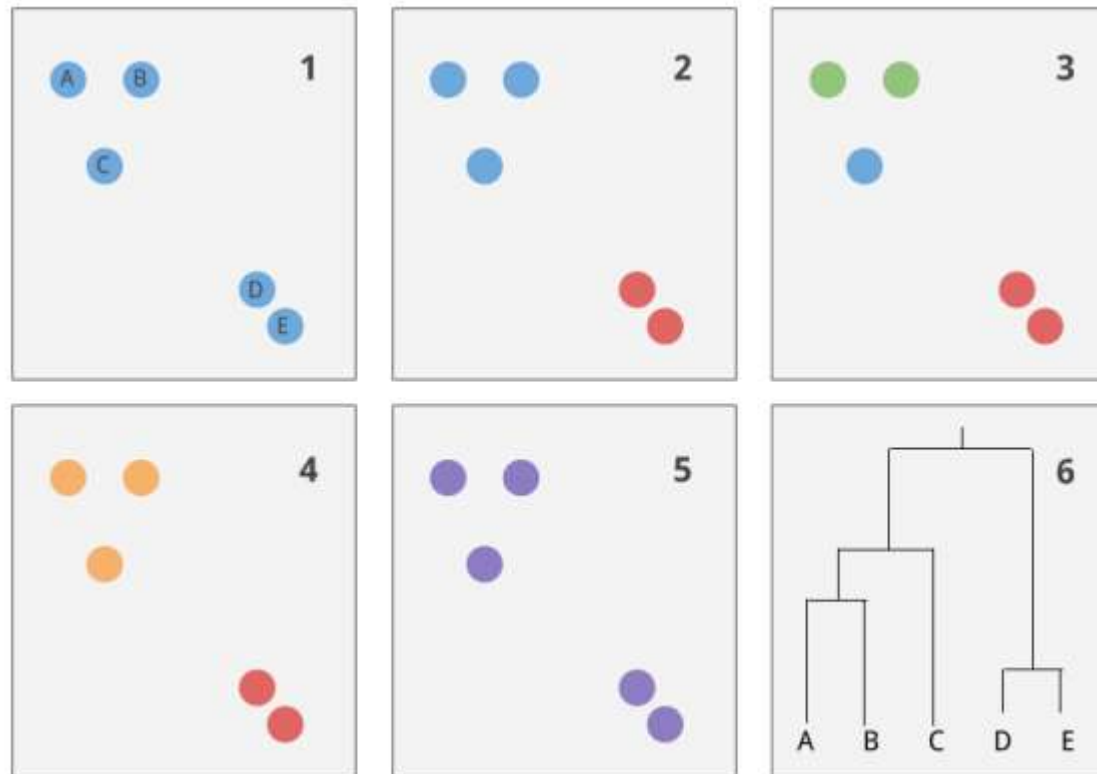


1. Say "Every point is its own cluster"
2. Find "most similar" pair of clusters
3. Merge it into a parent cluster
4. Repeat



Slide credit: Andrew Moore

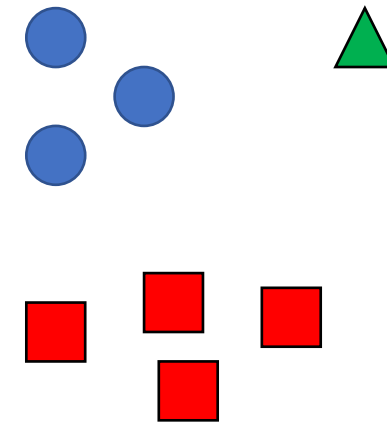
# Agglomerative clustering example



# Agglomerative clustering

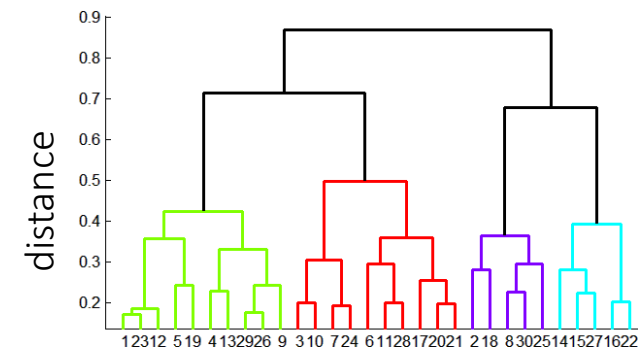
How to define cluster similarity?

- average distance between points
- maximum distance
- minimum distance
- distance between means or medoids



How many clusters?

- Clustering creates a dendrogram (a tree)
- Threshold based on max number of clusters or based on distance between merges

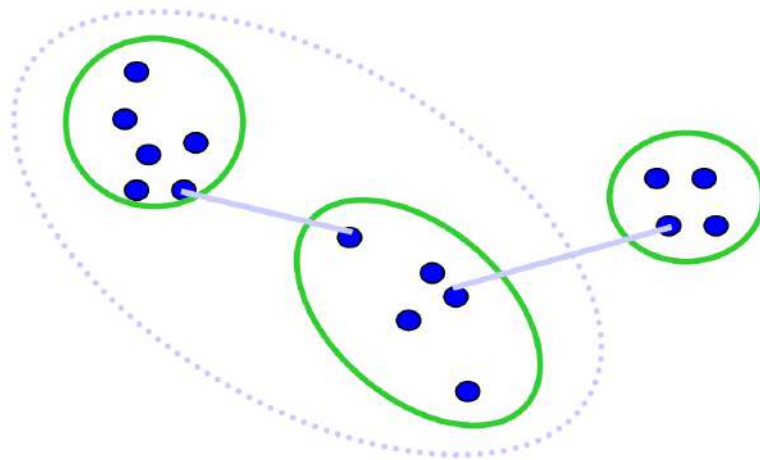




# Different measures of nearest clusters

## Single Link

- Distance between clusters is the minimum distance between their points

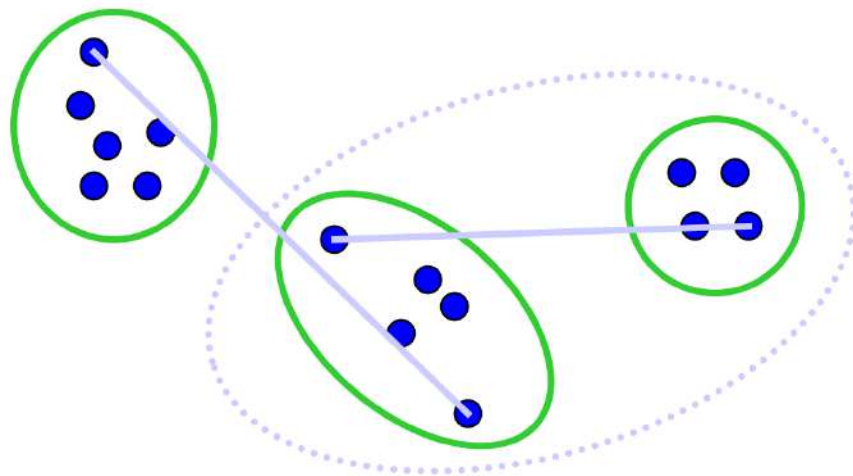


Long, skinny clusters

# Different measures of nearest clusters

## Complete Link

- Distance between clusters is the maximum distance between their points

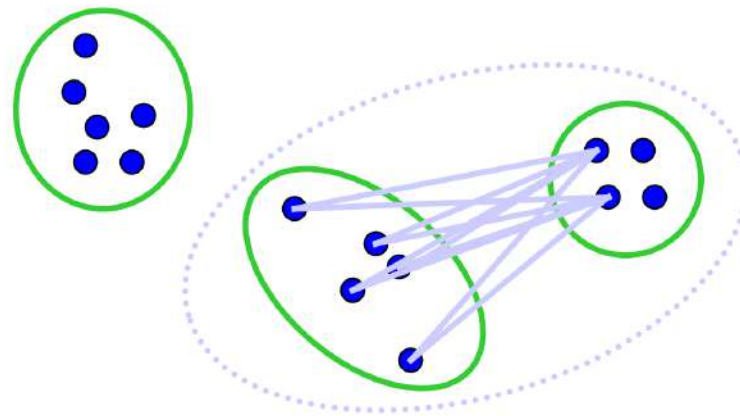


Tight clusters

# Different measures of nearest clusters

## Average Link

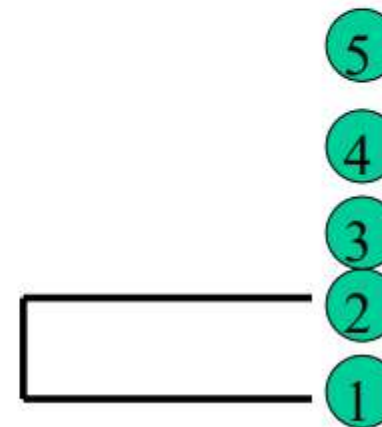
- Distance between clusters is the average distance between their points



Robust against noise

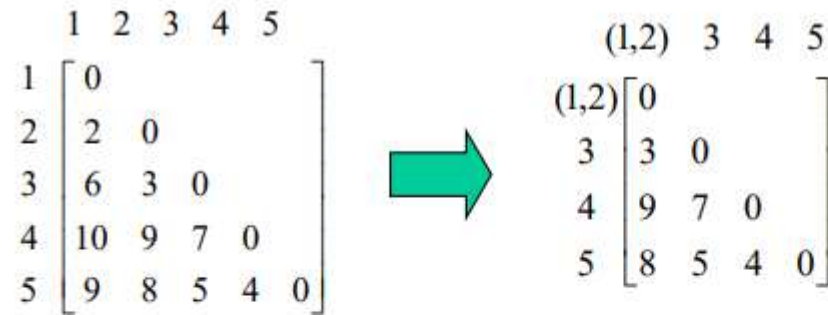
# Example – single link

	1	2	3	4	5
1	0				
2	2	0			
3	6	3	0		
4	10	9	7	0	
5	9	8	5	4	0





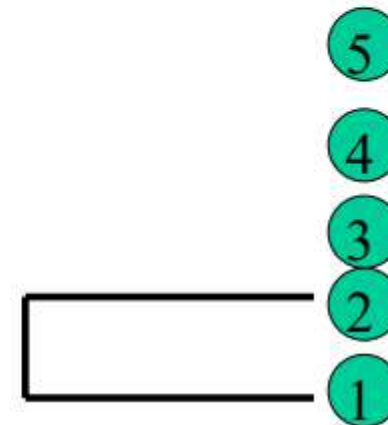
# Example – single link



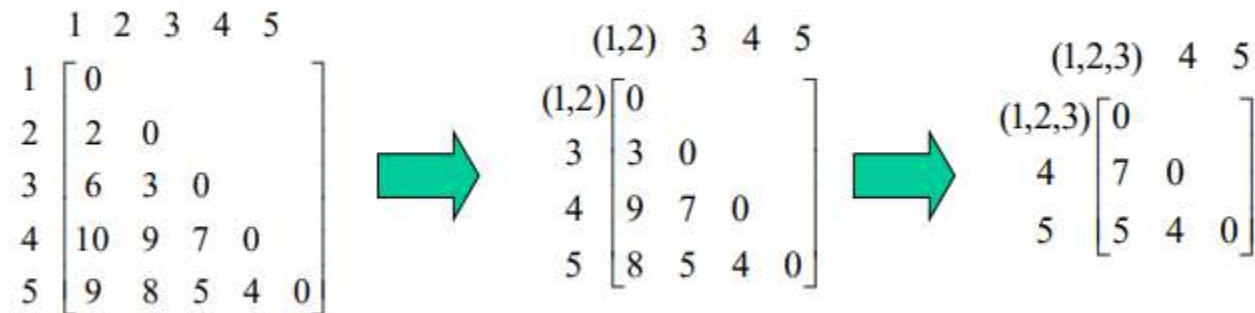
$$d_{(1,2),3} = \min\{d_{1,3}, d_{2,3}\} = \min\{6,3\} = 3$$

$$d_{(1,2),4} = \min\{d_{1,4}, d_{2,4}\} = \min\{10,9\} = 9$$

$$d_{(1,2),5} = \min\{d_{1,5}, d_{2,5}\} = \min\{9,8\} = 8$$

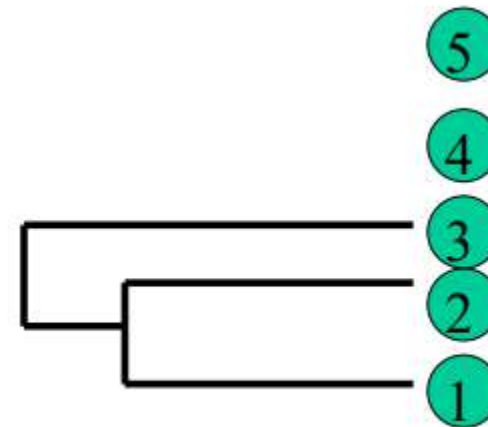


## Example – single link

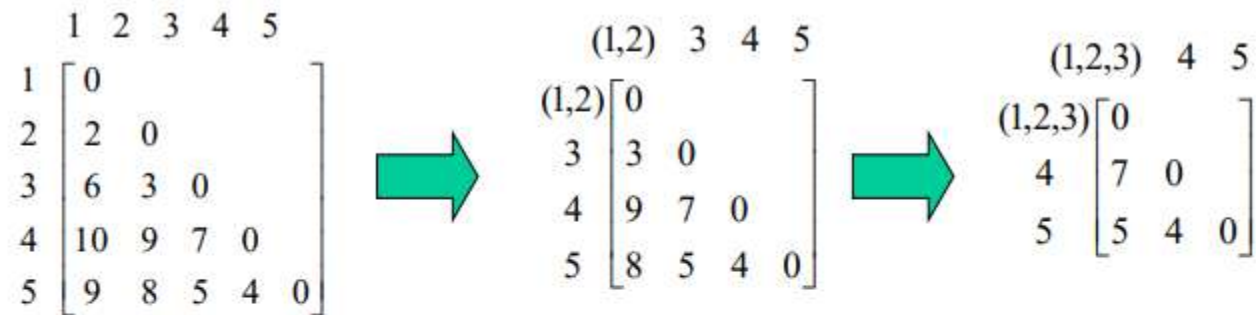


$$d_{(1,2,3),4} = \min\{d_{(1,2),4}, d_{3,4}\} = \min\{9, 7\} = 7$$

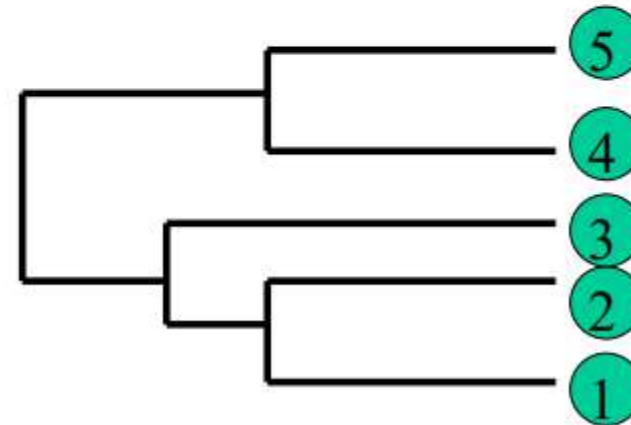
$$d_{(1,2,3),5} = \min\{d_{(1,2),5}, d_{3,5}\} = \min\{8, 5\} = 5$$



## Example – single link

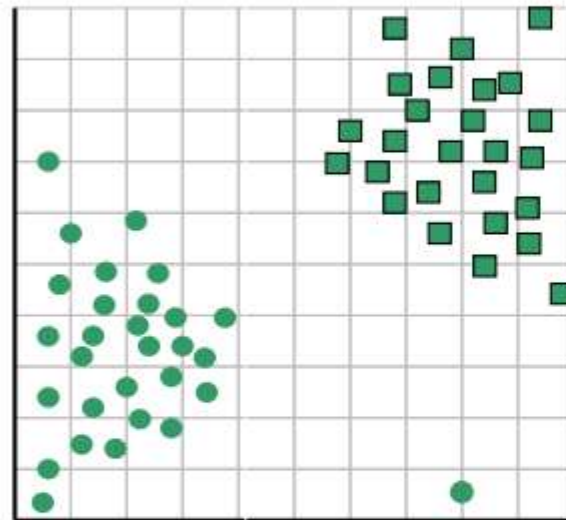


$$d_{(1,2,3),(4,5)} = \min\{d_{(1,2,3),4}, d_{(1,2,3),5}\} = 5$$

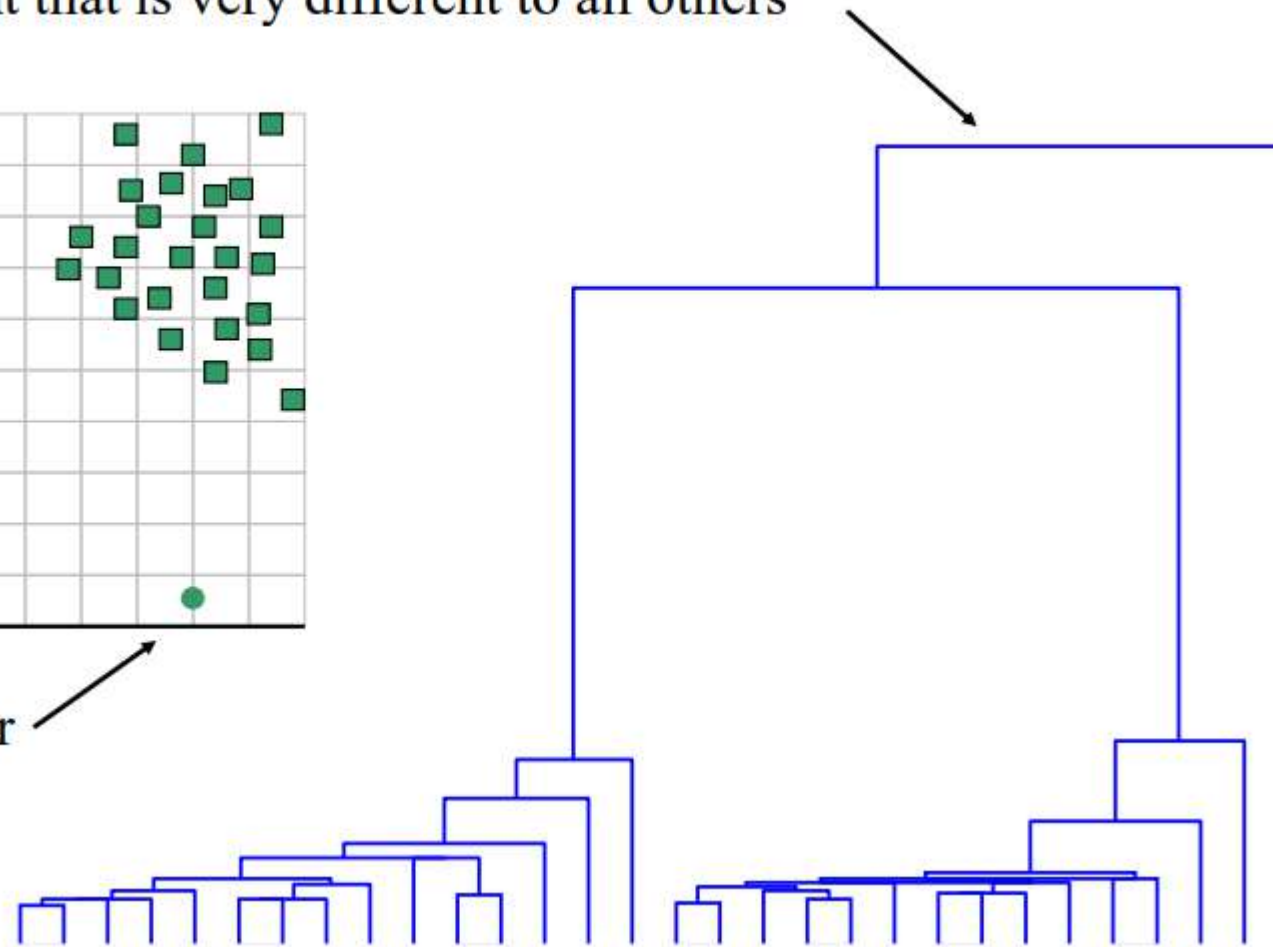


# Outliers

The single isolated branch is suggestive of a data point that is very different to all others



Outlier





# Conclusions: Agglomerative Clustering

## Good

- Simple to implement, widespread application.
- Clusters have adaptive shapes.
- Provides a hierarchy of clusters.
- Can avoid specifying number of clusters in advance.

## Bad

- May have imbalanced clusters.
- Still have to choose number of clusters or threshold to use them.
- Does not scale well. Runtime of  $O(n^3)$ .

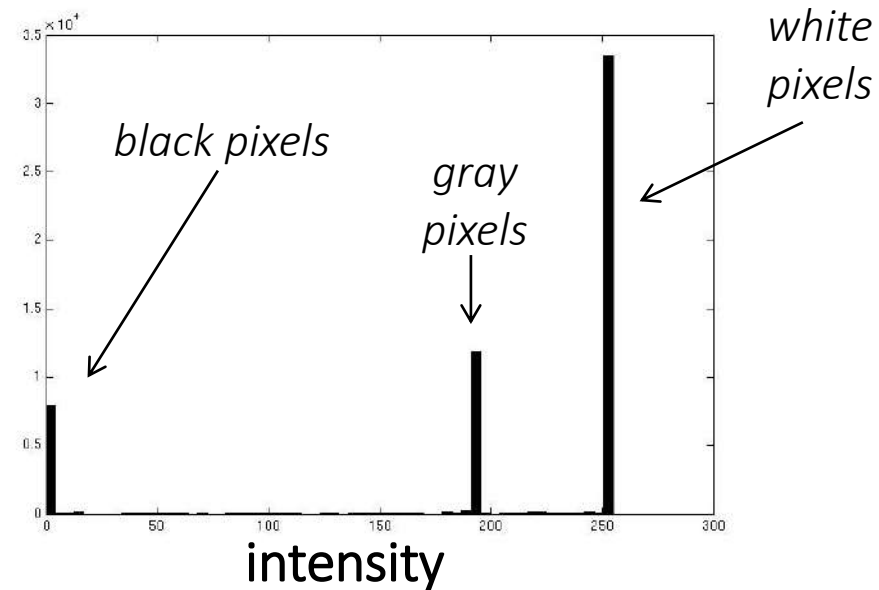
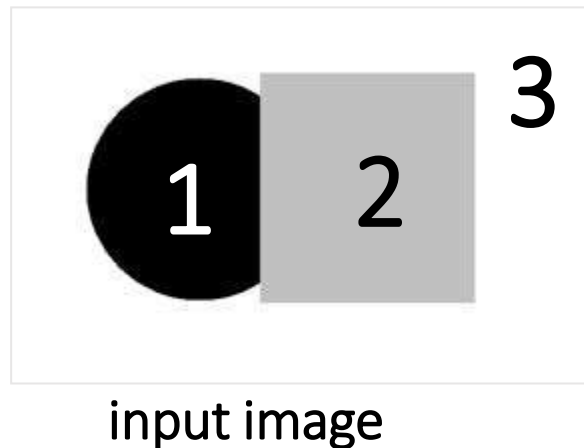
# Today's Agenda

- Visual Recognition Tasks
- Introduction to segmentation and clustering
- Agglomerative clustering
- K-means clustering
- Mean-shift clustering
- Efficient Graph-based image segmentation

**Reading:** Forsyth Chapter 9

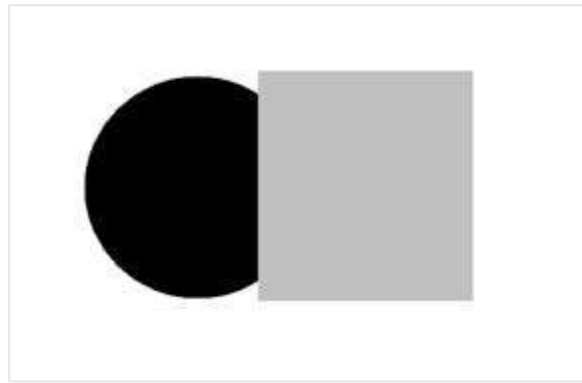
D. Comaniciu and P. Meer, [Mean Shift: A Robust Approach toward Feature Space Analysis](#), TPAMI 2002

# Image Segmentation: Toy Example

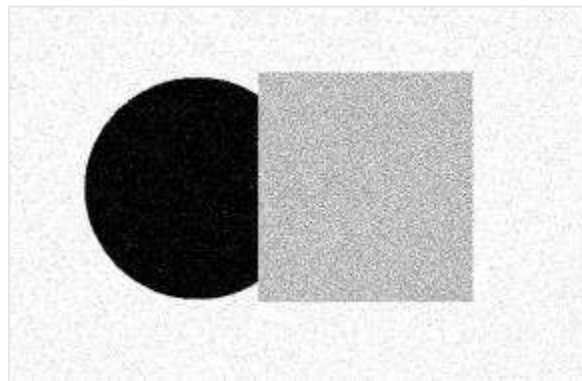
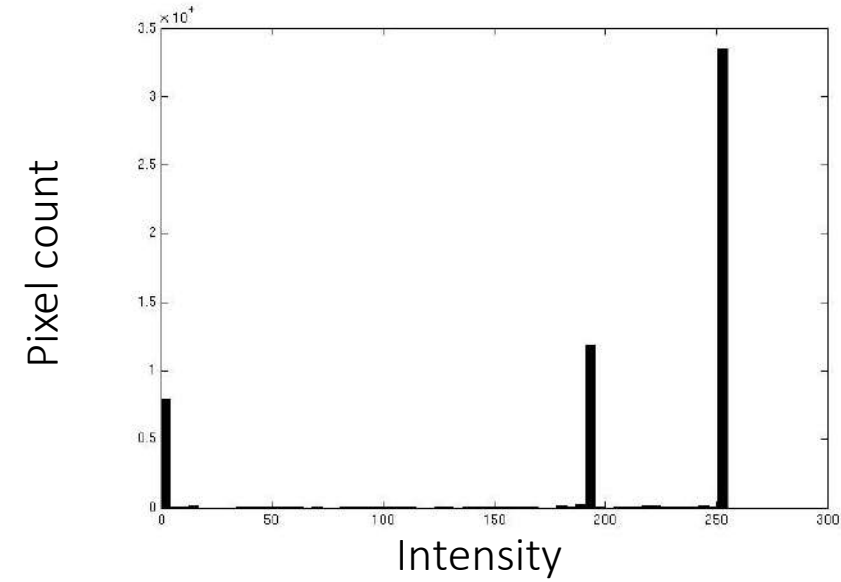


- These intensities define the three groups.
- We could label every pixel in the image according to which of these primary intensities it is.
  - i.e., segment the image based on the image intensity feature.
- What if the image isn't quite so simple?

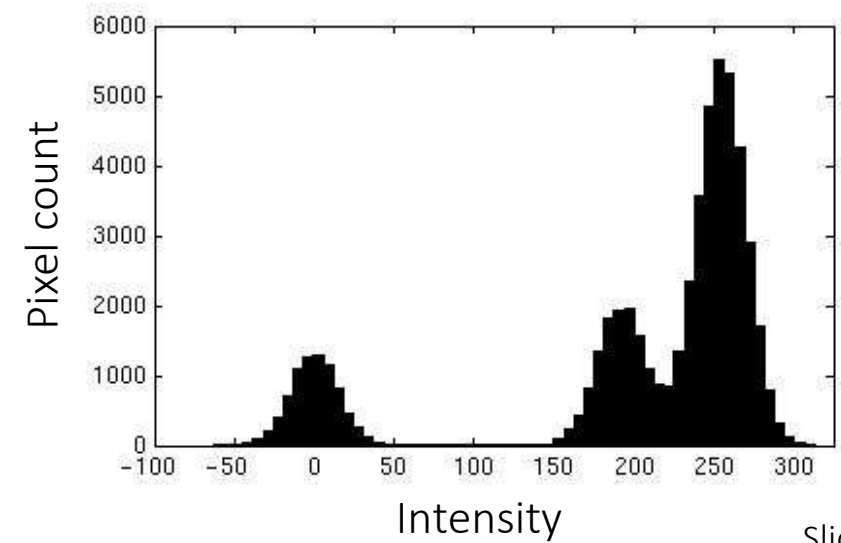
Slide credit: Kristen Grauman



Input image

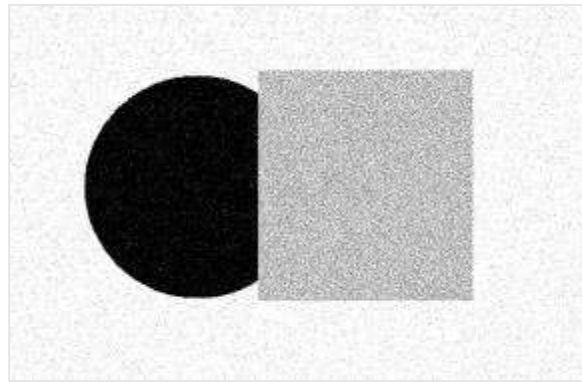


Input image

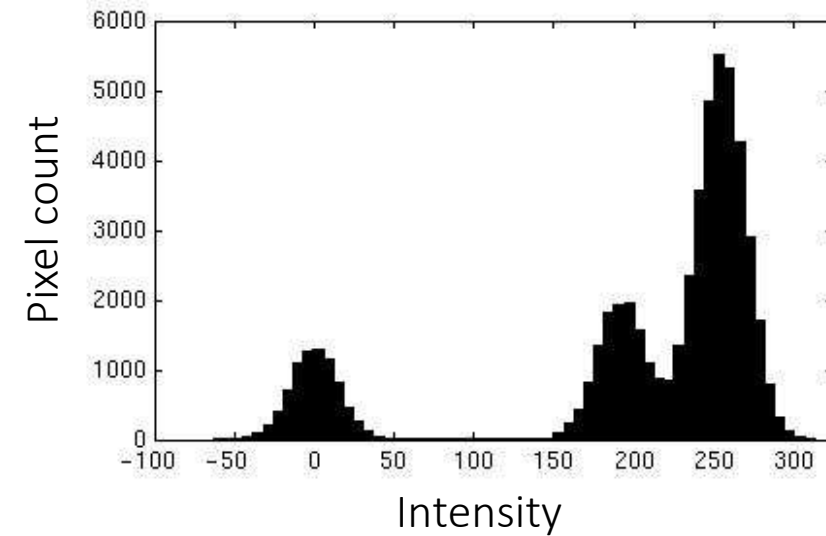


Slide credit: Kristen Grauman



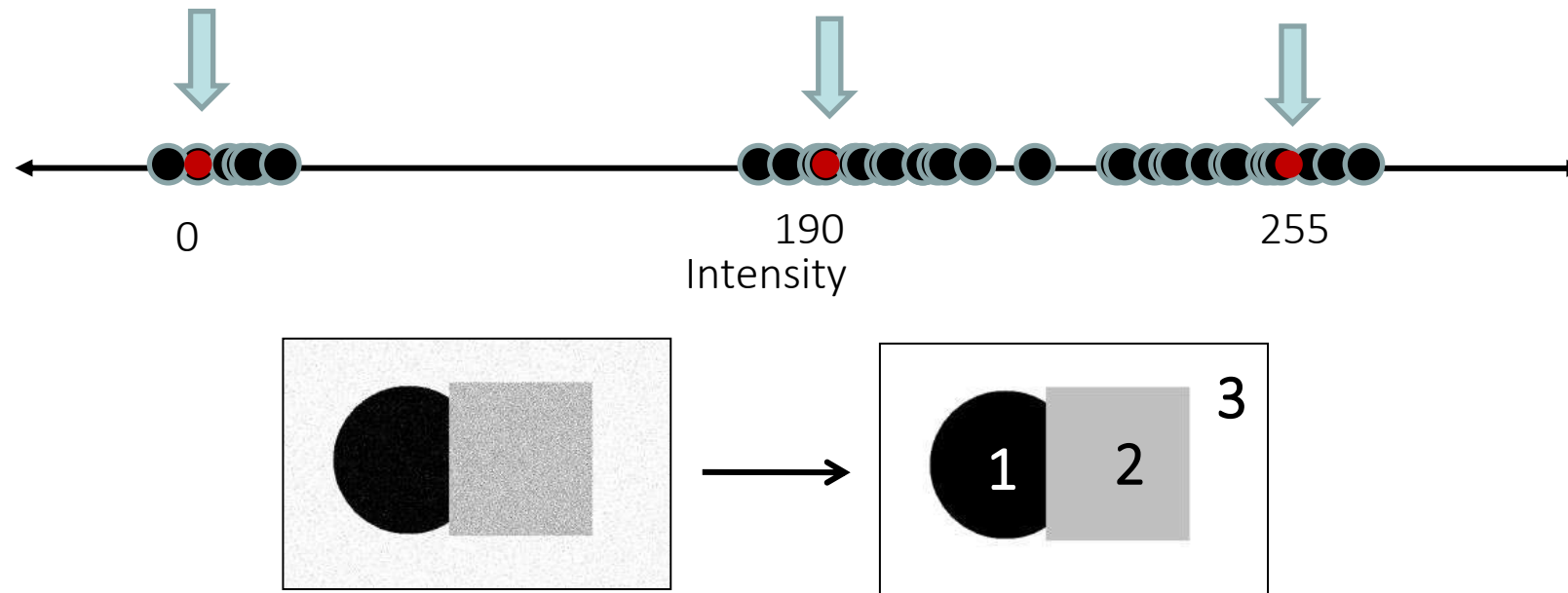


Input image



- Now how to determine the three main intensities that define our groups?
- We need to cluster

Slide credit: Kristen Grauman



- Goal: choose three “centers” as the representative intensities and label every pixel according to which of these centers it is nearest to.
- Best cluster centers are those that minimize Sum of Square Distance (SSD) between all points and their nearest cluster center  $c_i$ :

$$SSD = \sum_{cluster\ i} \sum_{x \in cluster\ i} (x - c_i)^2$$

Slide credit: Kristen Grauman

# Clustering for Summarization

Goal: cluster to minimize variance in data, given clusters

$$c^*, \delta^* = \arg \min_{c, \delta} \frac{1}{N} \sum_j^N \sum_i^K \delta_{ij} (c_i - x_j)^2$$

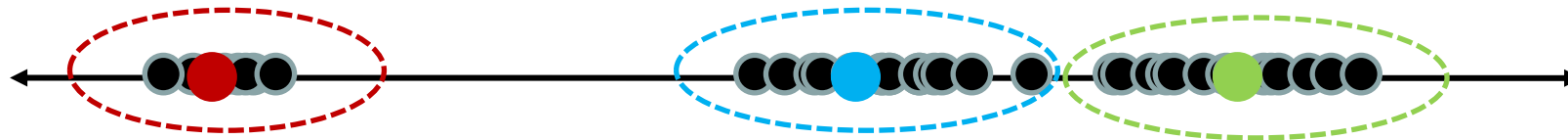
Cluster center
Data

Whether  $x_j$  is assigned to  $c_i$

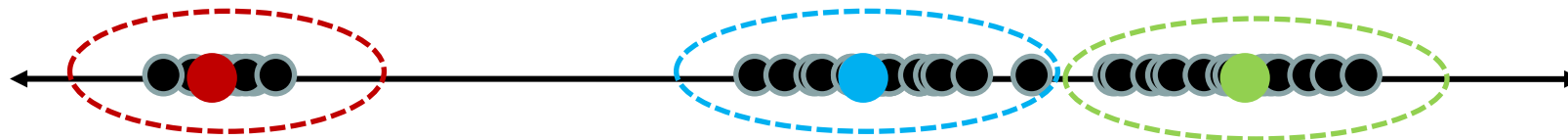
Slide credit: Derek Hoiem

# Clustering

- With this objective, it is a “chicken and egg” problem:
  - If we knew the *cluster centers*, we could allocate points to groups by assigning each to its closest center.



- If we knew the *group memberships*, we could get the centers by computing the mean per group.



Slide credit: Kristen Grauman



# K-means clustering

1. Initialize ( $t = 0$ ): cluster centers  $c_1, \dots, c_K$

2. Compute  $\delta^t$ : assign each point to the closest center

- $\delta^t$  denotes the set of assignments for each  $x_j$  to cluster  $c_i$  at iteration  $t$

$$\delta^t = \operatorname{argmin}_{\delta} \frac{1}{N} \sum_j \sum_i \delta_{ij}^t (c_i^{t-1} - x_j)^2$$

3. Compute  $c^t$ : update cluster centers as the mean of the points

$$c^t = \operatorname{argmin}_c \frac{1}{N} \sum_j \sum_i \delta_{ij}^t (c_i^{t-1} - x_j)^2$$

4. Update  $t = t + 1$ , Repeat Step 2-3 till stopped



Slide credit: Derek Hoiem

# K-means clustering

1. Initialize ( $t = 0$ ): cluster centers  $c_1, \dots, c_K$

- Commonly used: random initialization
- Or greedily choose  $K$  to minimize residual

2. Compute  $\delta^t$ : assign each point to the closest center

- $\delta^t$  denotes the set of assignments for each  $x_j$  to cluster  $c_i$  at iteration  $t$
- Typical distance measure:
  - Euclidean
  - Cosine

$$\delta^t = \operatorname{argmin}_{\delta} \frac{1}{N} \sum_j \sum_i \delta_{ij}^t (c_i^{t-1} - x_j)^2$$

3. Compute  $c^t$ : update cluster centers as the mean of the points

$$c^t = \operatorname{argmin}_c \frac{1}{N} \sum_j \sum_i \delta_{ij}^t (c_i^{t-1} - x_j)^2$$

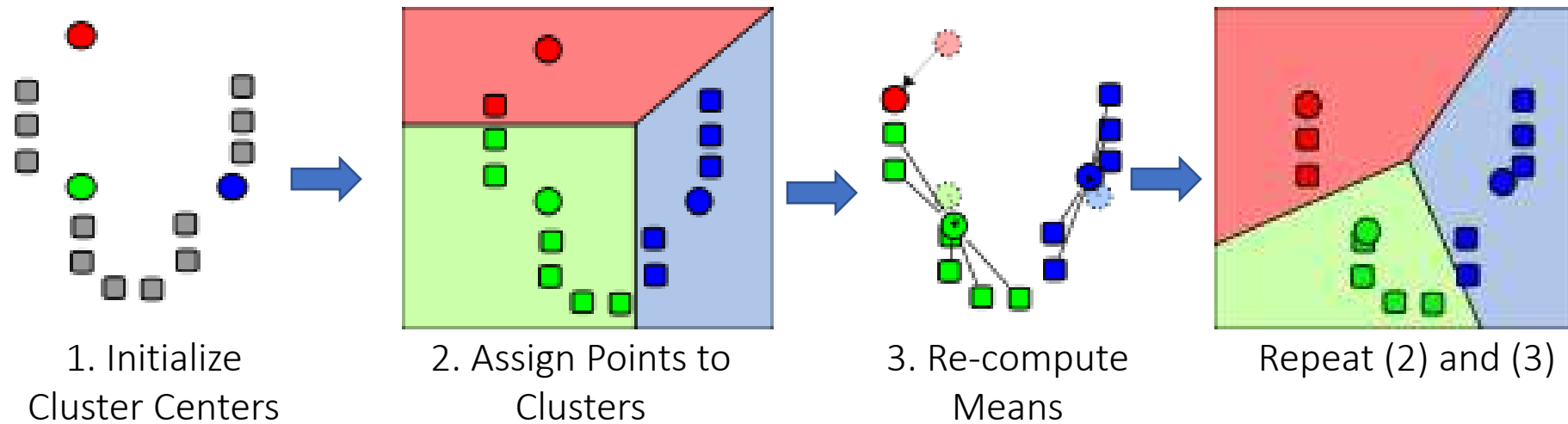
4. Update  $t = t + 1$ , Repeat Step 2-3 till stopped

- $c^t$  doesn't change anymore.



Slide credit: Derek Hoiem

# K-means clustering



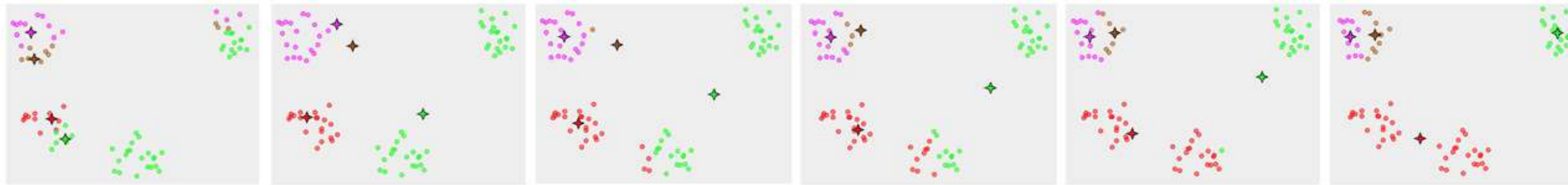
## Demo

<https://stanford.edu/class/engr108/visualizations/kmeans/kmeans.html>

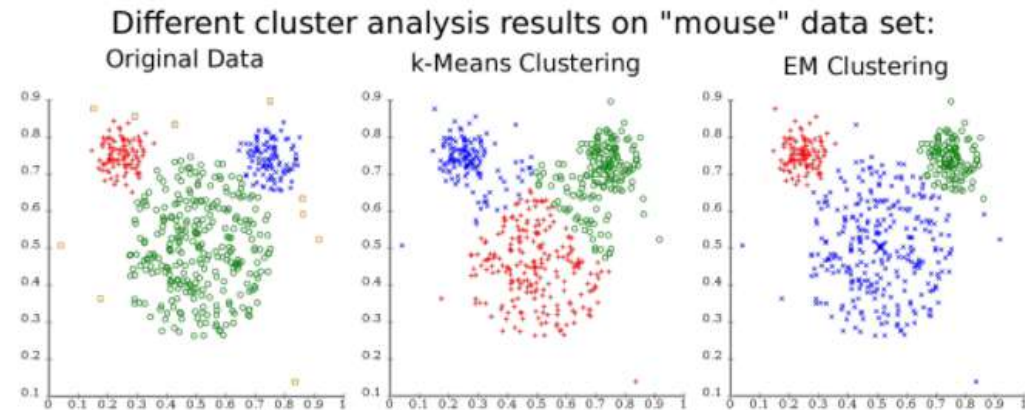
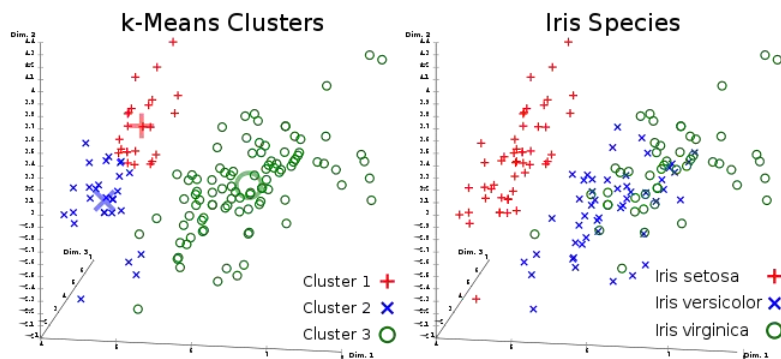
Illustration Source: wikipedia

# K-means clustering

- Converges to a *local minimum* solution
  - Initialize multiple runs



- Better fit for spherical data



- Need to pick K (# of clusters)



## Segmentation as Clustering



Original image



2 clusters



3 clusters

# Feature Space

- Depending on what we choose as the *feature space*, we can group pixels in different ways.
- Grouping pixels based on **intensity** similarity

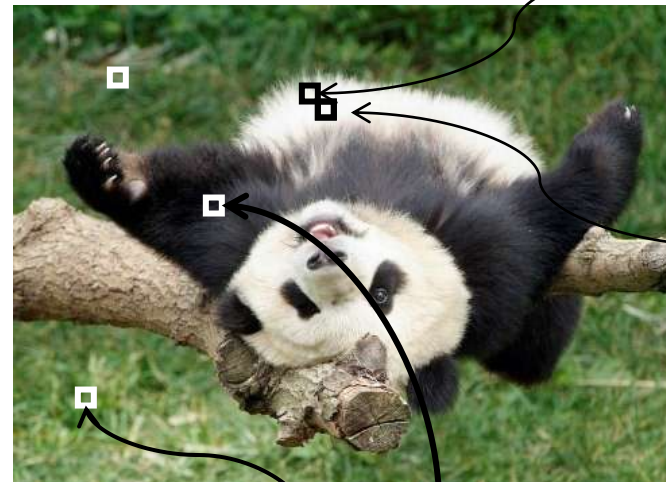
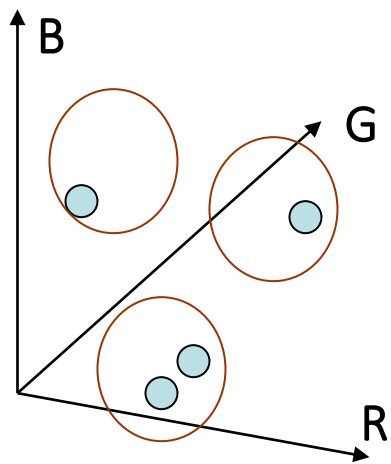


- Feature space: intensity value (1D)

Slide credit: Kristen Grauman

# Feature Space

- Depending on what we choose as the *feature space*, we can group pixels in different ways.
- Grouping pixels based on **color** similarity



R=255  
G=200  
B=250

R=245  
G=220  
B=248

R=15  
G=189  
B=2

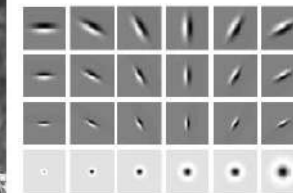
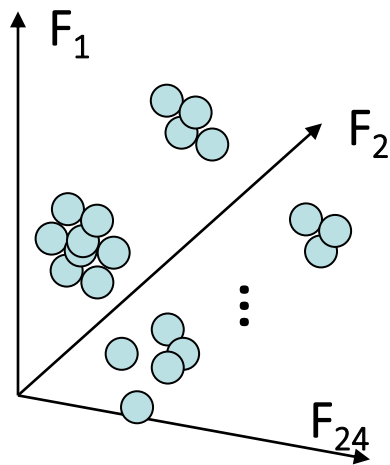
R=3  
G=12  
B=2

- Feature space: color value (3D)

Slide credit: Kristen Grauman

# Feature Space

- Depending on what we choose as the *feature space*, we can group pixels in different ways.
- Grouping pixels based on **texture** similarity



Filter bank of 24 filters

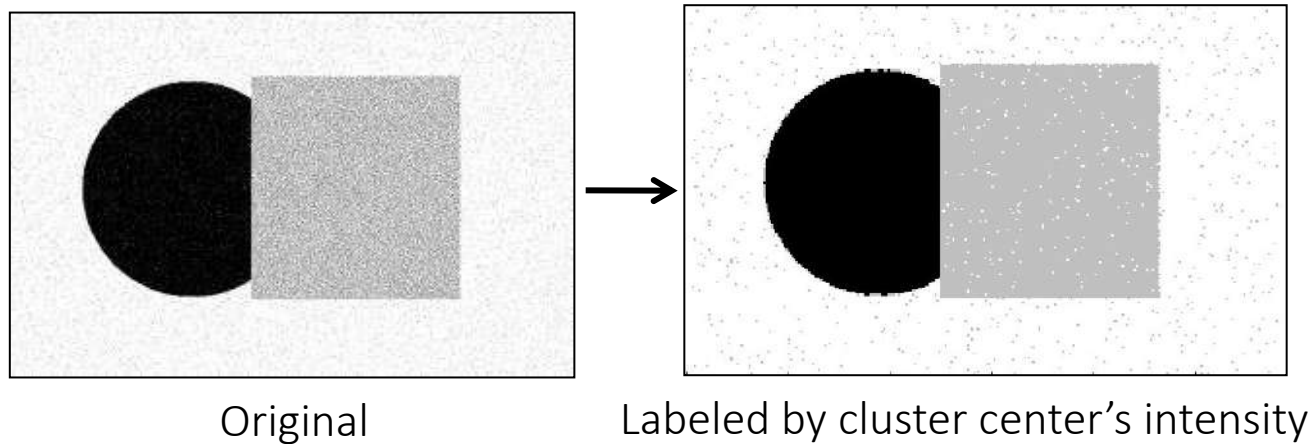
- Feature space: filter bank responses (e.g., 24D)

Slide credit: Kristen Grauman

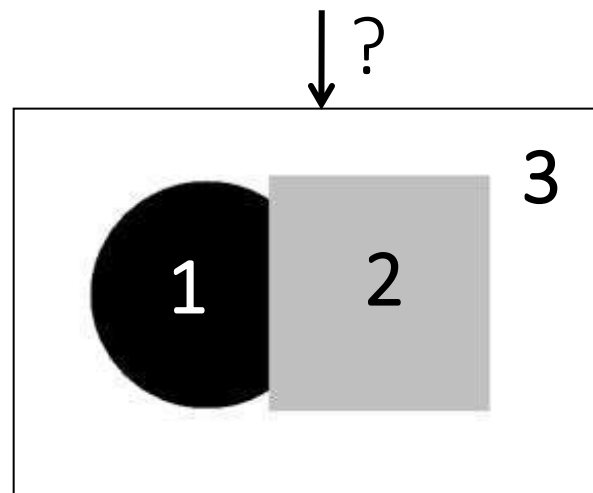


# Smoothing Out Cluster Assignments

- Assigning a cluster label per pixel may yield outliers:



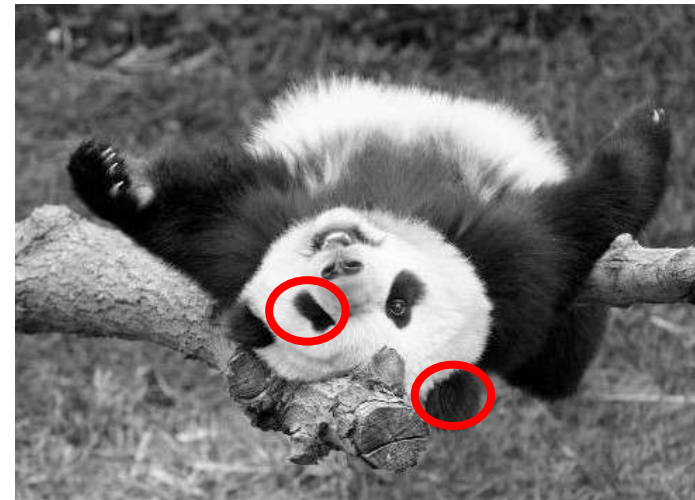
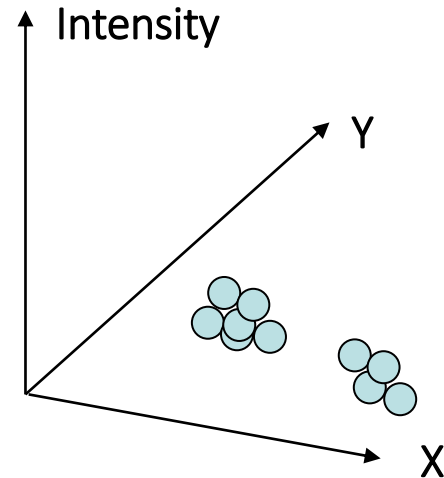
- How can we ensure they are spatially smooth?



Slide credit: Kristen Grauman

# Segmentation as Clustering

- Depending on what we choose as the *feature space*, we can group pixels in different ways.
- Grouping pixels based on *intensity+position* similarity



⇒ Way to encode both *similarity* and *proximity*.

Slide credit: Kristen Grauman

# K-Means Clustering Results

- K-means clustering based on intensity or color is essentially vector quantization of the image attributes
  - Clusters don't have to be spatially coherent



Image source: Forsyth & Ponce

# K-Means Clustering Results

- K-means clustering based on intensity or color is essentially vector quantization of the image attributes
  - Clusters don't have to be spatially coherent
- Clustering based on  $(r,g,b,x,y)$  values enforces more spatial coherence

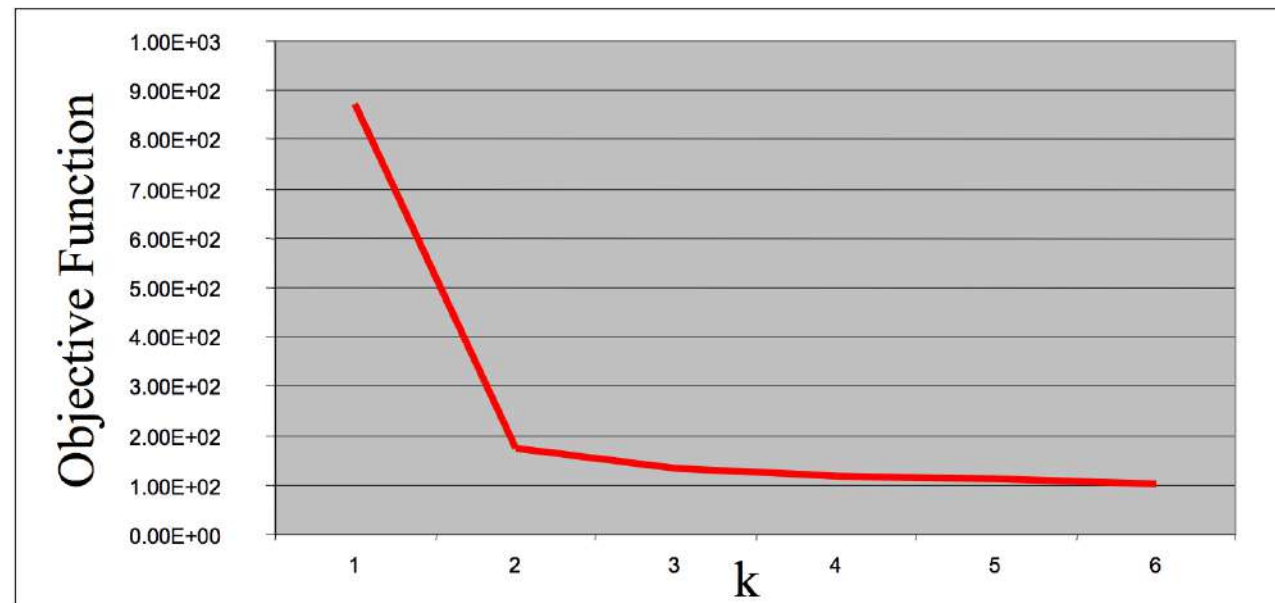


# How to choose the number of clusters?

Try different numbers of clusters in a validation set and look at performance.

We can plot the objective function values for  $k$  equals 1 to 6...

The abrupt change at  $k = 2$ , is highly suggestive of two clusters in the data. This technique for determining the number of clusters is known as “knee finding” or “elbow finding”.



Slide credit: Derek Hoiem

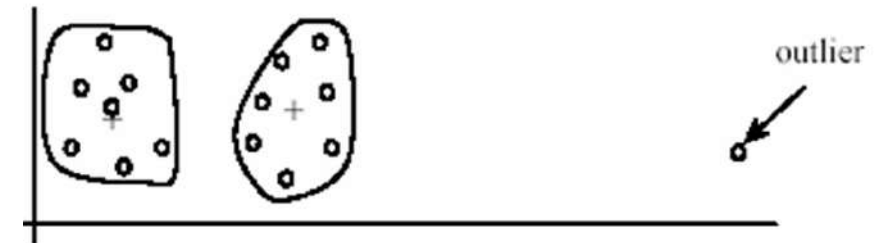
# K-Means pros and cons

- Pros

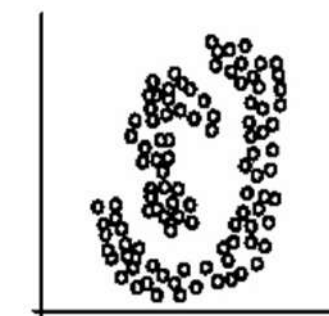
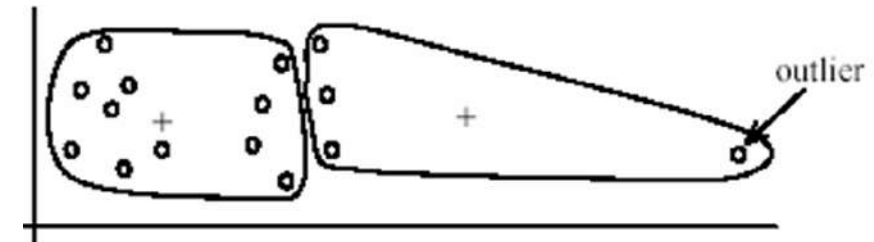
- Finds cluster centers that minimize conditional variance (good representation of data)
- Simple and fast, Easy to implement

- Cons

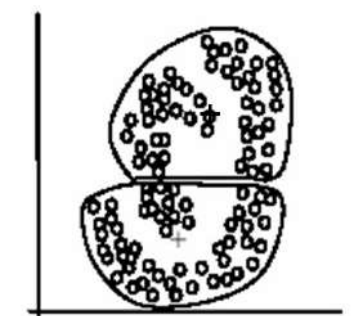
- Need to choose K
- Sensitive to outliers
- Prone to local minima
- All clusters have the same parameters (e.g., distance measure is non-adaptive)
- Distance computation in N-dimensional space could be slow



(B): Ideal clusters



(A): Two natural clusters



(B):  $k$ -means clusters

# Today's Agenda

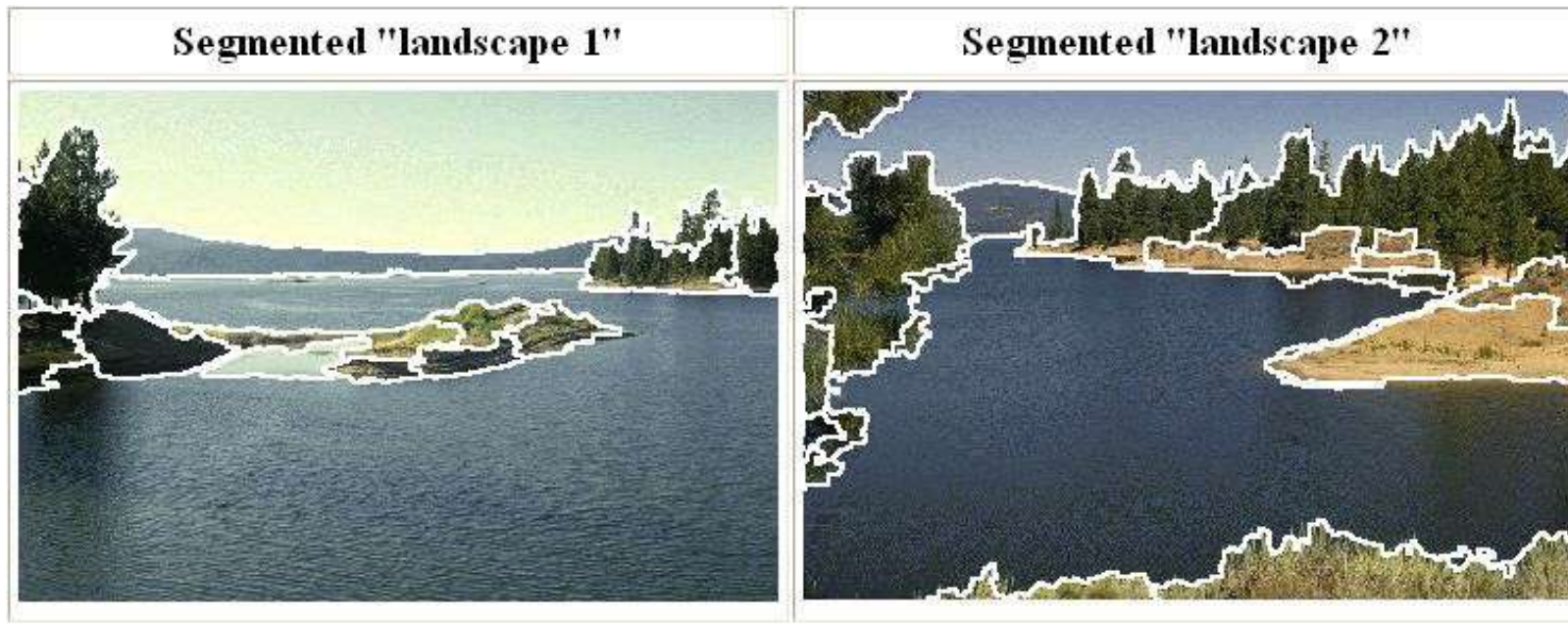
- Visual Recognition Tasks
- Introduction to segmentation and clustering
- Agglomerative clustering
- K-means clustering
- Mean-shift clustering
- Efficient Graph-based image segmentation

**Reading:** Forsyth Chapter 9

D. Comaniciu and P. Meer, [Mean Shift: A Robust Approach toward Feature Space Analysis](#), TPAMI 2002

# Mean-Shift Segmentation

- An advanced and versatile technique for clustering-based segmentation

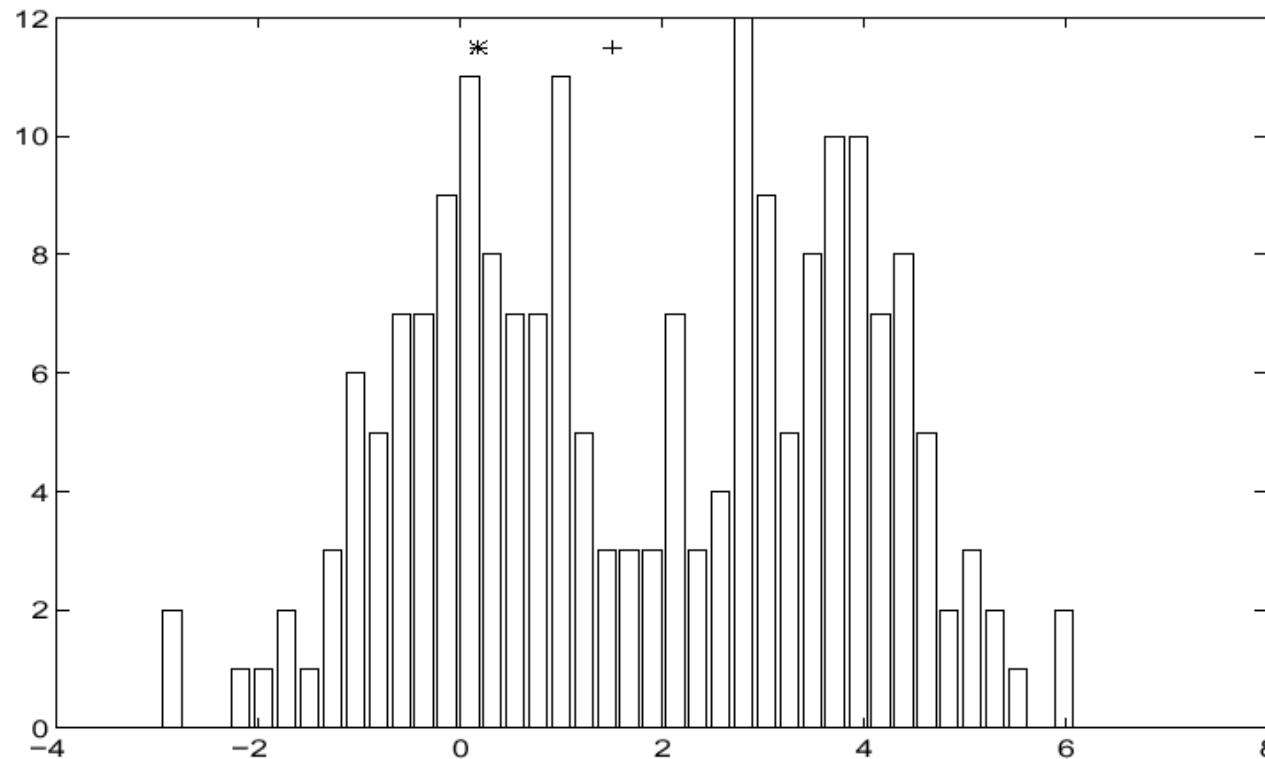


D. Comaniciu and P. Meer, [Mean Shift: A Robust Approach toward Feature Space Analysis](#), TPAMI 2002

Slide credit: Svetlana Lazebnik



# Mean-Shift Algorithm



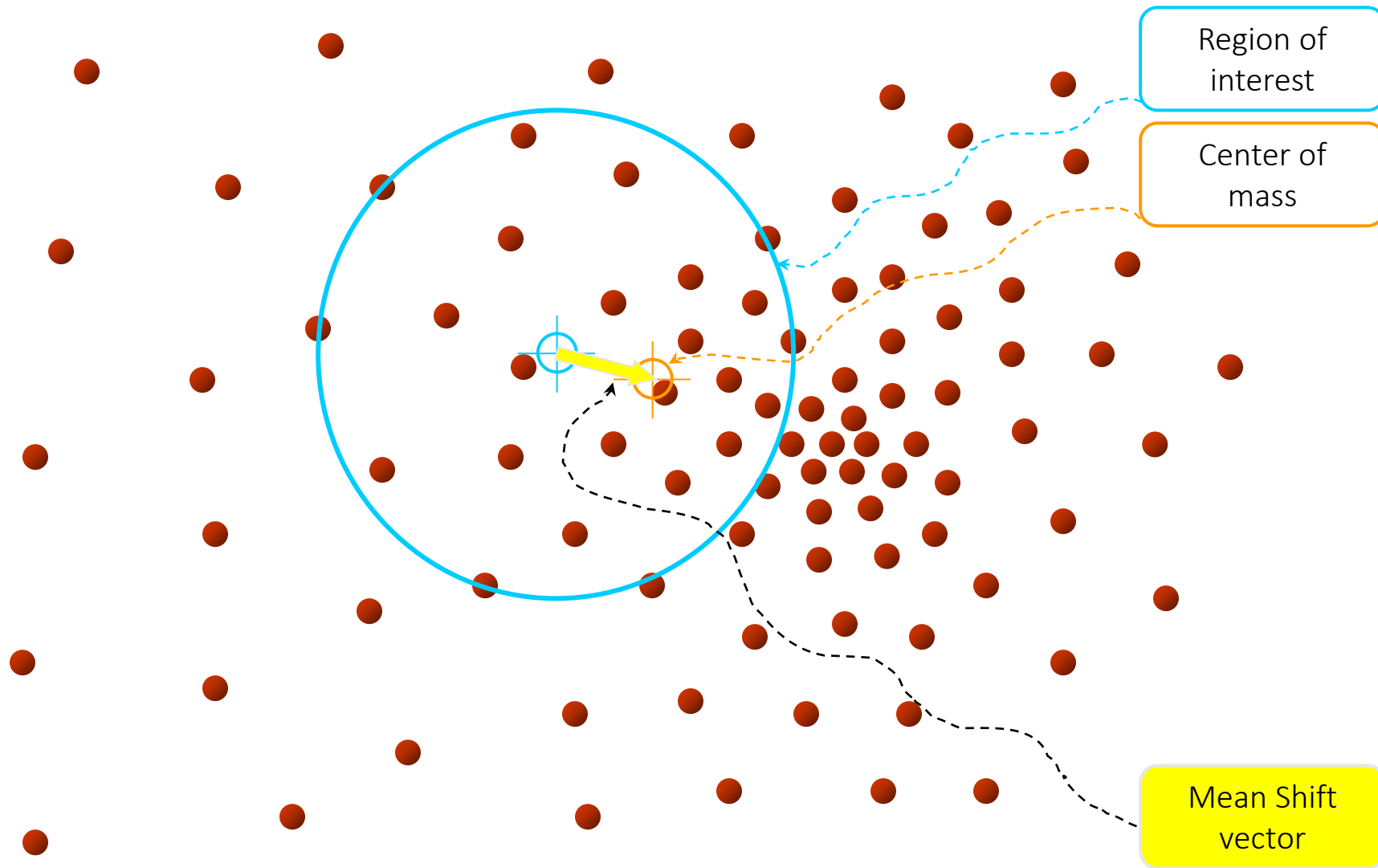
- Iterative Mode Search

1. Initialize random seed, and window  $W$
2. Calculate center of gravity (the “mean”) of  $W$ :
3. Shift the search window to the mean
4. Repeat Step 2 until convergence

$$\sum_{x \in W} xH(x)$$

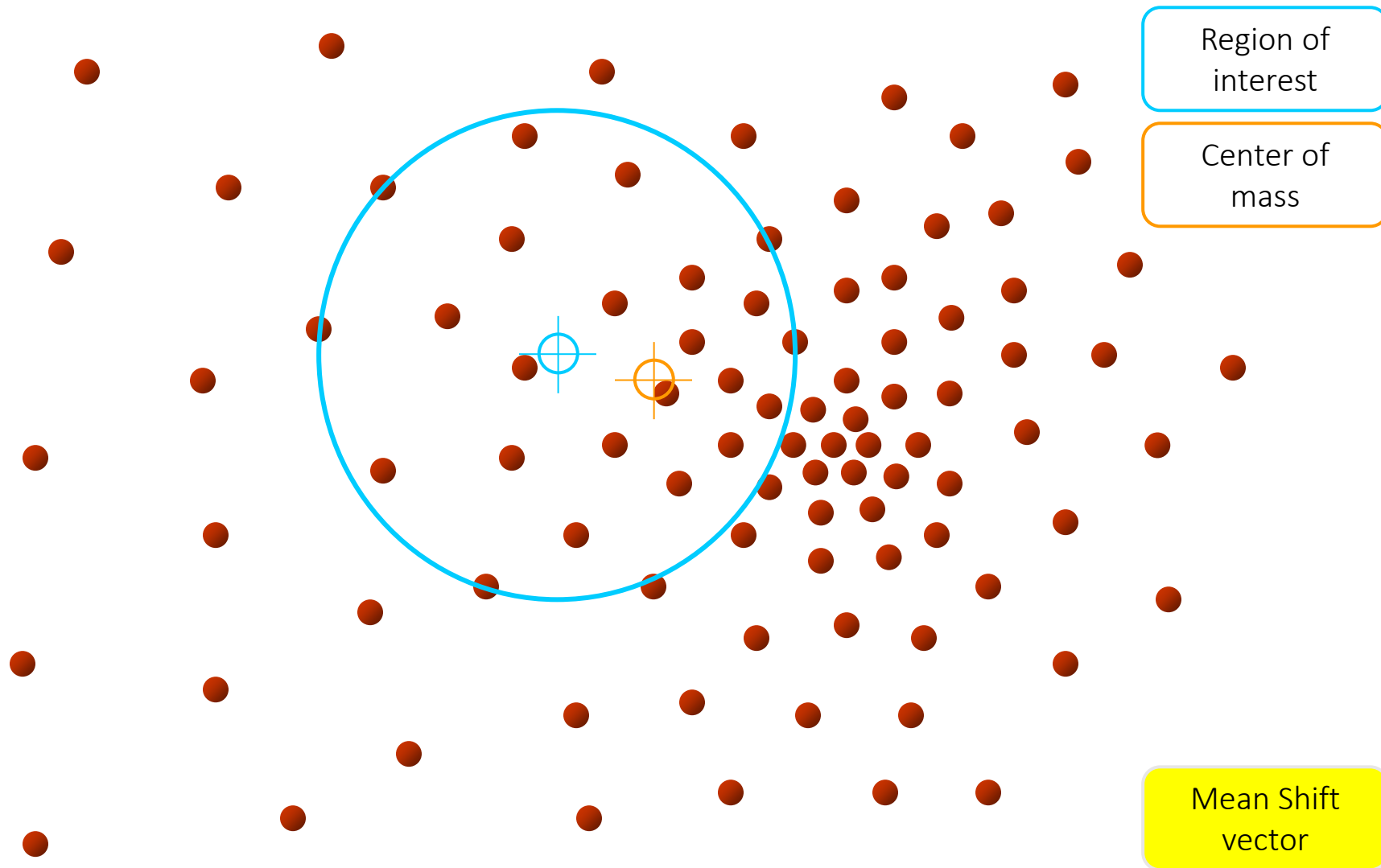
Slide credit: Steve Seitz

## Mean-Shift



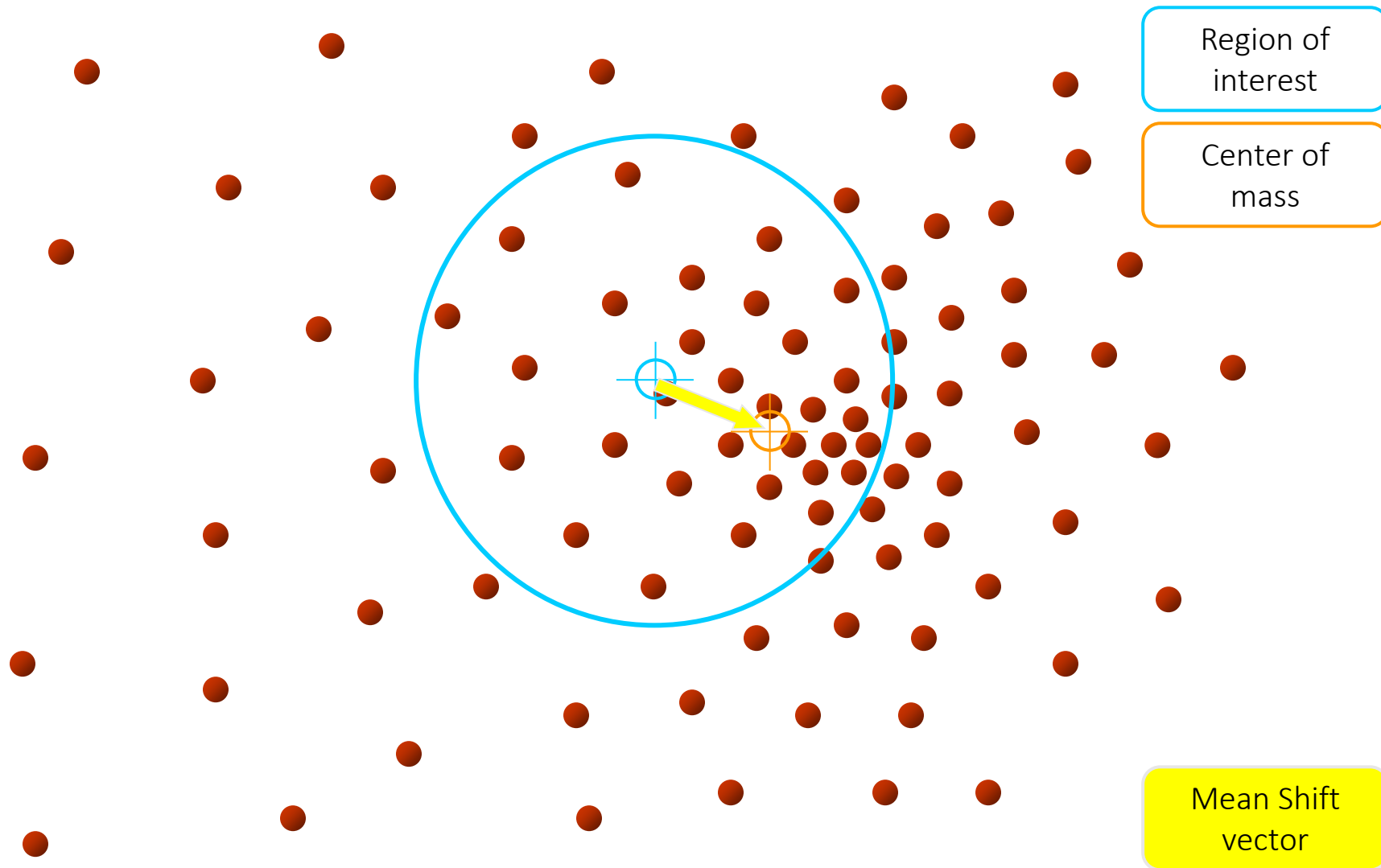
Slide credit: Y. Ukrainitz & B. Sarel

# Mean-Shift



Slide credit: Y. Ukrainitz & B. Sarel

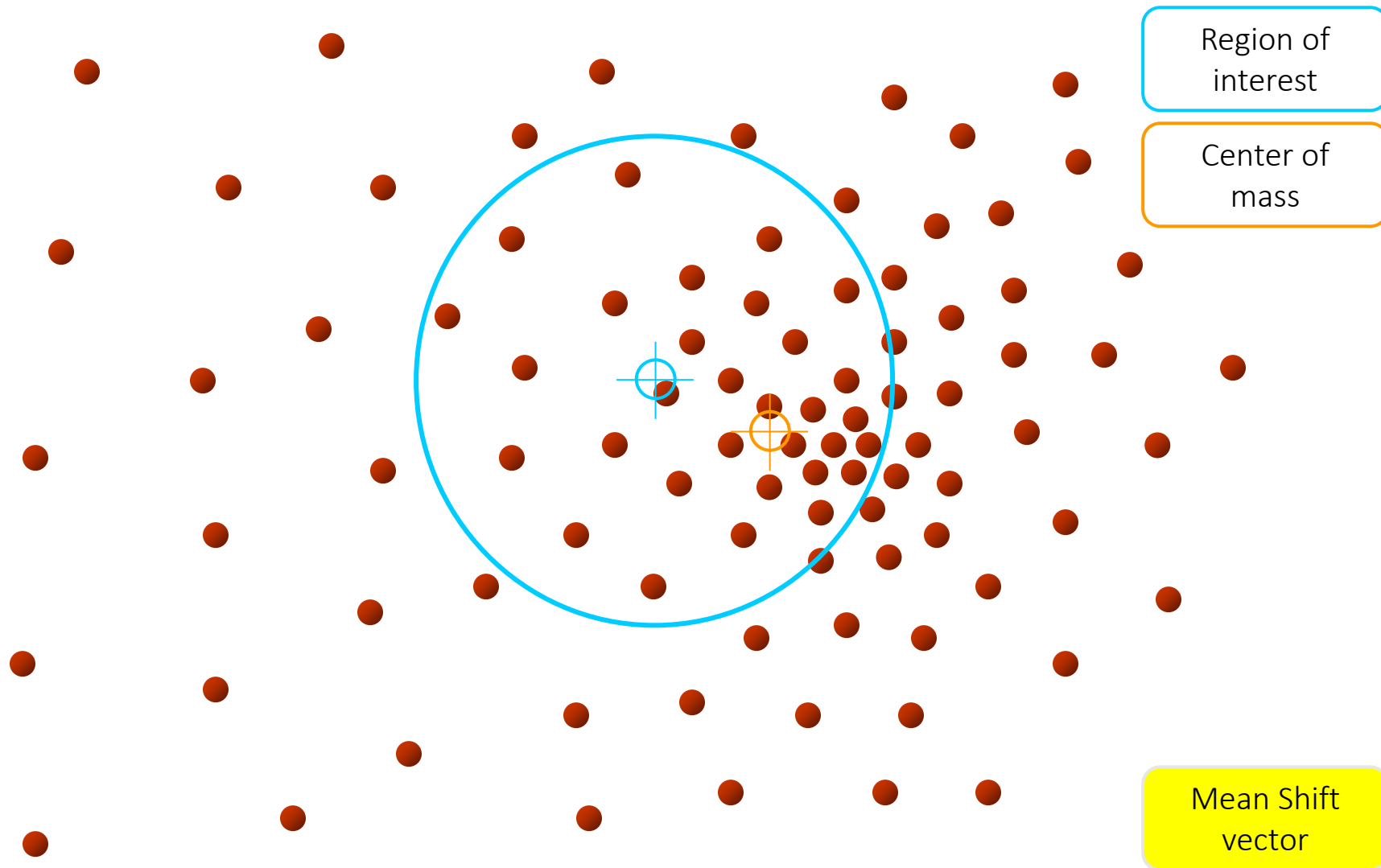
# Mean-Shift



Slide credit: Y. Ukrainitz & B. Sarel

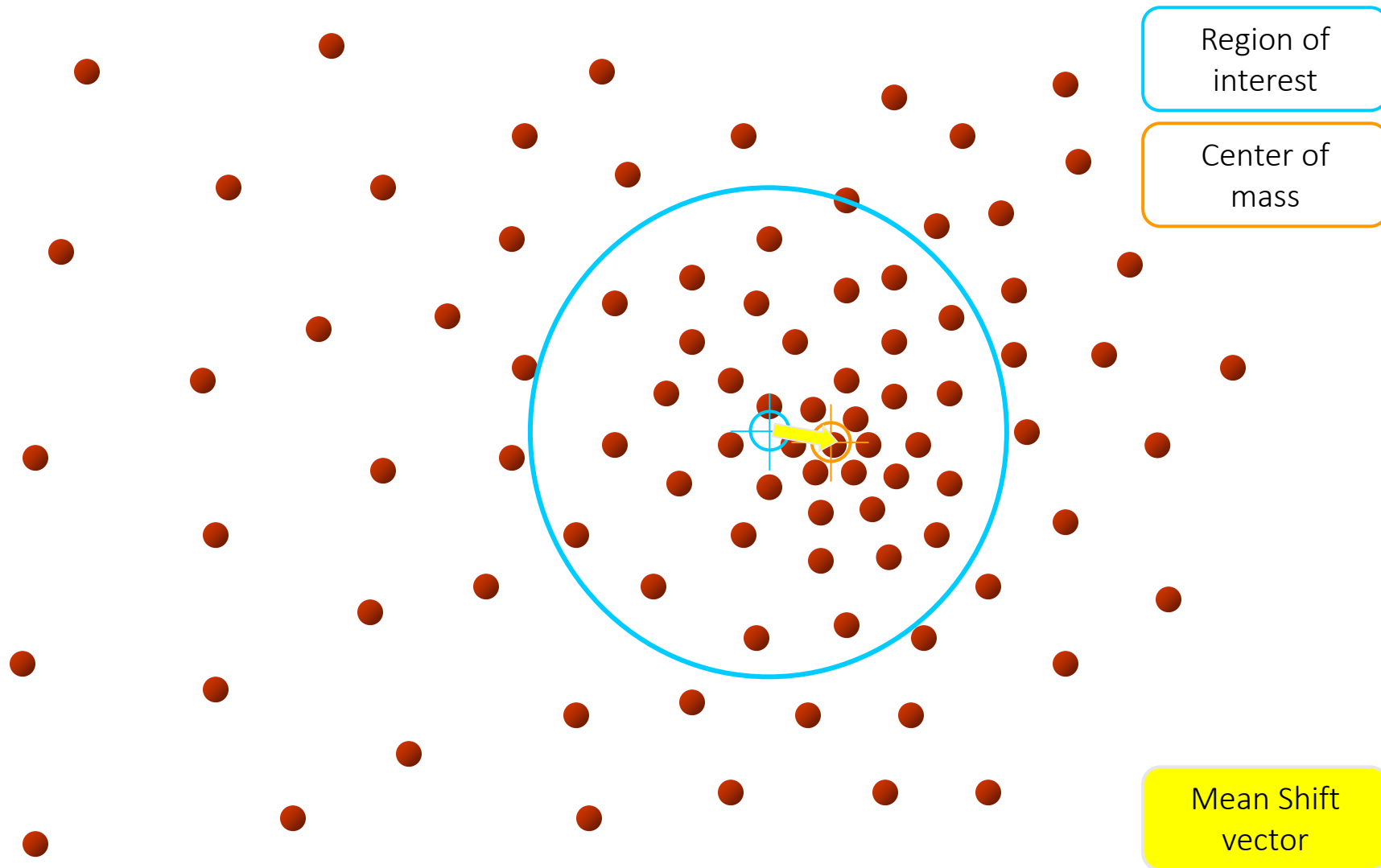


# Mean-Shift



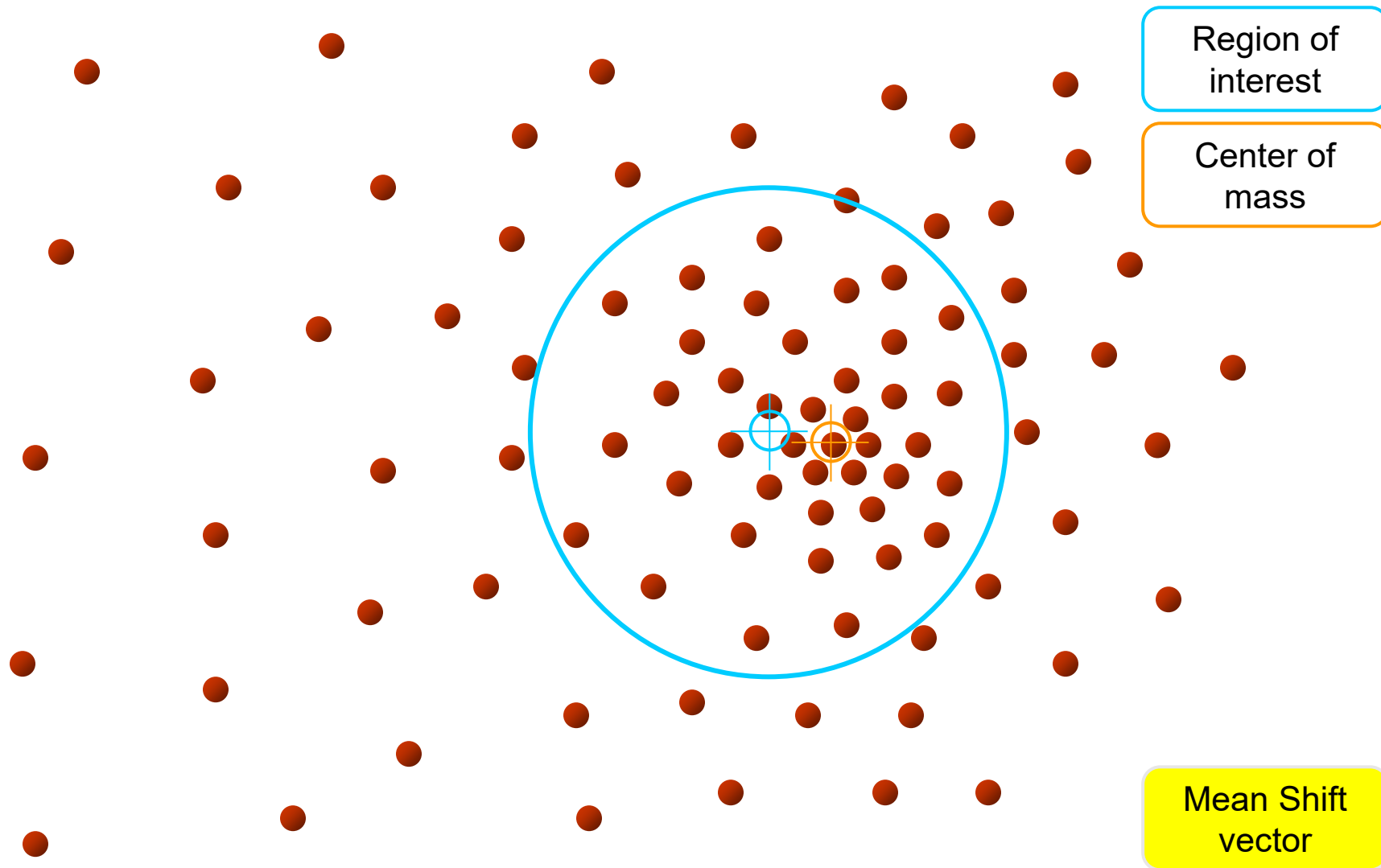
Slide credit: Y. Ukrainitz & B. Sarel

# Mean-Shift



Slide credit: Y. Ukrainitz & B. Sarel

# Mean-Shift



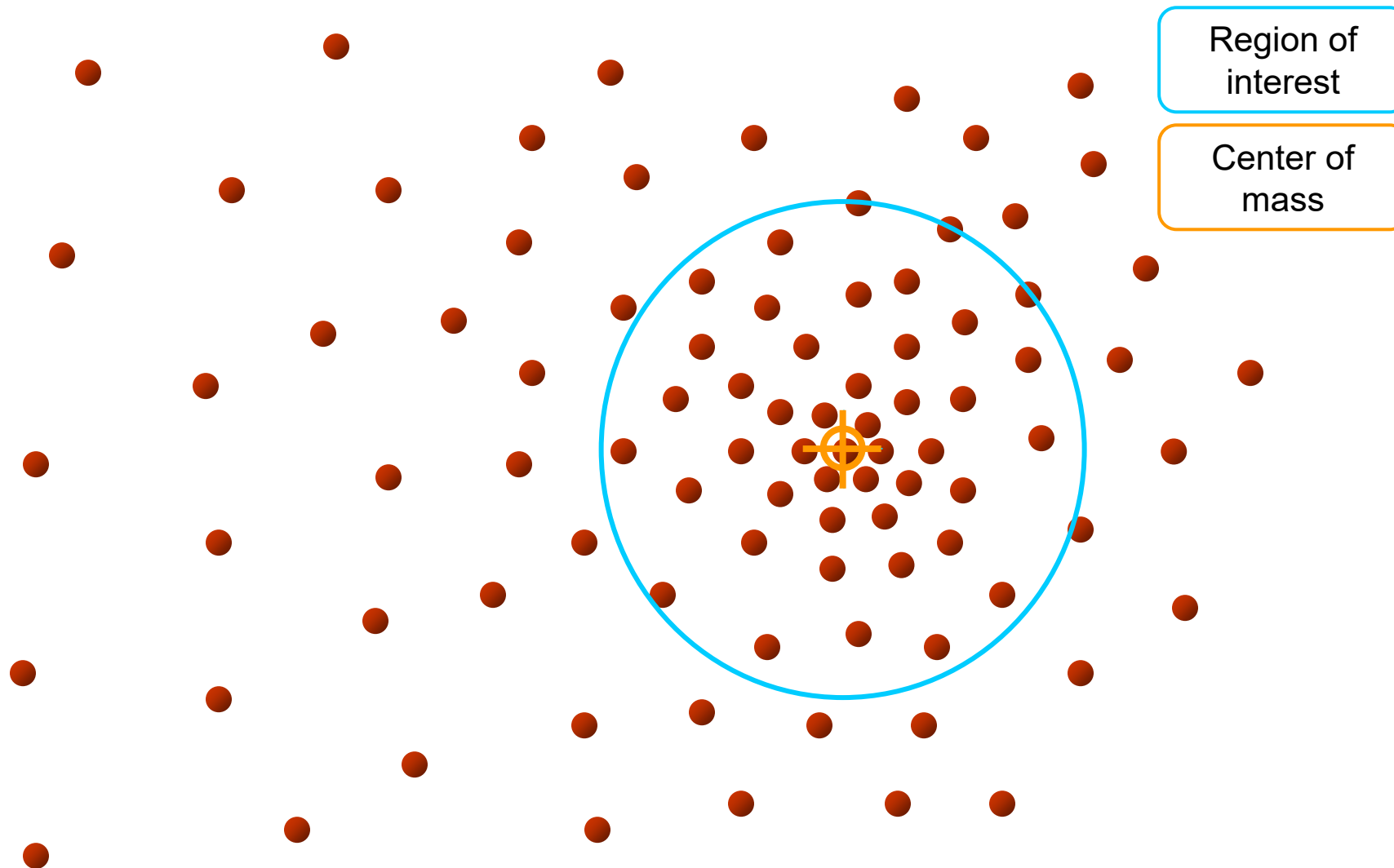
Region of interest

Center of mass

Mean Shift vector

Slide credit: Y. Ukrainitz & B. Sarel

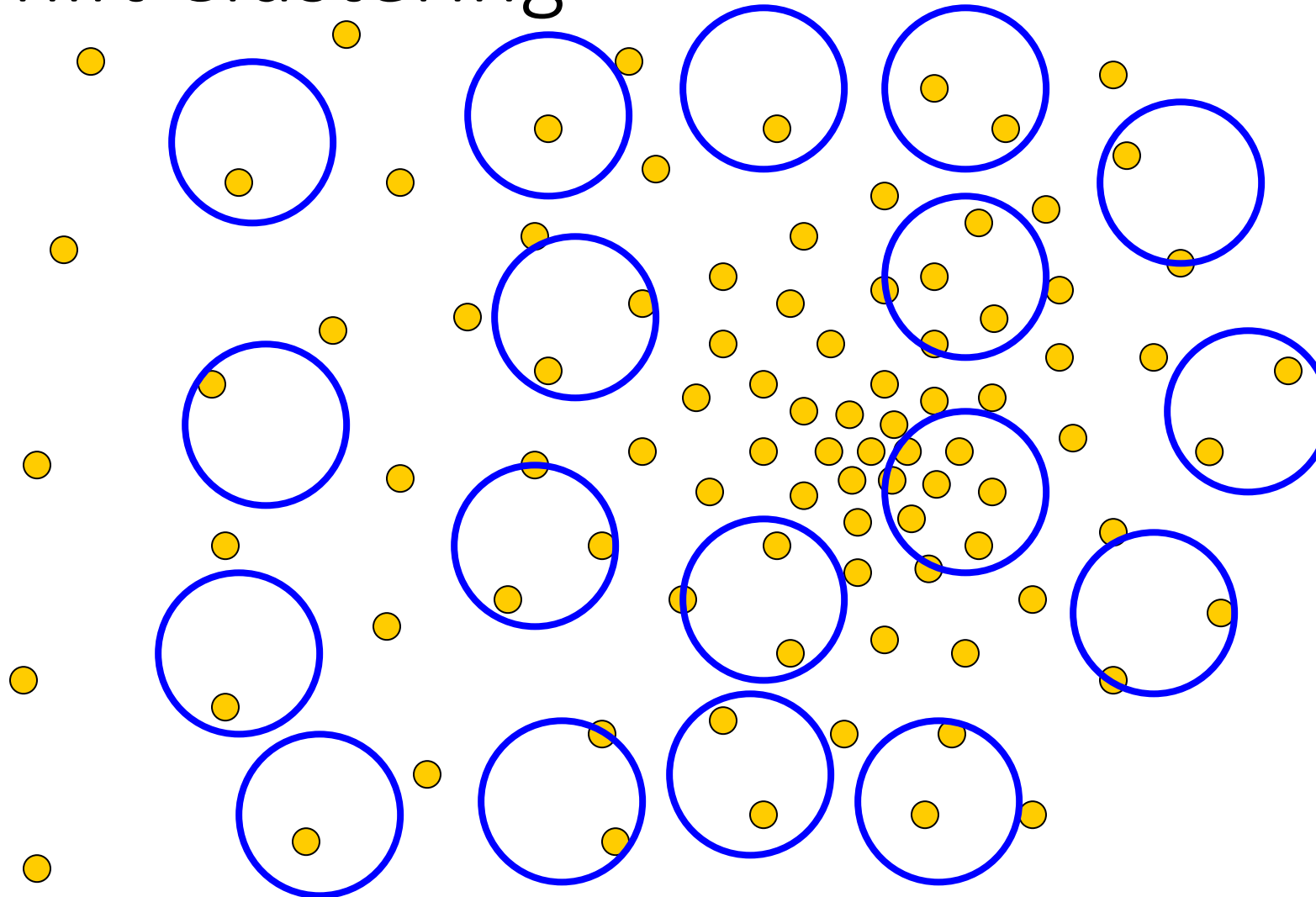
## Mean-Shift



Slide credit: Y. Ukrainitz & B. Sarel



# Mean-Shift Clustering

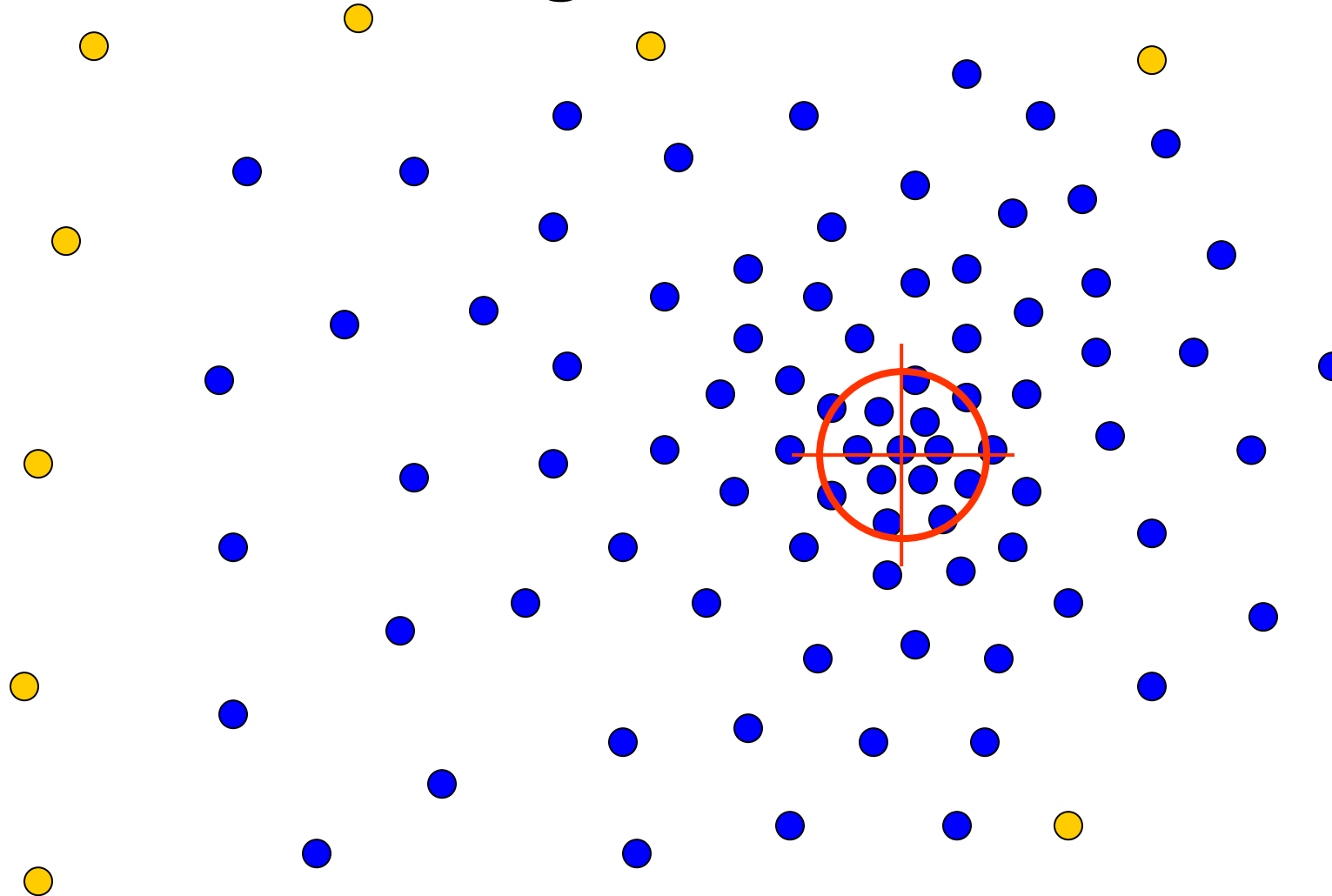


Initialize multiple windows in the space

Run the procedure in parallel

Slide credit: Y. Ukrainitz & B. Sarel

# Mean-Shift Clustering

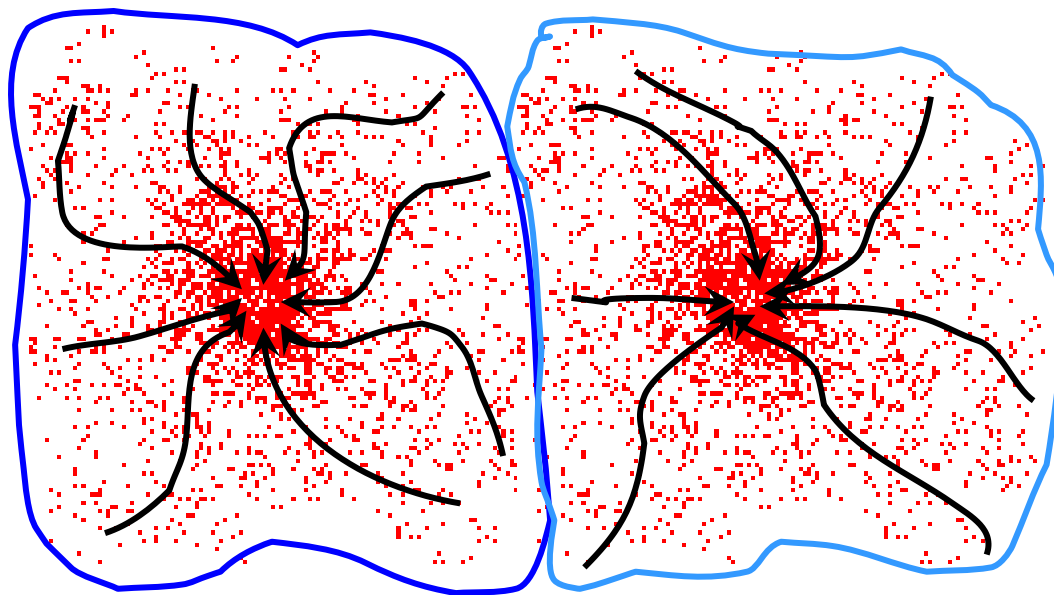


The **blue** data points were traversed by the windows towards the mode.

Slide credit: Y. Ukrainitz & B. Sarel

# Mean-Shift Clustering

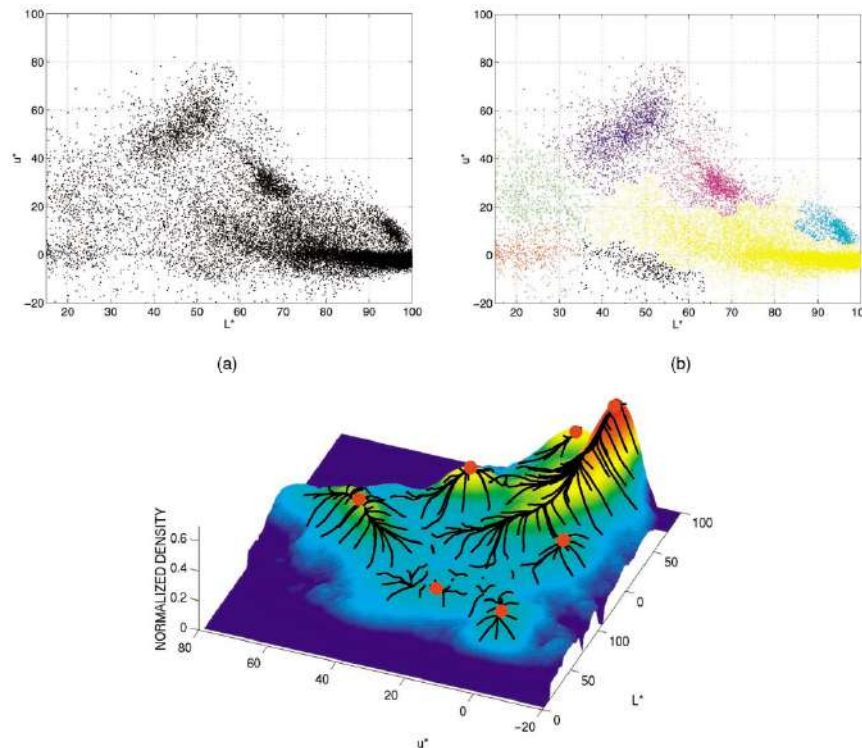
- Cluster: all data points in the attraction basin of a mode
- Attraction basin: the region for which all trajectories lead to the same mode



Slide credit: Y. Ukrainitz & B. Sarel

# Mean-Shift Clustering/Segmentation

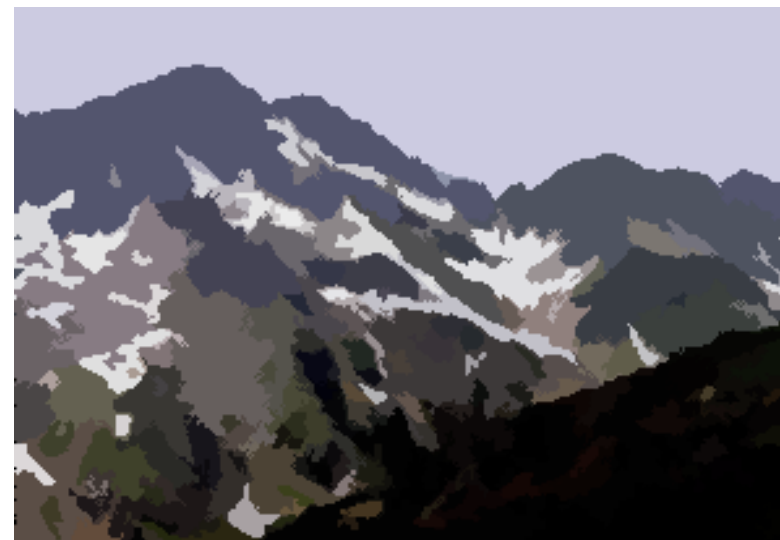
- Find features (color, gradients, texture, etc.)
- Initialize windows at individual pixel locations
- Perform mean shift for each window until convergence
- Merge windows that end up near the same “peak” or mode



Slide credit: Svetlana Lazebnik



## Mean-Shift Segmentation Results



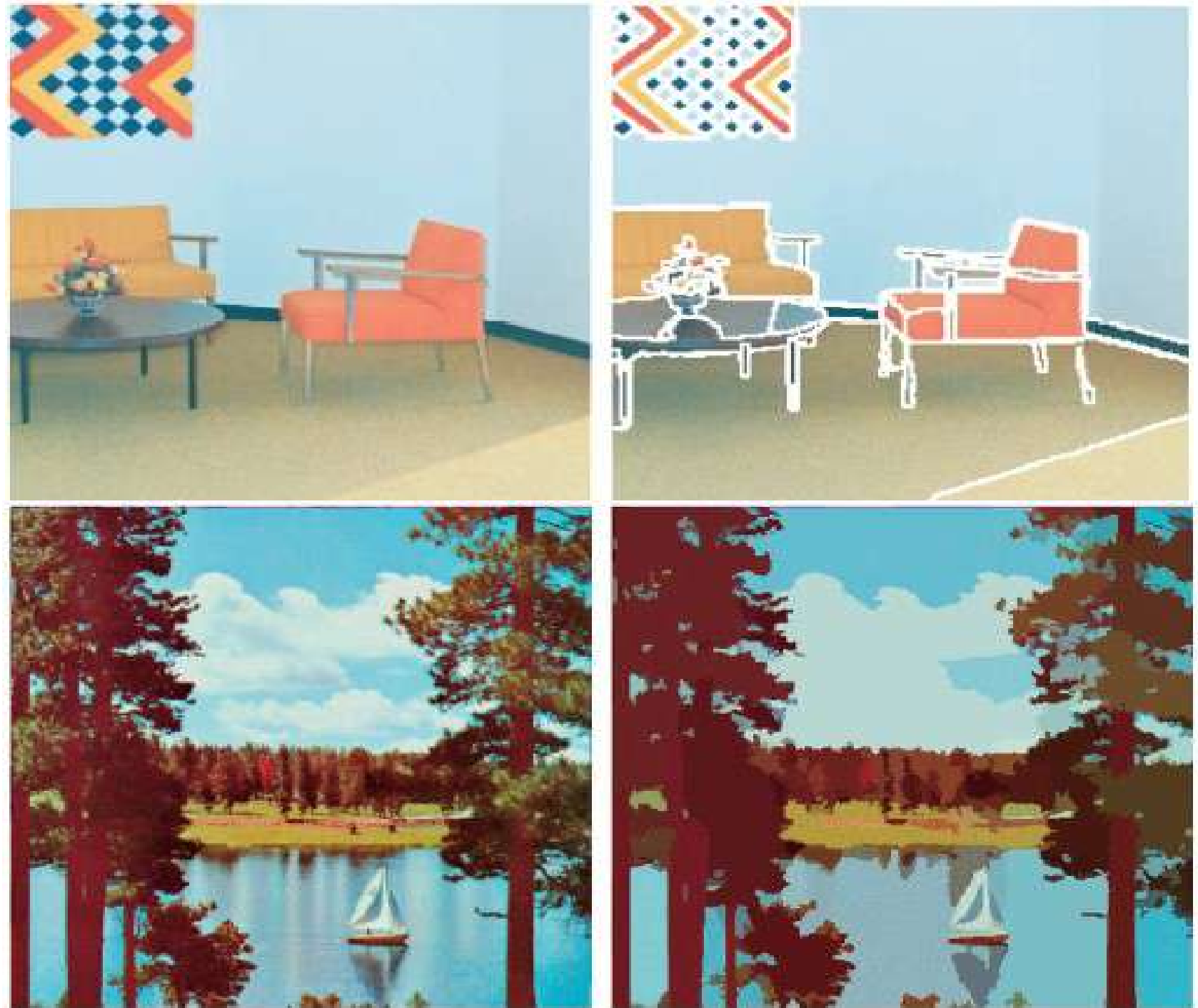
Slide credit: Svetlana Lazebnik

## More Results



Slide credit: Svetlana Lazebnik

## More Results



# Mean-Shift pros and cons

- Pros

- General, application-independent tool
- Model-free, does not assume any prior shape (spherical, elliptical, etc.) on data clusters
- Just a single parameter (window size  $h$ )
  - $h$  has a physical meaning (unlike  $k$ -means)
- Finds variable number of modes
- Robust to outliers

- Cons

- Output depends on window size
- Window size selection is not trivial
- Computationally (relatively) expensive ( $\sim 2s/\text{image}$ )
- Does not scale well with dimensionality of feature space

Slide credit: Svetlana Lazebnik



# Today's Agenda

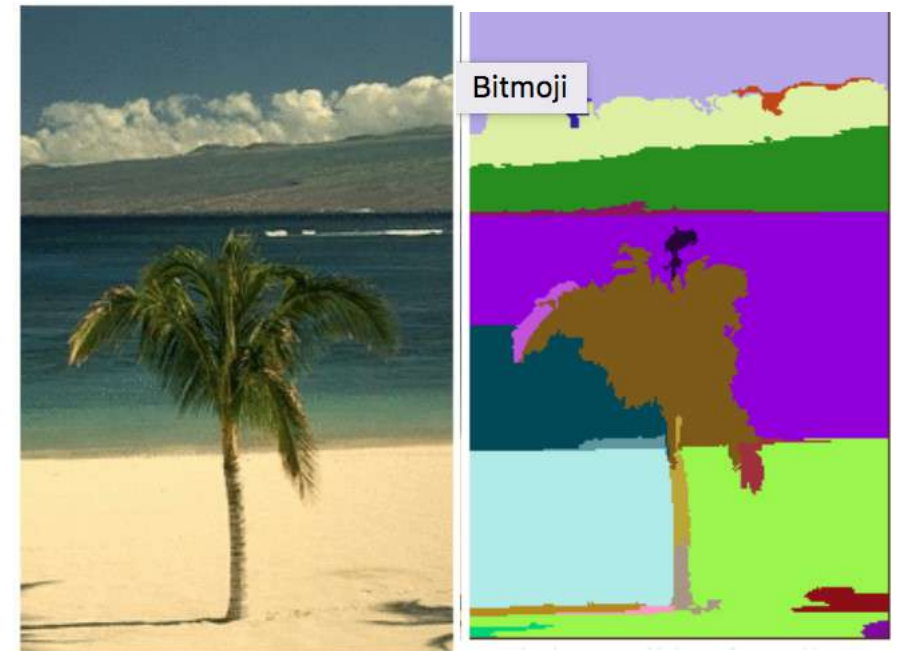
- Visual Recognition Tasks
- Introduction to segmentation and clustering
- Agglomerative clustering
- K-means clustering
- Mean-shift clustering
- Efficient Graph-based image segmentation

**Reading:** Forsyth Chapter 9

D. Comaniciu and P. Meer, [Mean Shift: A Robust Approach toward Feature Space Analysis](#), TPAMI 2002

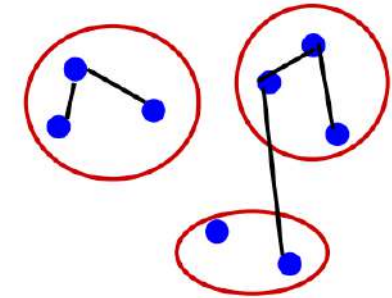
# Efficient Graph-based Image Segmentation

Oversegmentation algorithm introduced by *Felzenszwalb and Huttenlocher* in the paper titled [Efficient Graph-Based Image Segmentation](#)



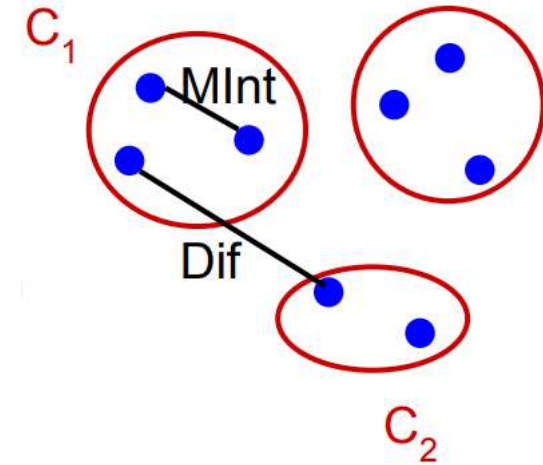
# Problem Formulation

- Graph  $G = (V, E)$
- $V$  is set of nodes (i.e. pixels)
- $E$  is a set of undirected edges between pairs of pixels
- $w(v_i, v_j)$  is the weight of the edge between nodes  $v_i$  and  $v_j$ .
- $S$  is a segmentation of a graph  $G$  such that  $G' = (V, E')$  where  $E' \subset E$ .
- $S$  divides  $G$  into  $G'$  such that it contains distinct clusters  $C$ .



# Predicate for segmentation

- Predicate D determines whether there is a boundary for segmentation.



$$Merge(C_1, C_2) = \begin{cases} True & \text{if } dif(C_1, C_2) < in(C_1, C_2) \\ False & \text{otherwise} \end{cases}$$

Where

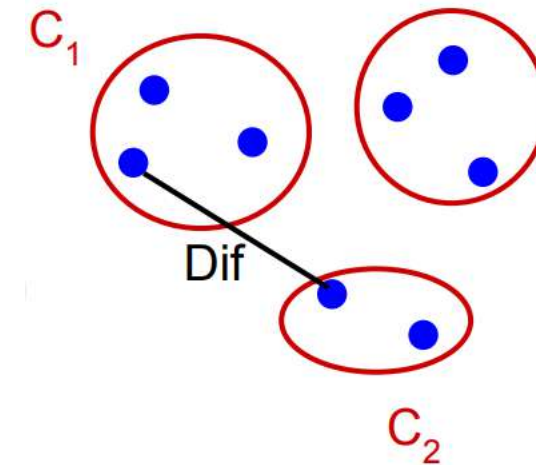
- $dif(C_1, C_2)$  is the difference between two clusters.
- $in(C_1, C_2)$  is the internal difference in the clusters  $C_1$  and  $C_2$





# Predicate for Segmentation

- Predicate D determines whether there is a boundary for segmentation.



$$\text{Merge}(C_1, C_2) = \begin{cases} \text{True} & \text{if } \text{dif}(C_1, C_2) < \text{in}(C_1, C_2) \\ \text{False} & \text{otherwise} \end{cases}$$

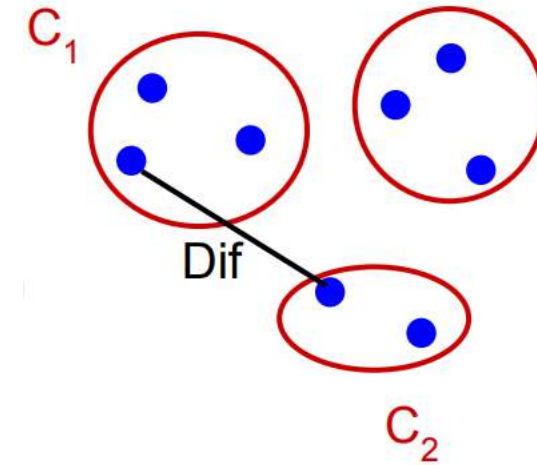
$$\text{dif}(C_1, C_2) = \min_{v_i \in C_1, v_j \in C_2, (C_1, C_2) \in E} w(v_i, v_j)$$

The difference between two components is the minimum weight edge that connects a node  $v_i$  in cluster  $C_1$  to node  $v_j$  in  $C_2$



# Predicate for Segmentation

- Predicate D determines whether there is a boundary for segmentation.



$$Merge(C_1, C_2) = \begin{cases} True & \text{if } dif(C_1, C_2) < in(C_1, C_2) \\ False & \text{otherwise} \end{cases}$$

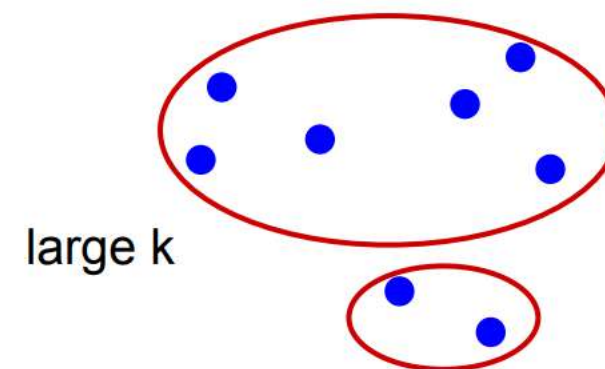
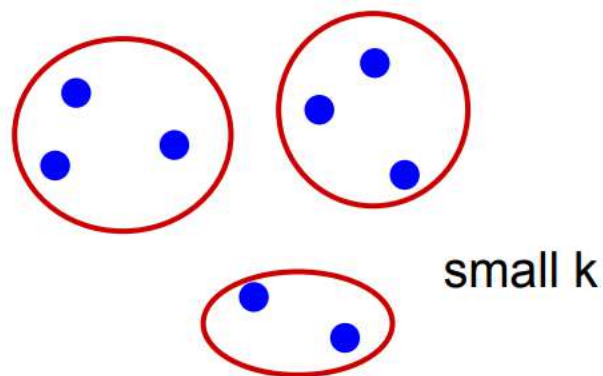
$$dif(C_1, C_2) = \min_{v_i \in C_1, v_j \in C_2, (C_1, C_2) \in E} w(v_i, v_j)$$

$$in(C_1, C_2) = \min_{C \in \{C_1, C_2\}} \left[ \max_{v_i, v_j \in C} \left[ w(v_i, v_j) + \frac{k}{|C|} \right] \right]$$

$in(C_1, C_2)$  is the maximum weight edge that connects two nodes in the same component.

# Predicate for Segmentation

- $k/|C|$  sets the threshold by which the components need to be different from the internal nodes in a component.
- Properties of constant  $k$ :
  - If  $k$  is large, it causes a preference for larger components.
  - $k$  does not set a minimum size for components.



# Algorithm for Segmentation

The input is a graph  $G = (V, E)$ , with  $n$  vertices and  $m$  edges. The output is a segmentation of  $V$  into components  $S = (C_1, \dots, C_r)$ .

0. Sort  $E$  into  $\pi = (o_1, \dots, o_m)$ , by non-decreasing edge weight.
1. Start with a segmentation  $S^0$ , where each vertex  $v_i$  is in its own component.
2. Repeat step 3 for  $q = 1, \dots, m$ .
3. Construct  $S^q$  given  $S^{q-1}$  as follows. Let  $v_i$  and  $v_j$  denote the vertices connected by the  $q$ -th edge in the ordering, i.e.,  $o_q = (v_i, v_j)$ . If  $v_i$  and  $v_j$  are in disjoint components of  $S^{q-1}$  and  $w(o_q)$  is small compared to the internal difference of both those components, then merge the two components otherwise do nothing. More formally, let  $C_i^{q-1}$  be the component of  $S^{q-1}$  containing  $v_i$  and  $C_j^{q-1}$  the component containing  $v_j$ . If  $C_i^{q-1} \neq C_j^{q-1}$  and  $w(o_q) \leq MInt(C_i^{q-1}, C_j^{q-1})$  then  $S^q$  is obtained from  $S^{q-1}$  by merging  $C_i^{q-1}$  and  $C_j^{q-1}$ . Otherwise  $S^q = S^{q-1}$ .
4. Return  $S = S^m$ .



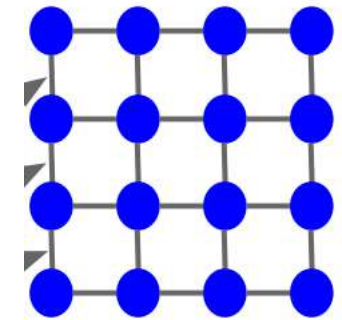
# Features and weights

How to build the graph ? Two options:

1. Grid-graph: Every pixel is connected to its 8 neighboring pixels and the weights are determined by the difference in intensities.

2. NN-graph: Project every pixel into **feature space** defined by

- $(x, y, r, g, b)$ .
- Weights between pixels are determined using L2 (Euclidian) distance in feature space.
- Edges are chosen for only top ten nearest neighbors in feature space to ensure run time of  $O(n \log n)$  where  $n$  is number of pixels.



## Results

With 8-neighbor grid graph  
Edge weight: intensity difference

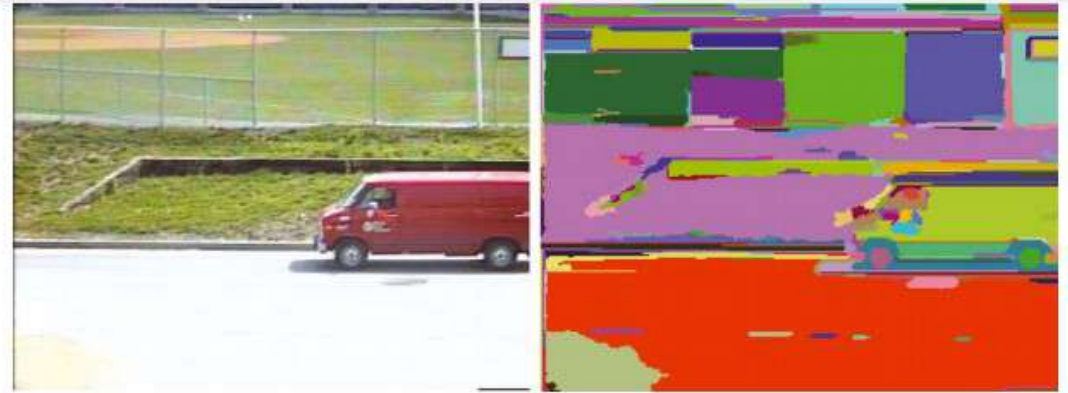


Figure 2. A street scene (320 × 240 color image), and the segmentation results produced by our algorithm ( $\sigma = 0.8, k = 300$ ).



Figure 3. A baseball scene (432 × 294 grey image), and the segmentation results produced by our algorithm ( $\sigma = 0.8, k = 300$ ).



Figure 4. An indoor scene (image 320 × 240, color), and the segmentation results produced by our algorithm ( $\sigma = 0.8, k = 300$ ).

# Results

With nearest neighbor graph  
 Edge weight: L2 distance in feature space

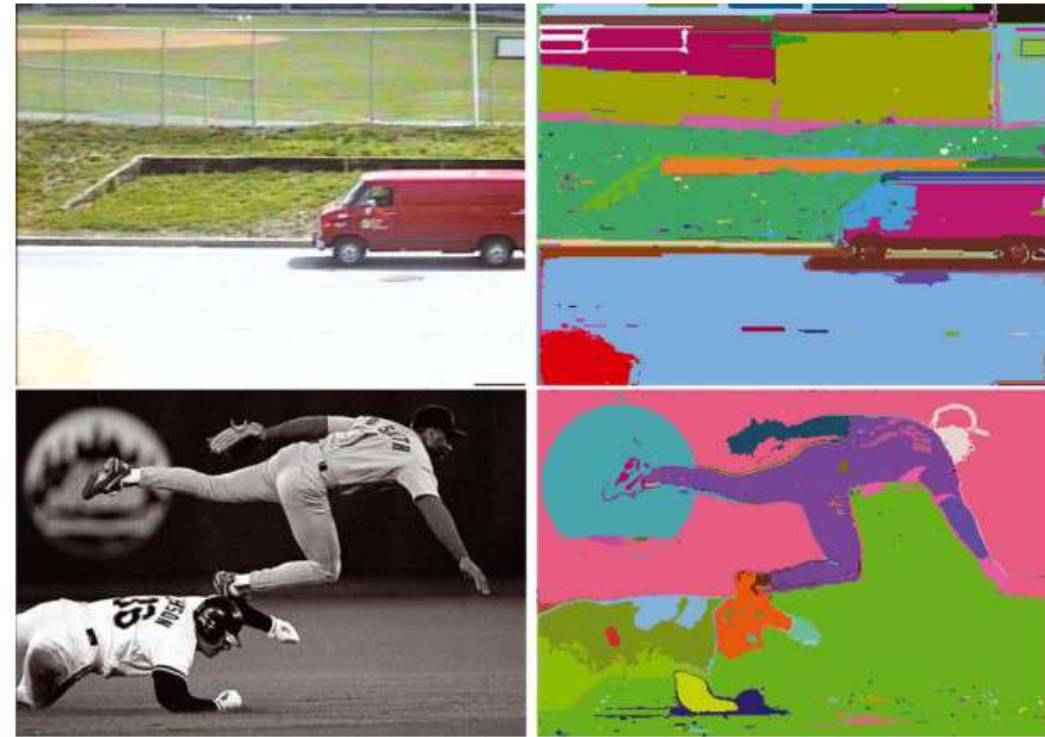


Figure 7. Segmentation of the street and baseball player scenes from the previous section, using the nearest neighbor graph rather than the grid graph ( $\sigma = 0.8, k = 300$ ).



Figure 8. Segmentation using the nearest neighbor graph can capture spatially non-local regions ( $\sigma = 0.8, k = 300$ ).



## Results – close up





**MAI4CAREU**

Master programmes in Artificial  
Intelligence 4 Careers in Europe



# Thank you.





University of Cyprus – MSc Artificial Intelligence

# MAI644 – COMPUTER VISION

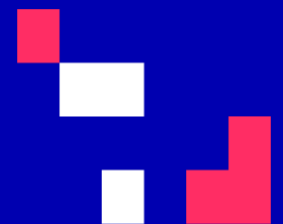
## Lecture 11: Visual Recognition – Image Classification

**Melinos Averkiou**

CYENS Centre of Excellence

University of Cyprus - Department of Computer Science

[m.averkiou@cyens.org.cy](mailto:m.averkiou@cyens.org.cy)



# Last time

- Visual Recognition Tasks
- Introduction to segmentation and clustering
- Agglomerative clustering
- K-means clustering
- Mean-shift clustering
- Efficient Graph-based image segmentation

# Today's Agenda

- A simple Image Classification pipeline
  - Classification overview
- K-nearest neighbor algorithm
  - kNN: algorithm
  - kNN: analysis

[material based on Niebles-Krishna]





# Today's Agenda

- A simple Image Classification pipeline
  - Classification overview
- K-nearest neighbor algorithm
  - kNN: algorithm
  - kNN: analysis



# Visual recognition: a classification framework

- Apply a prediction function to a feature representation of the image to get the desired output:

$$f(\text{apple image}) = \text{"apple"}$$

$$f(\text{tomato image}) = \text{"tomato"}$$

$$f(\text{cow image}) = \text{"cow"}$$

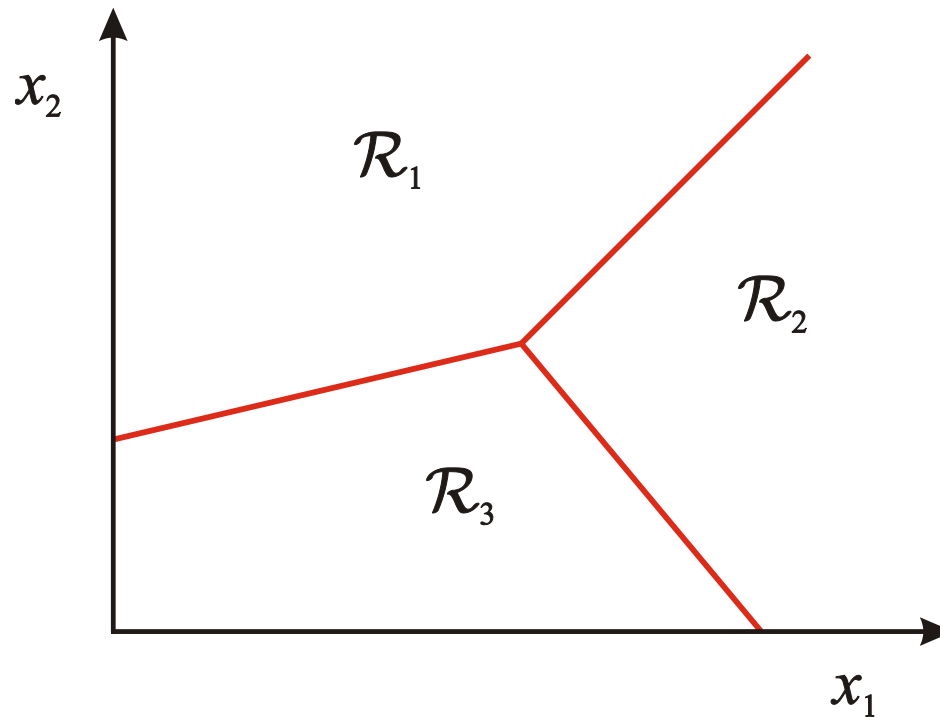
Dataset: ETH-80, by B. Leibe

Slide credit: L. Lazebnik



# Classification

- Assign input feature vector to one of two or more classes
- Any decision rule divides input space into *decision regions* separated by *decision boundaries*



Slide credit: L. Lazebnik



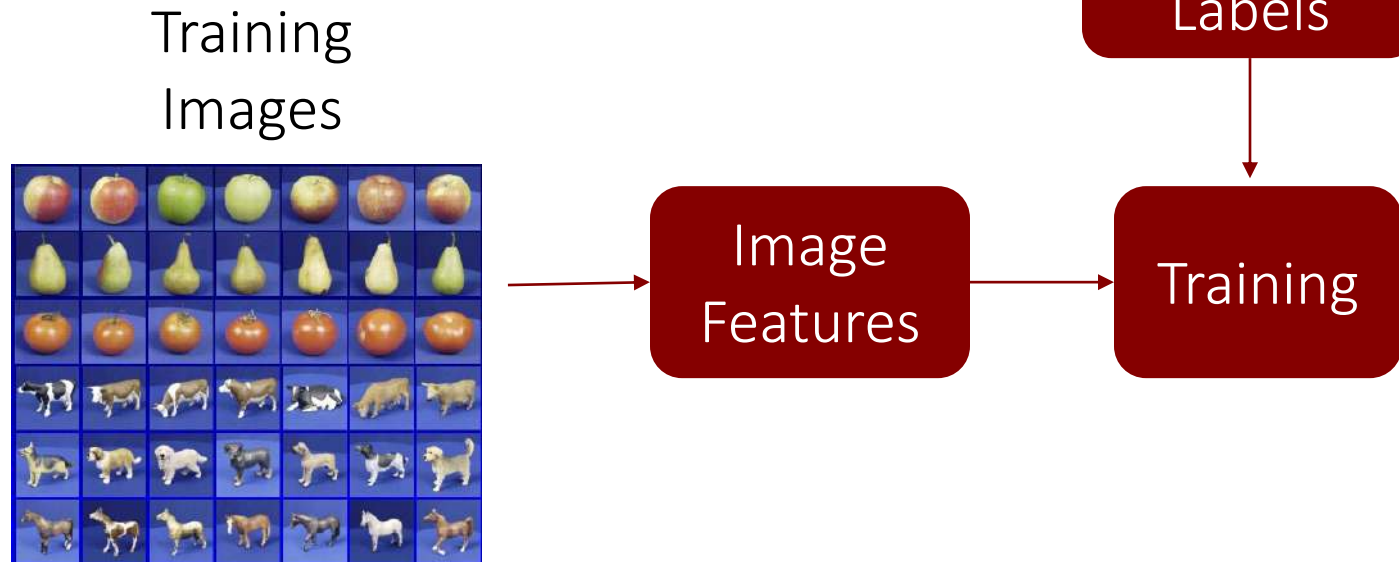
# A simple pipeline - Training

Training Images

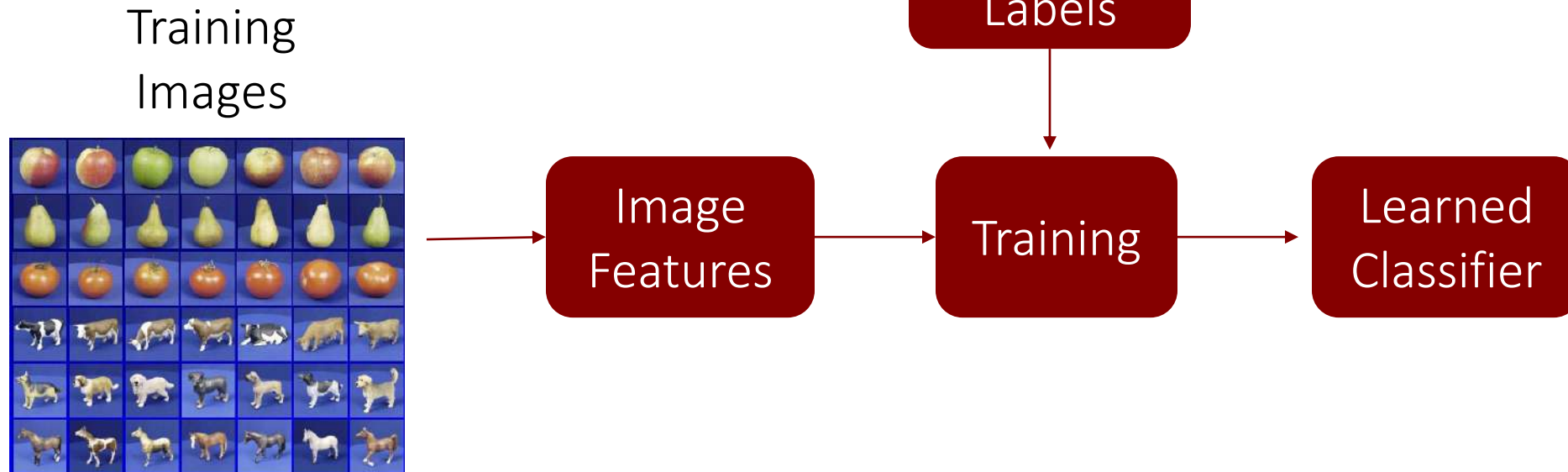


Image Features

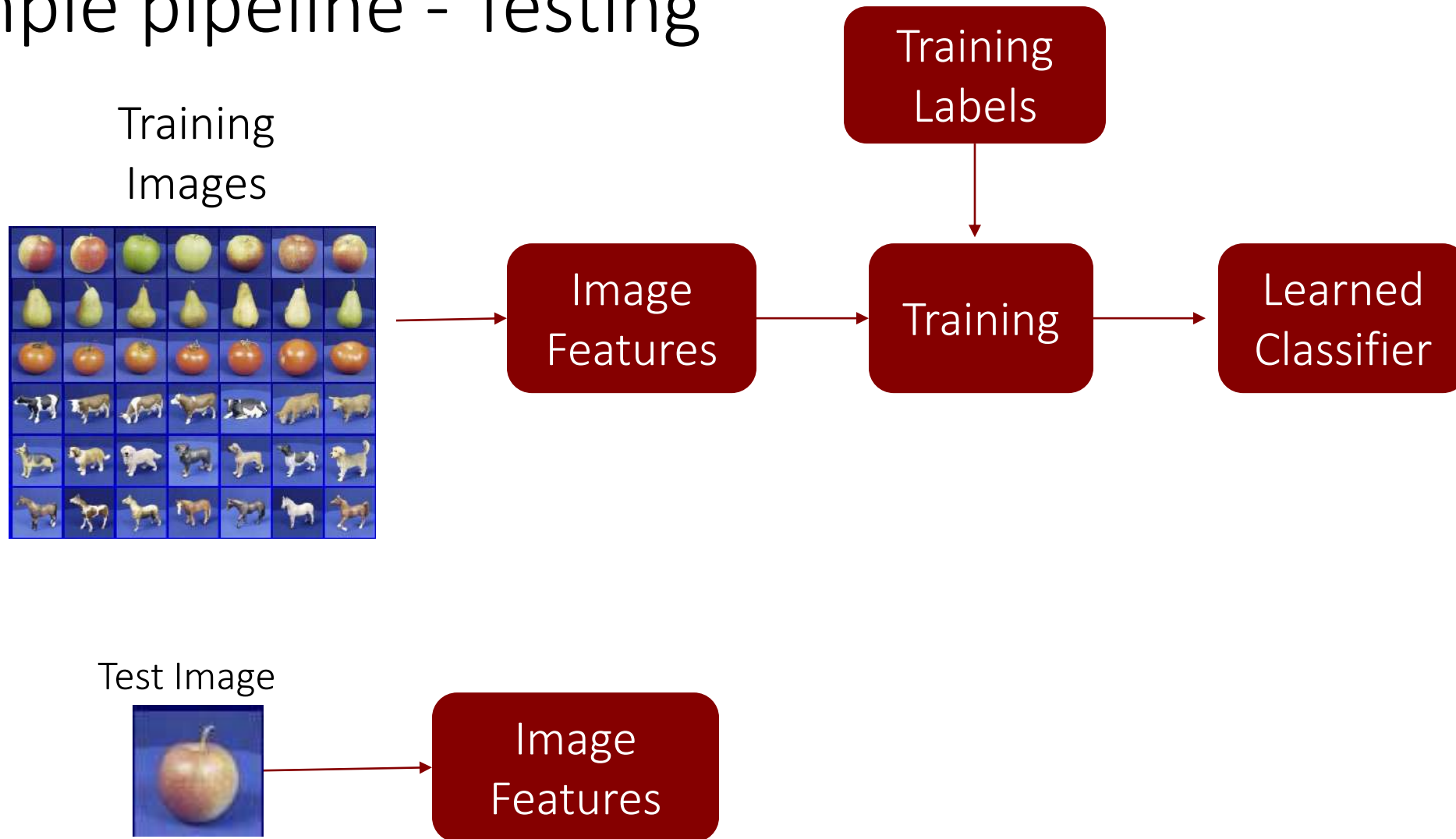
# A simple pipeline - Training



# A simple pipeline - Training

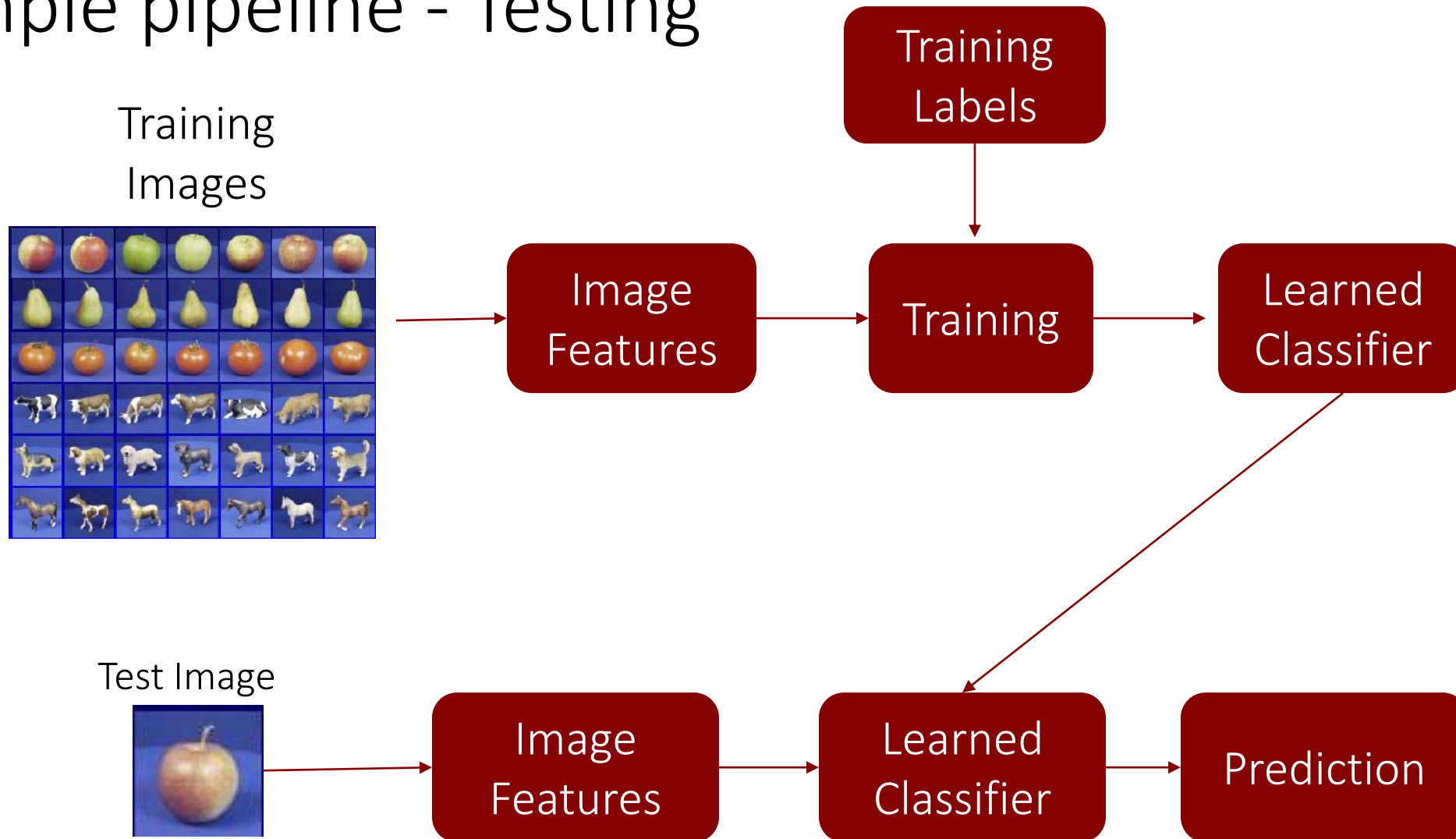


# A simple pipeline - Testing

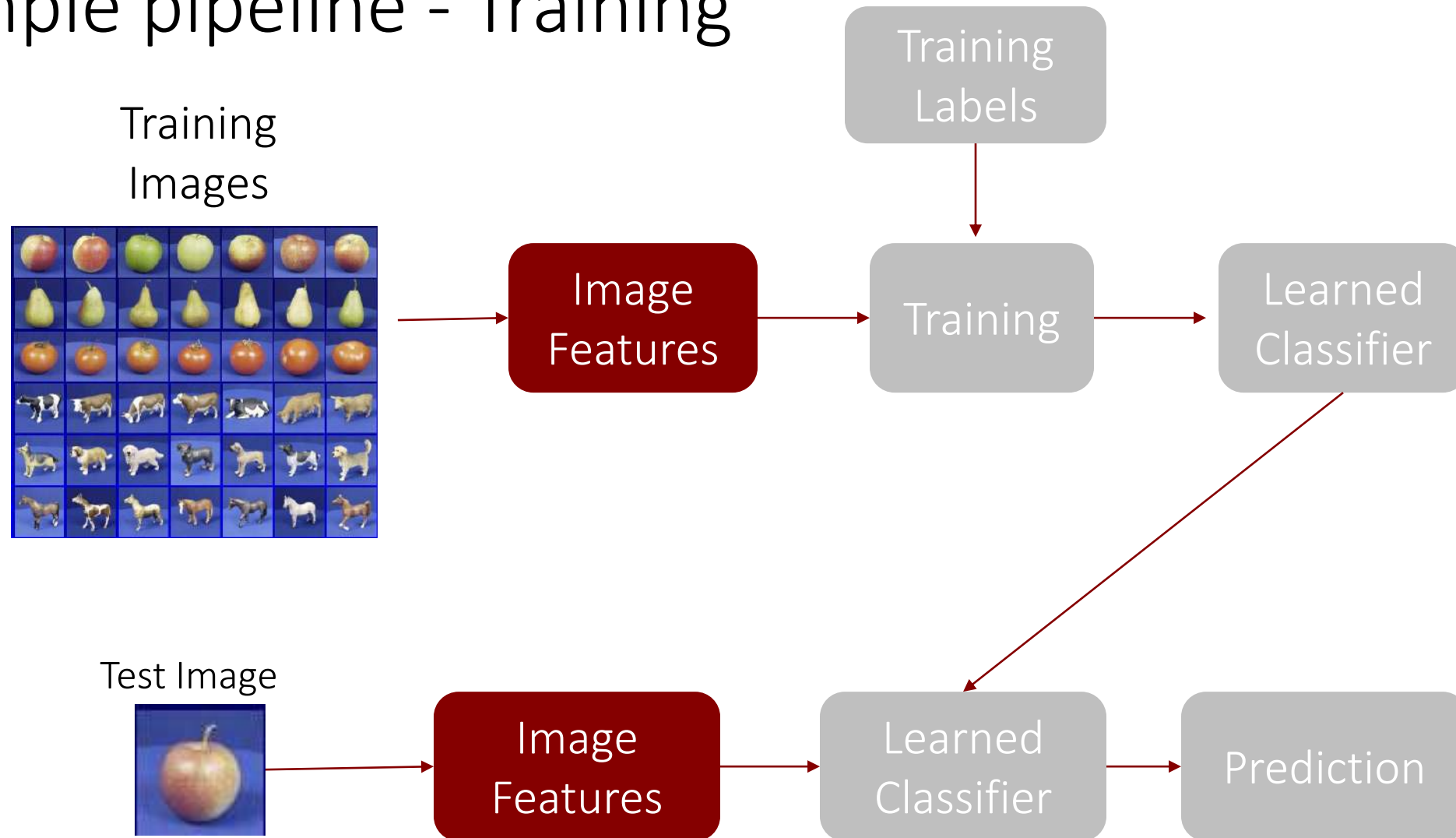




## A simple pipeline - Testing



# A simple pipeline - Training



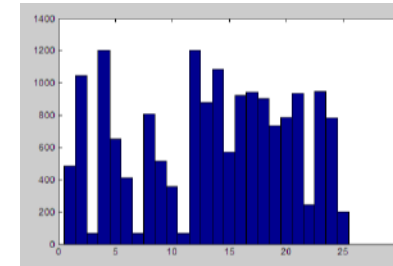
# Image Features

Input image



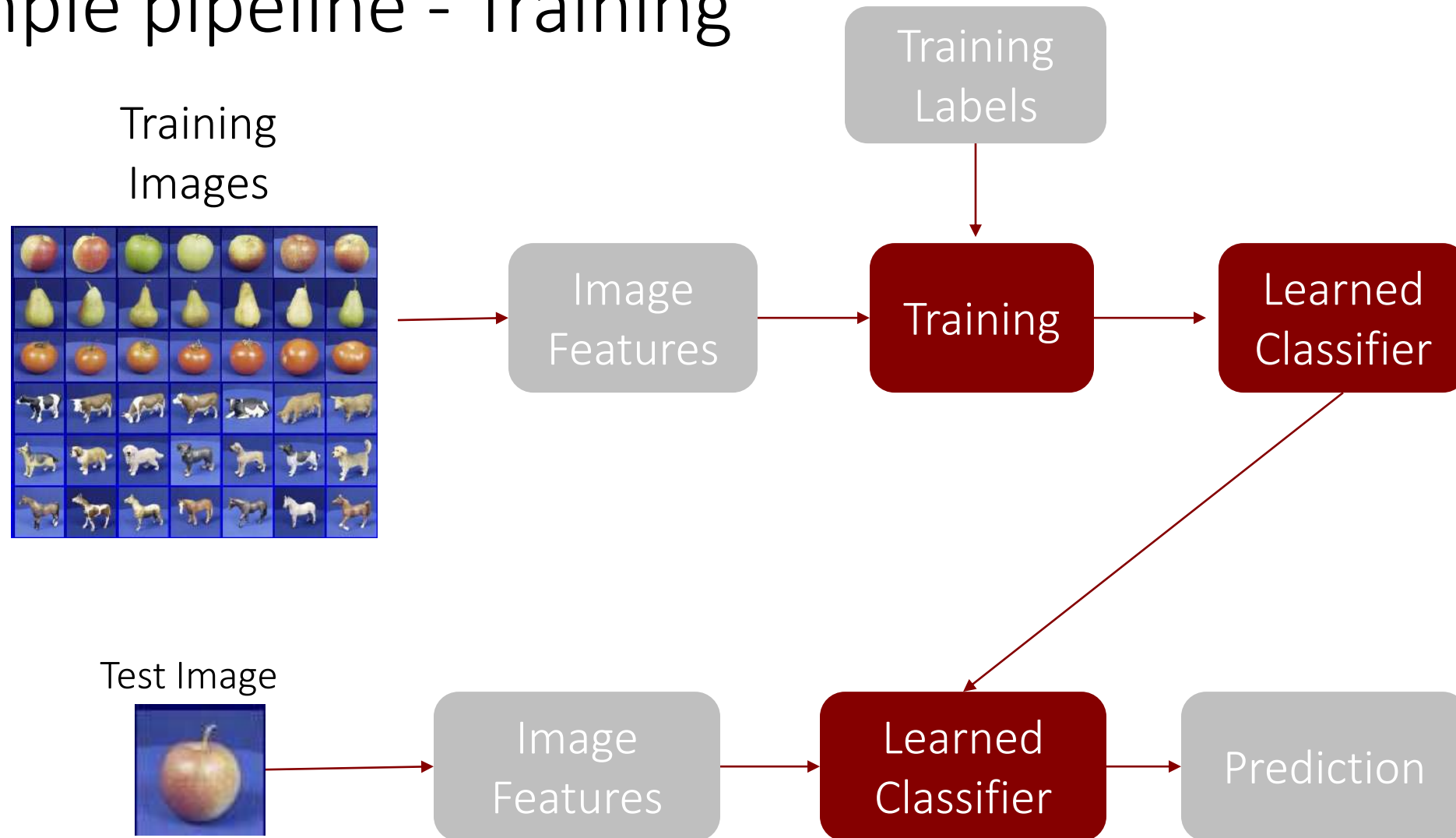
Many features to choose from

- Histogram of Color
- Histogram of Gradients
- SIFT
- Bag of words
- etc...



[Image From ETH-80 dataset: Analyzing Appearance and Contour Based Methods for Object Categorization - by Leibe et al 2003](#)

# A simple pipeline - Training





# Many classifiers to choose from

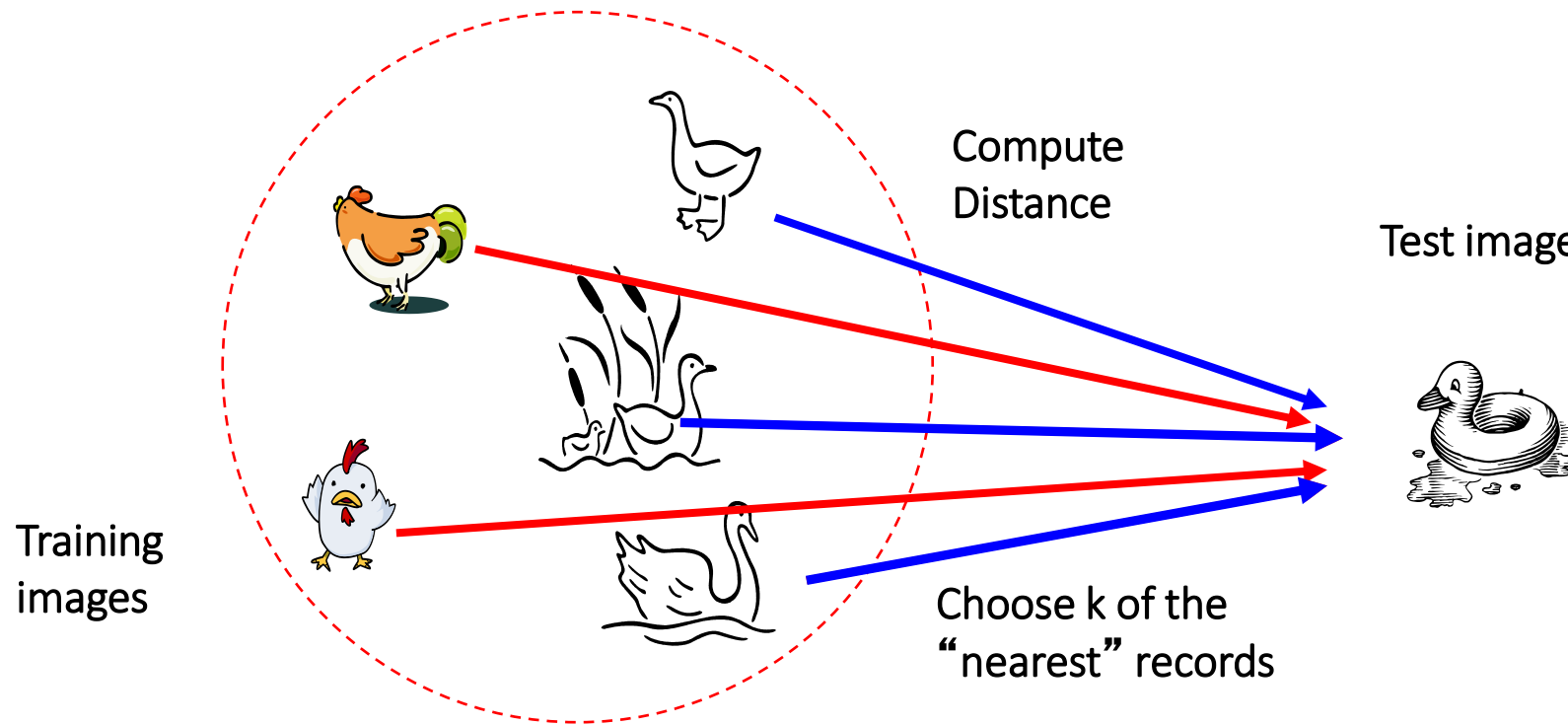
- **K-nearest neighbor**
- SVM
- Neural networks
- Naïve Bayes
- Bayesian network
- Logistic regression
- Randomized Forests
- Boosted Decision Trees
- Restricted Boltzmann Machines
- Etc.

Which is the best one?

Slide credit: D. Hoiem

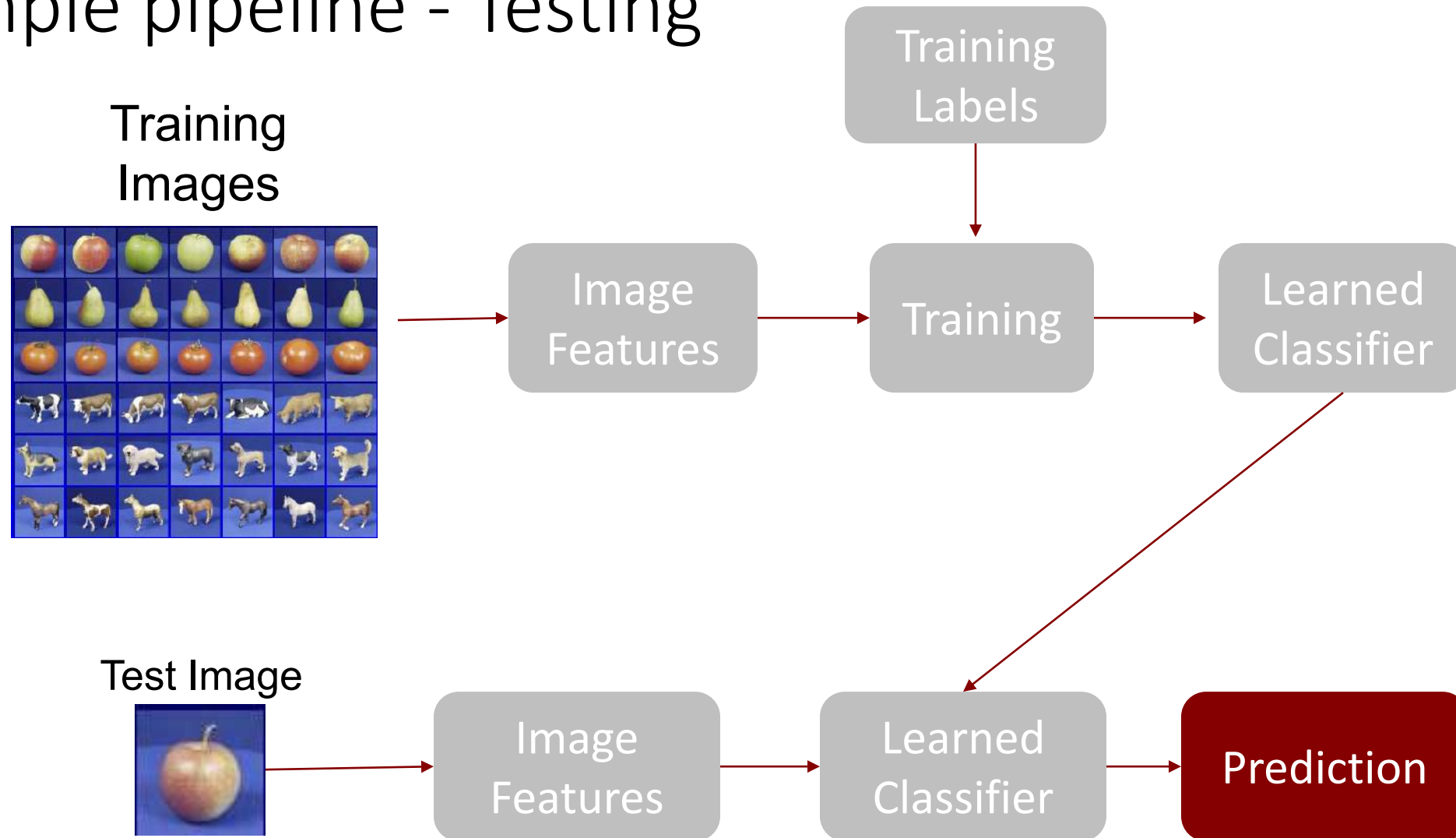
# Classifiers: Nearest neighbor

Assign label of nearest training data point to each test data point

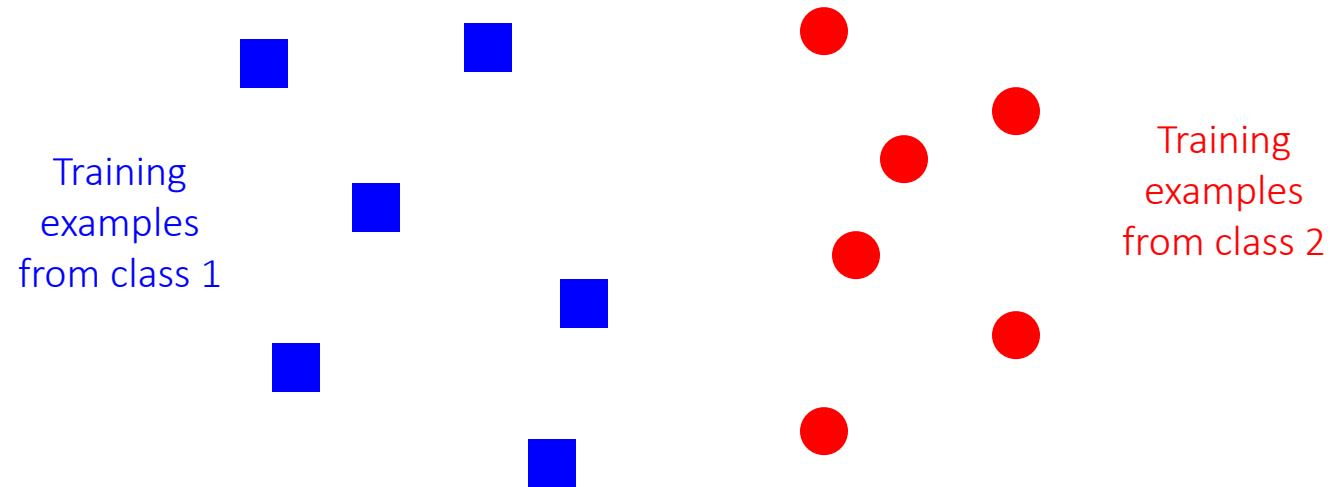


Source: N. Goyal

## A simple pipeline - Testing



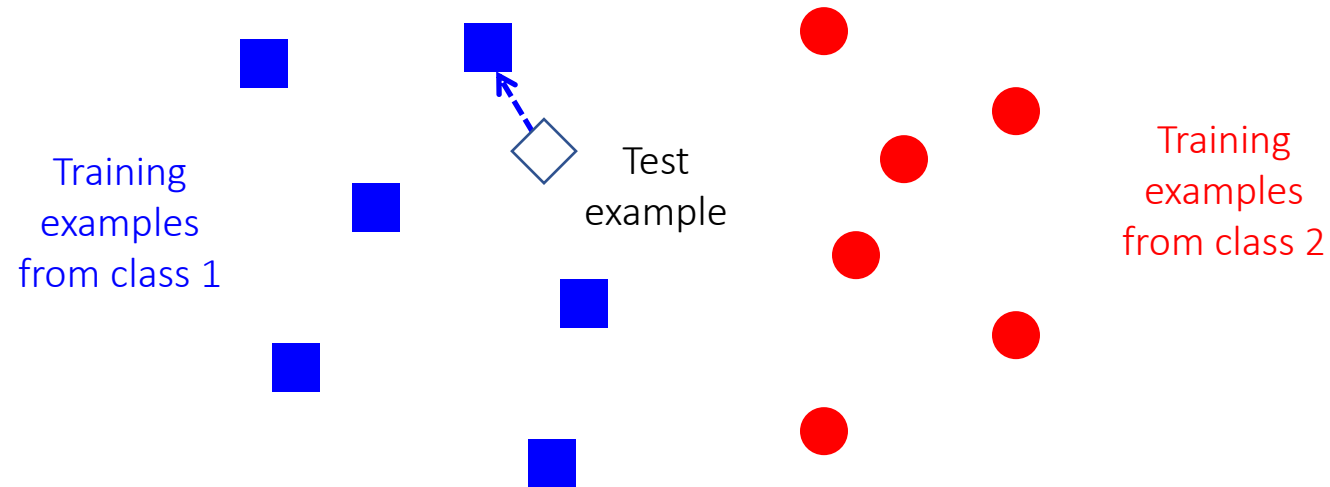
# Classifiers: Nearest neighbor



Slide credit: L. Lazebnik



# Classifiers: Nearest neighbor



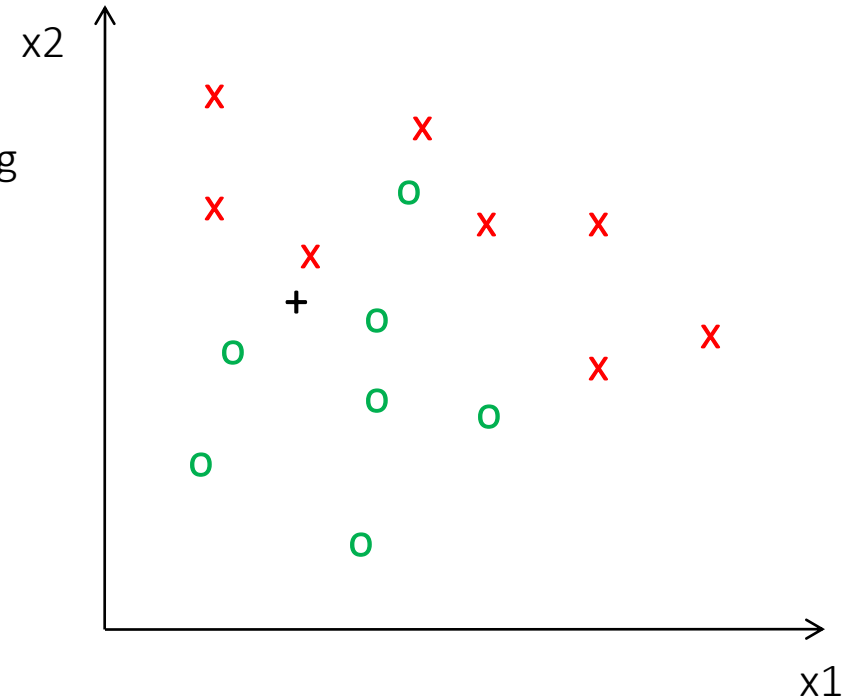
Slide credit: L. Lazebnik

# Today's Agenda

- A simple Image Classification pipeline
  - Classification overview
- K-nearest neighbor algorithm
  - kNN: algorithm
  - kNN: analysis

# K-nearest neighbor

- Algorithm (training):
  1. Store all training data points  $x_i$  with their corresponding category labels  $y_i$
- Algorithm (testing):
  1. We are given a new test point  $x$
  2. Compute distance to all training data points
  3. Select  $k$  training points closest to  $x$
  4. Assign  $x$  to label  $y$  that is most common among the  $k$  nearest neighbors.

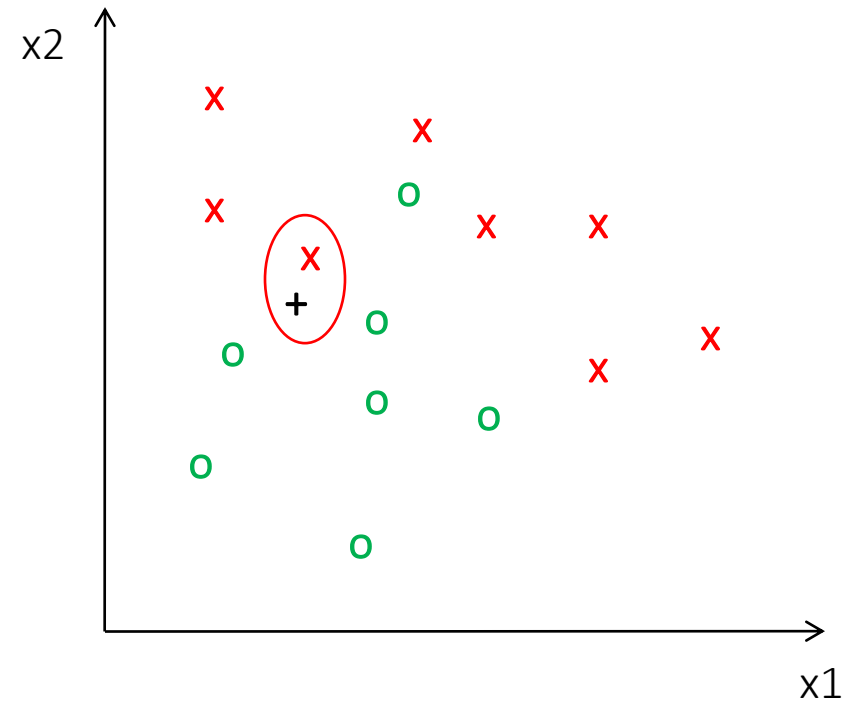


- Distance measurement: – Euclidean

$$Dist(X^n, X^m) = \sqrt{\sum_{i=1}^D (X_i^n - X_i^m)^2}$$

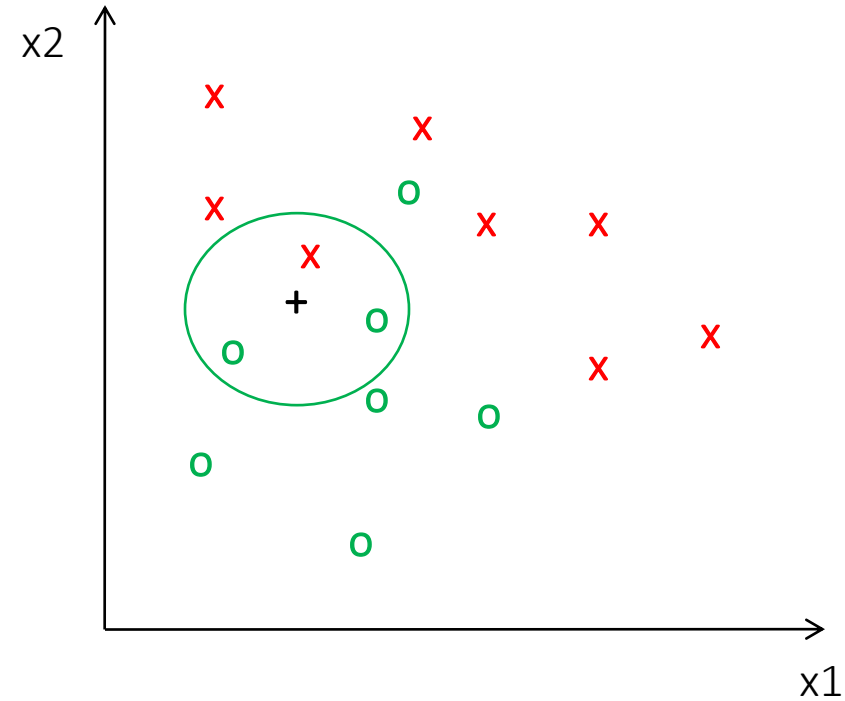
Where  $X^n$  and  $X^m$  are the n-th and m-th data points

# 1-nearest neighbor

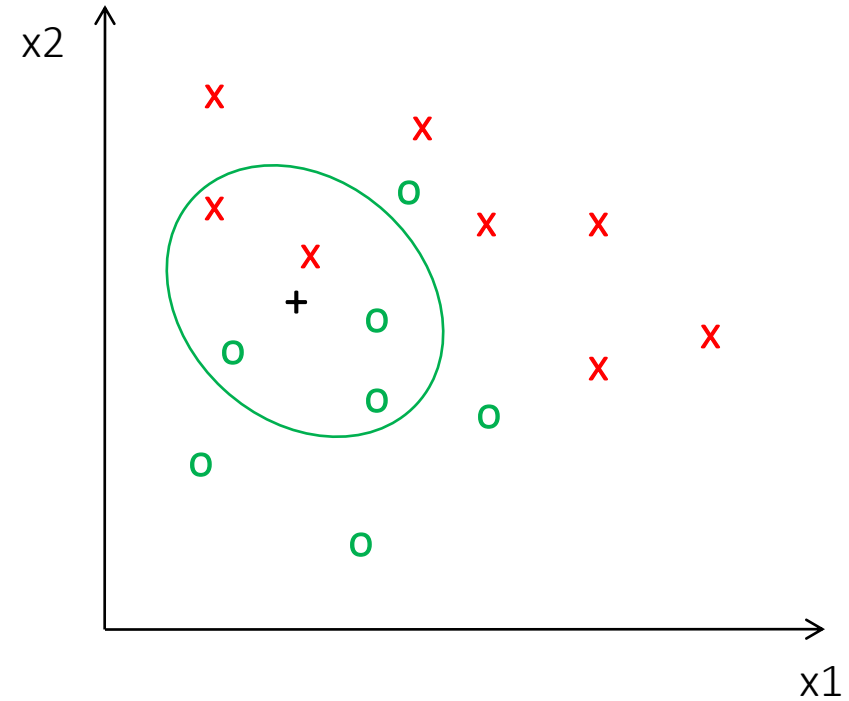




# 3-nearest neighbor



# 5-nearest neighbor





# Today's Agenda

- A simple Image Classification pipeline
  - Classification overview
- K-nearest neighbor algorithm
  - kNN: algorithm
  - kNN: analysis



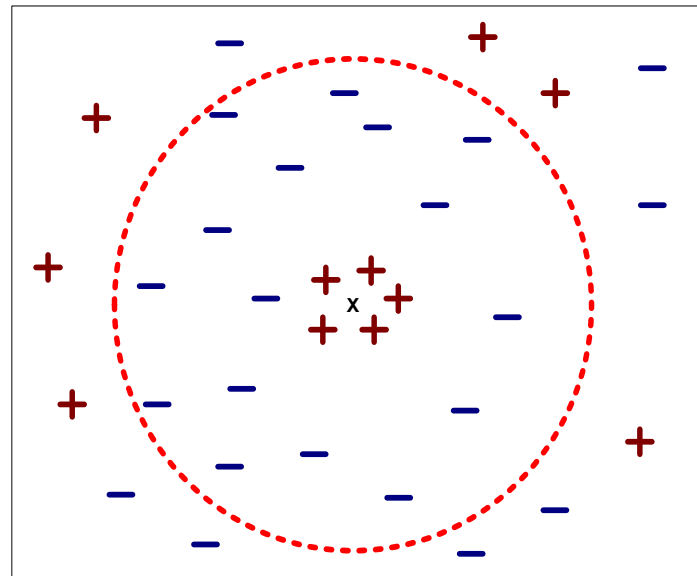
# K-NN: a very useful algorithm

- Simple, a good one to try first
- Very flexible decision boundaries



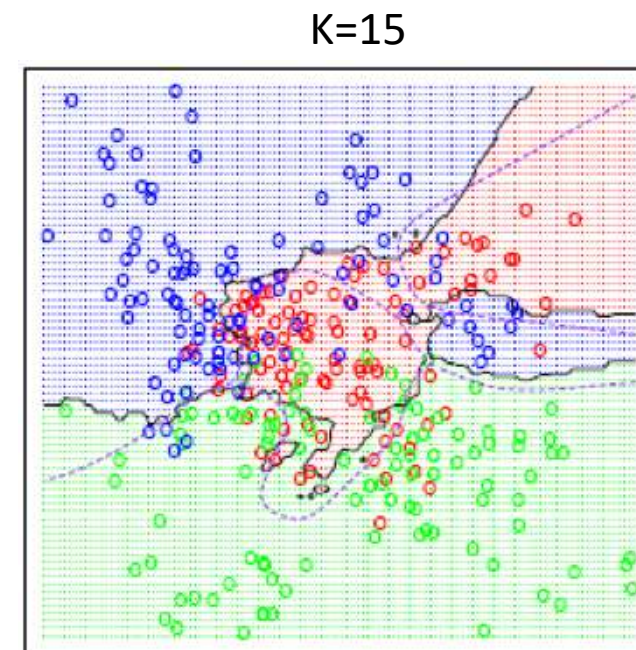
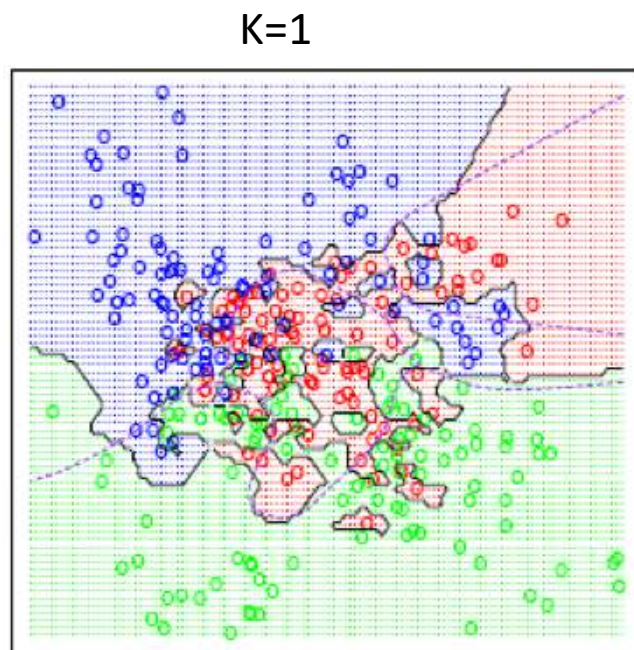
# K-NN: issues to keep in mind

- Choosing the value of k:
  - If too small, sensitive to noise points
  - If too large, neighborhood may include points from other classes



# K-NN: issues to keep in mind

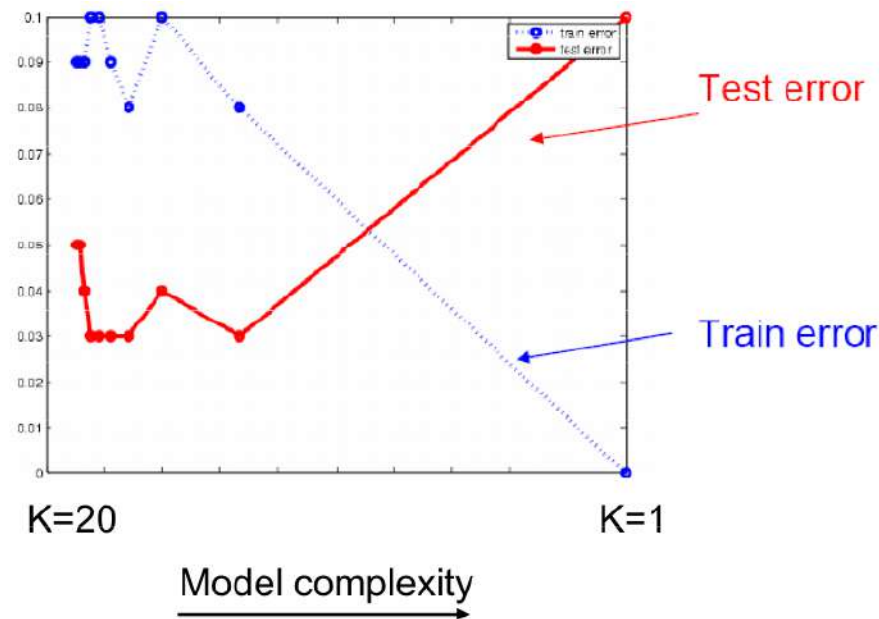
- Choosing the value of  $k$ :
  - If too small, sensitive to noise points
  - If too large, neighborhood may include points from other classes





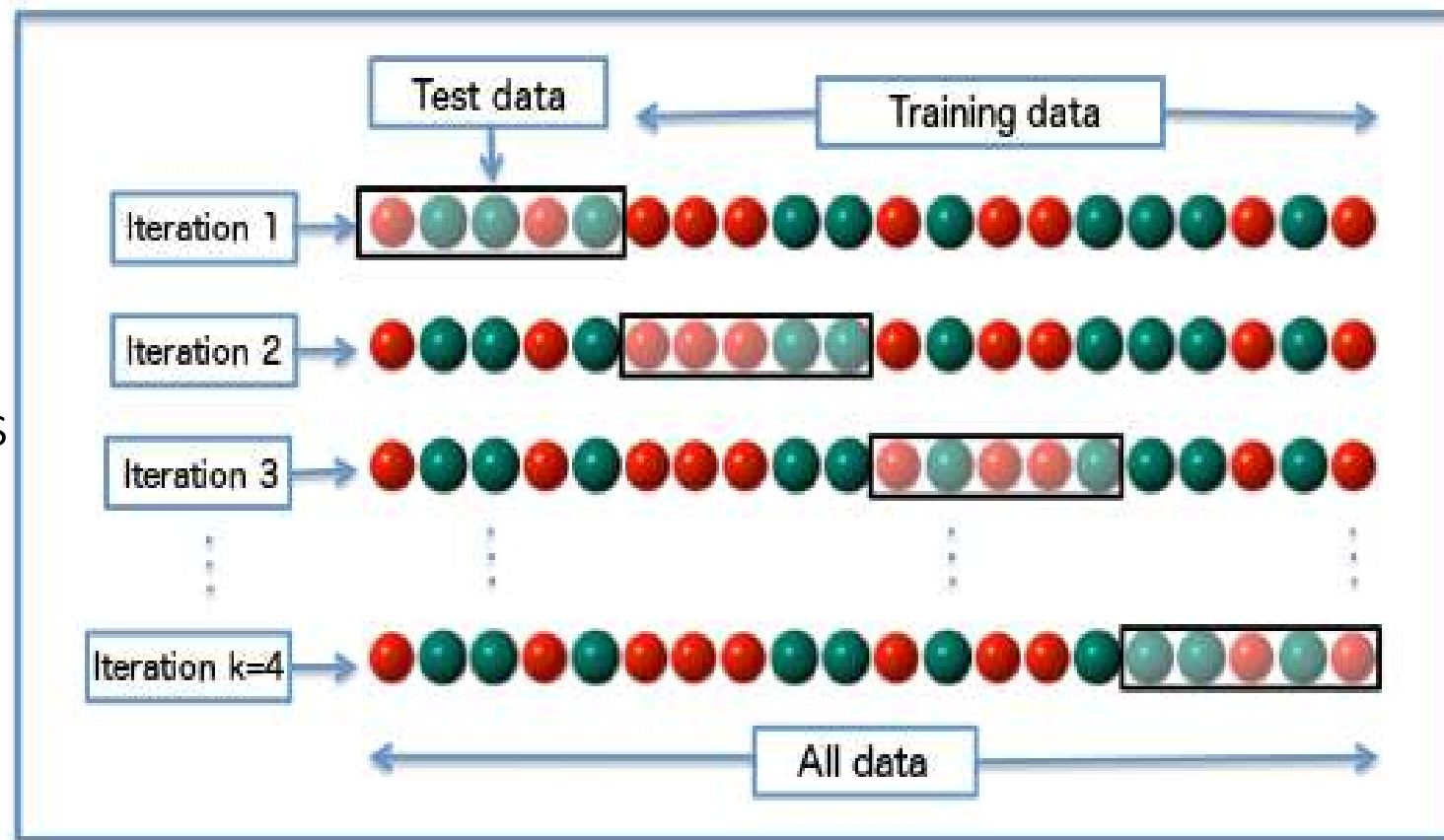
# K-NN: issues to keep in mind

- Choosing the value of  $k$ :
  - If too small, sensitive to noise points
  - If too large, neighborhood may include points from other classes
  - **Solution**: cross validate!



# Cross validation

- For each value of  $k$  in the nearest neighbors algorithm:
- Create multiple train/test splits
  - For each split:
    - Measure performance
- Average performance over all splits
  
- Select  $k$  with best average performance





# K-NN: issues to keep in mind

- Choosing the value of  $k$ :
  - If too small, sensitive to noise points
  - If too large, neighborhood may include points from other classes
  - **Solution**: cross validate!
- Can produce counter-intuitive results (using Euclidean measure)

# Euclidean measure

1 1 1 1 1 1 1 1 1 1 0

0 1 1 1 1 1 1 1 1 1 1

$d = 1.4142$

VS

1 0 0 0 0 0 0 0 0 0 0 0

0 0 0 0 0 0 0 0 0 0 0 1

$d = 1.4142$

# K-NN: issues to keep in mind

- Choosing the value of  $k$ :
  - If too small, sensitive to noise points
  - If too large, neighborhood may include points from other classes
  - **Solution**: cross validate!
- Can produce counter-intuitive results (using Euclidean measure)
  - **Solution**: normalize the vectors to unit length

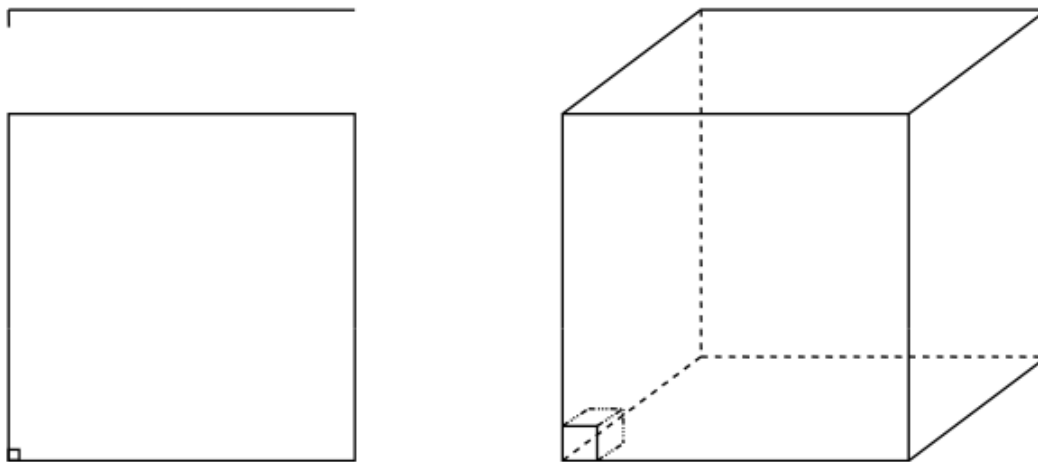
# K-NN: issues to keep in mind

- Choosing the value of  $k$ :
  - If too small, sensitive to noise points
  - If too large, neighborhood may include points from other classes
  - Solution: cross validate!
- Can produce counter-intuitive results (using Euclidean measure)
  - Solution: normalize the vectors to unit length
- Curse of Dimensionality



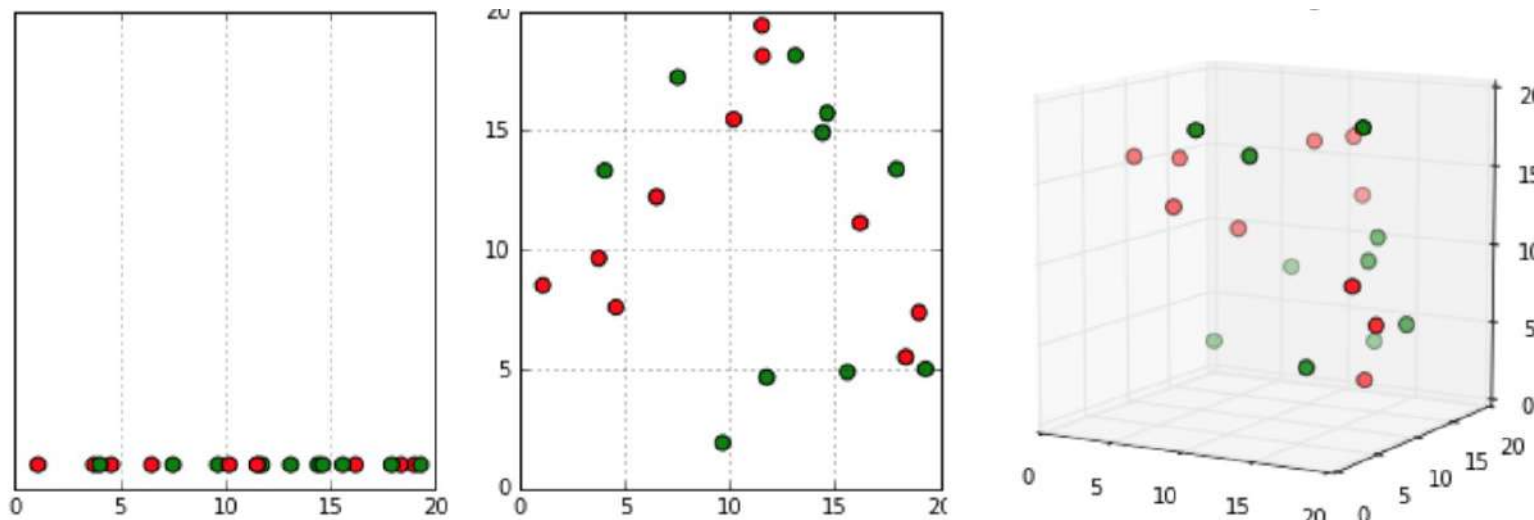
# Curse of dimensionality

- Assume 5000 points uniformly distributed in the unit (hyper)cube and we want to apply 5-NN. Suppose our query point is at the origin.
  - In 1-dimension, we must go a distance of  $5/5000=0.001$  on average to capture 5 nearest neighbors.
  - In 2 dimensions, we must go  $\sqrt{0.001}$  to get a square that contains 0.001 of the volume.
  - In  $d$  dimensions, we must go  $(0.001)^{1/d}$



# Curse of dimensionality

- Assume 5000 points uniformly distributed in the unit hypercube and we want to apply 5-NN. Suppose our query point is at the origin.
  - In 1-dimension, we must go a distance of  $5/5000=0.001$  on average to capture 5 nearest neighbors.
  - In 2 dimensions, we must go  $\sqrt{0.001}$  to get a square that contains 0.001 of the volume.
  - In d dimensions, we must go  $(0.001)^{1/d}$



# K-NN: issues to keep in mind

- Choosing the value of  $k$ :
  - If too small, sensitive to noise points
  - If too large, neighborhood may include points from other classes
  - **Solution**: cross validate!
- Can produce counter-intuitive results (using Euclidean measure)
  - **Solution**: normalize the vectors to unit length
- Curse of Dimensionality
  - **Solution**: no good one – need to get more data

**MAI4CAREU**

Master programmes in Artificial  
Intelligence 4 Careers in Europe



**CYENS**  
CENTRE OF EXCELLENCE



# Thank you.







University of Cyprus – MSc Artificial Intelligence

# MAI644 – COMPUTER VISION

## Lecture 12: Visual Bag of Words

**Melinos Averkiou**

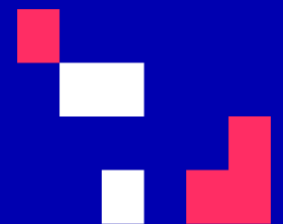
CYENS Centre of Excellence

University of Cyprus - Department of Computer Science

[m.averkiou@cyens.org.cy](mailto:m.averkiou@cyens.org.cy)



**CYENS**  
CENTRE OF EXCELLENCE



# Last time

- A simple Image Classification pipeline
  - Classification overview
- K-nearest neighbor algorithm
  - kNN: algorithm
  - kNN: analysis

# Today's Agenda

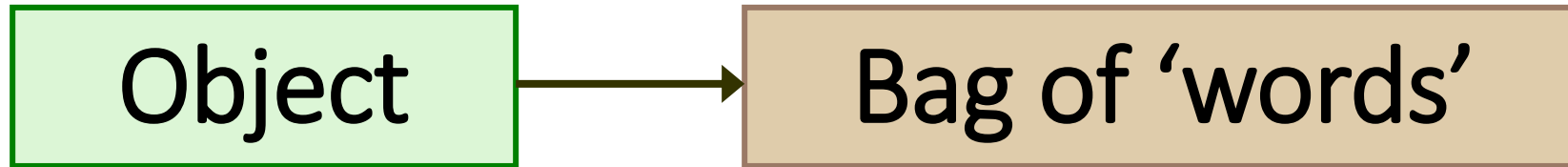
- Visual bag of words (BoW)
  - Background
  - Algorithm
- Applications
  - Image search
- Spatial Pyramid Matching

[material based on Niebles-Krishna]

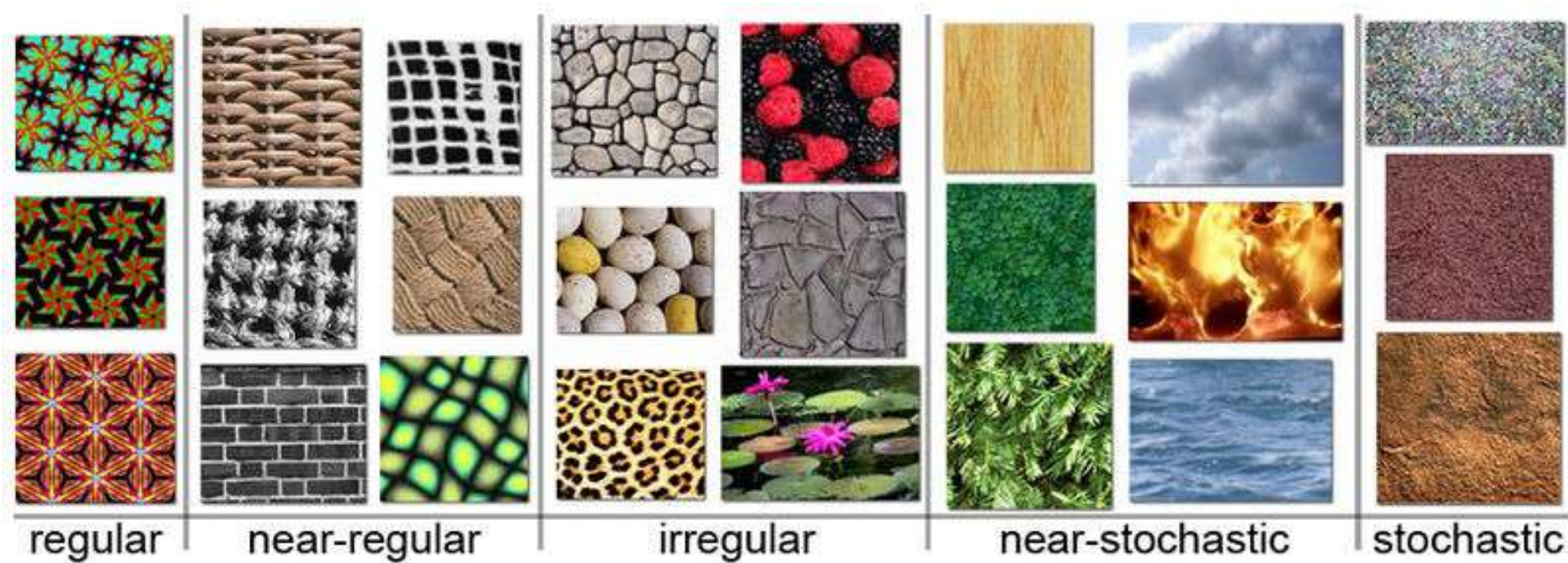
# Today's Agenda

- Visual bag of words (BoW)
  - Background
  - Algorithm
- Applications
  - Image search
- Spatial Pyramid Matching





# Origin 1: Texture Recognition

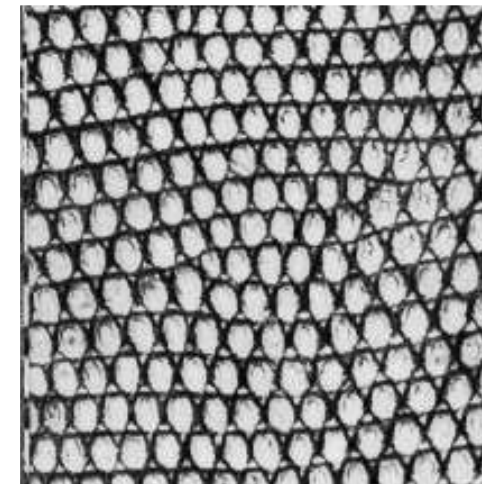
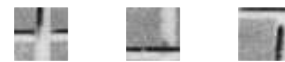
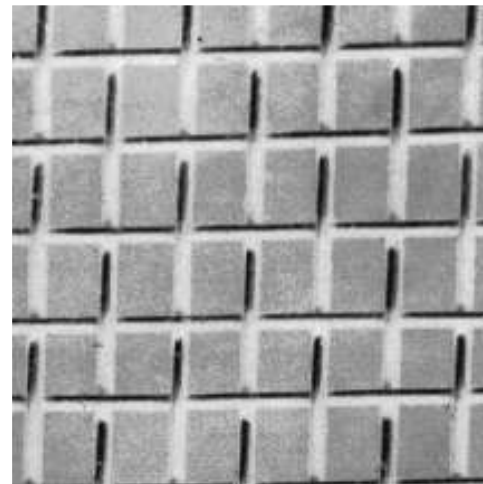
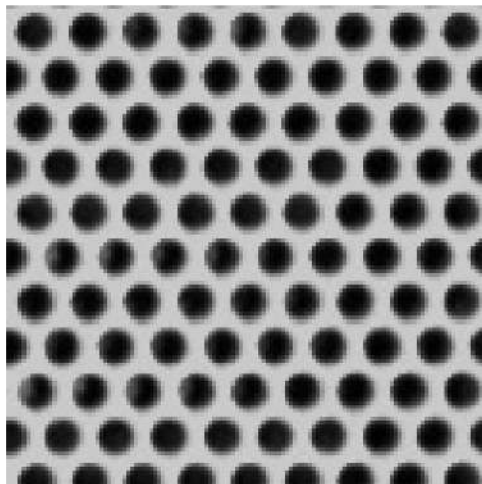


Example textures (from Wikipedia)



# Origin 1: Texture Recognition

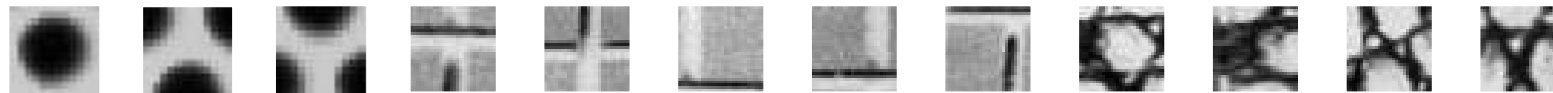
Texture is characterized by the repetition of basic elements or *textons*



Julesz, 1981; Cula & Dana, 2001; Leung & Malik 2001; Mori, Belongie & Malik, 2001; Schmid 2001; Varma & Zisserman, 2002, 2003; Lazebnik, Schmid & Ponce, 2003

# Origin 1: Texture Recognition

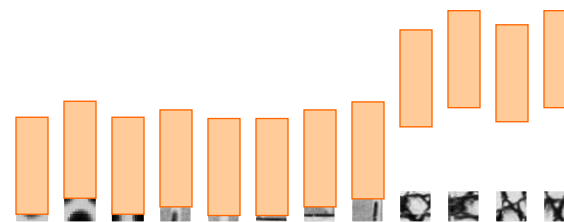
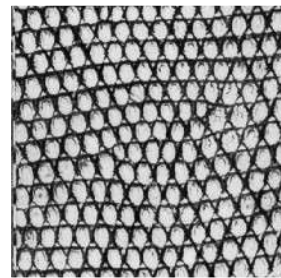
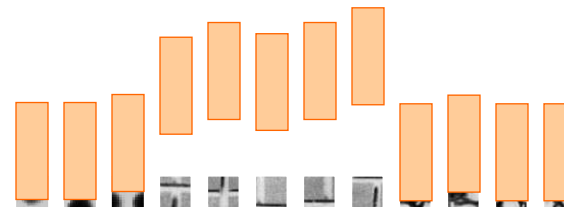
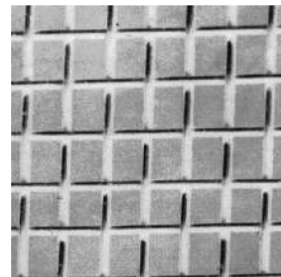
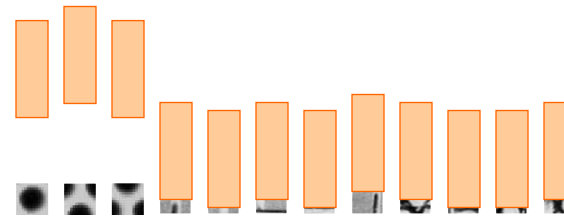
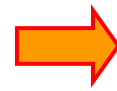
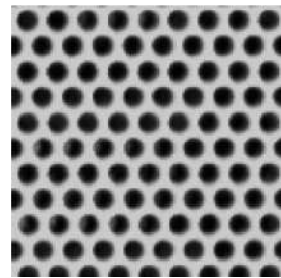
Recognition based on identity of the textons, not their spatial arrangement  
(although that is very important too!)



Universal texton dictionary



## Origin 1: Texture Recognition



# Origin 2: Bag-of-words models

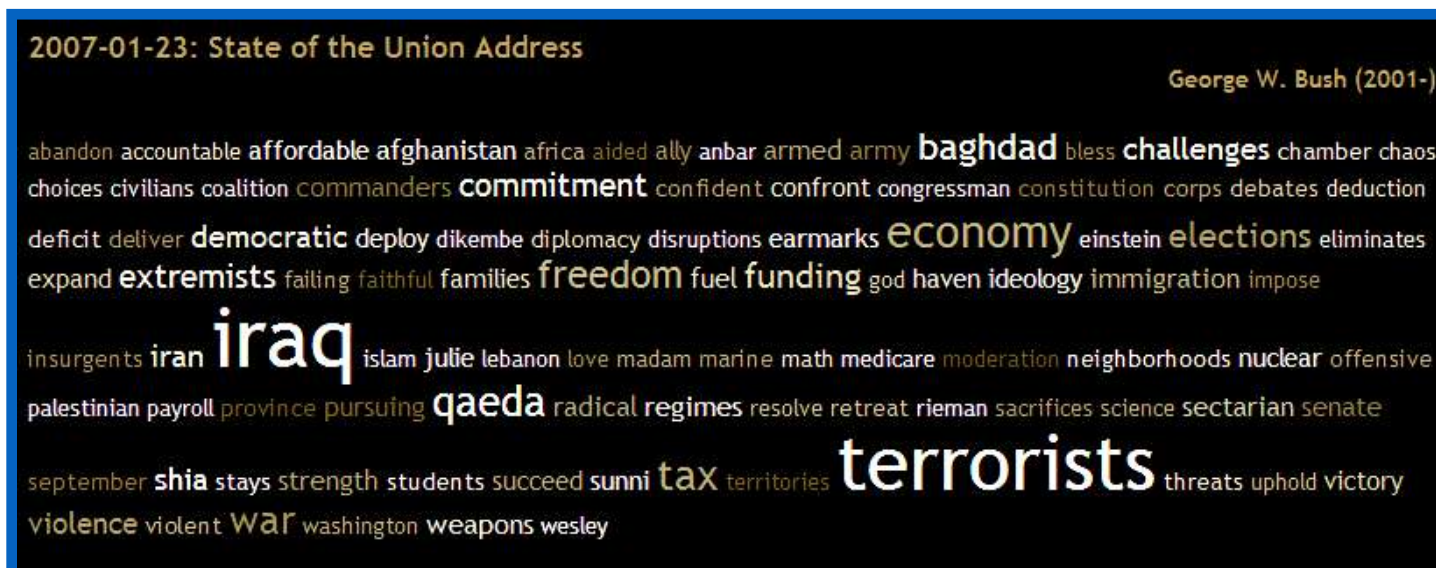
Orderless document representation: frequencies of words from a dictionary

Salton & McGill (1983)

# Origin 2: Bag-of-words models

Orderless document representation: frequencies of words from a dictionary

Salton & McGill (1983)



US Presidential Speeches Tag Cloud

<http://chir.ag/phernalia/preztags/>

# Origin 2: Bag-of-words models

Orderless document representation: frequencies of words from a dictionary

Salton & McGill (1983)



US Presidential Speeches Tag Cloud  
<http://chir.ag/phernalia/preztags/>



# Origin 2: Bag-of-words models

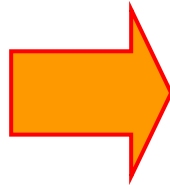
Orderless document representation: frequencies of words from a dictionary

Salton & McGill (1983)



US Presidential Speeches Tag Cloud  
<http://chir.ag/phernalia/preztags/>

# Bags of features for object recognition



face, flowers, building

Works pretty well for image-level classification and for recognizing object *instances*

Csurka et al. (2004), Willamowski et al. (2005), Grauman & Darrell (2005), Sivic et al. (2003, 2005)

# Bags of features for object recognition



class	bag of features	bag of features	Parts-and-shape model
	Zhang et al. (2005)	Willamowski et al. (2004)	Fergus et al. (2003)
airplanes	<b>98.8</b>	97.1	90.2
cars (rear)	98.3	<b>98.6</b>	90.3
cars (side)	<b>95.0</b>	87.3	88.5
faces	<b>100</b>	99.3	96.4
motorbikes	<b>98.5</b>	98.0	92.5
spotted cats	<b>97.0</b>	—	90.0

# Today's Agenda

- Visual bag of words (BoW)
  - Background
  - Algorithm
- Applications
  - Image search
- Spatial Pyramid Matching

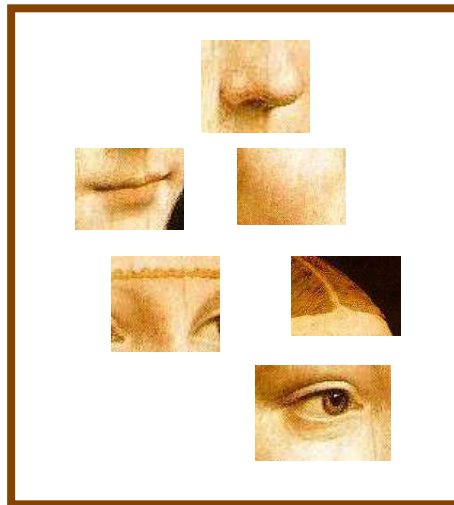


# Bag of features

- First, take a bunch of images, extract features, and build up a “dictionary” or “visual vocabulary” – a list of common features
- Given a new image, extract features and build a histogram – for each feature, find the closest visual word in the dictionary

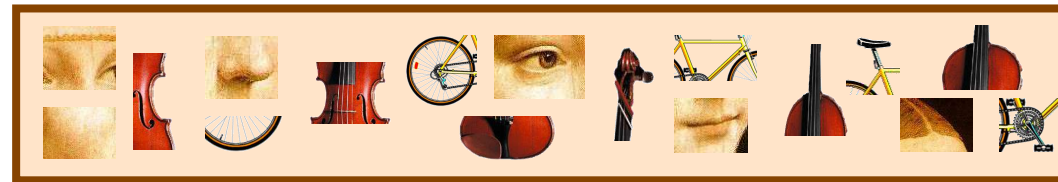
# Bag of features: outline

## 1. Extract features



# Bag of features: outline

1. Extract features
2. Learn “visual vocabulary”





# Bag of features: outline

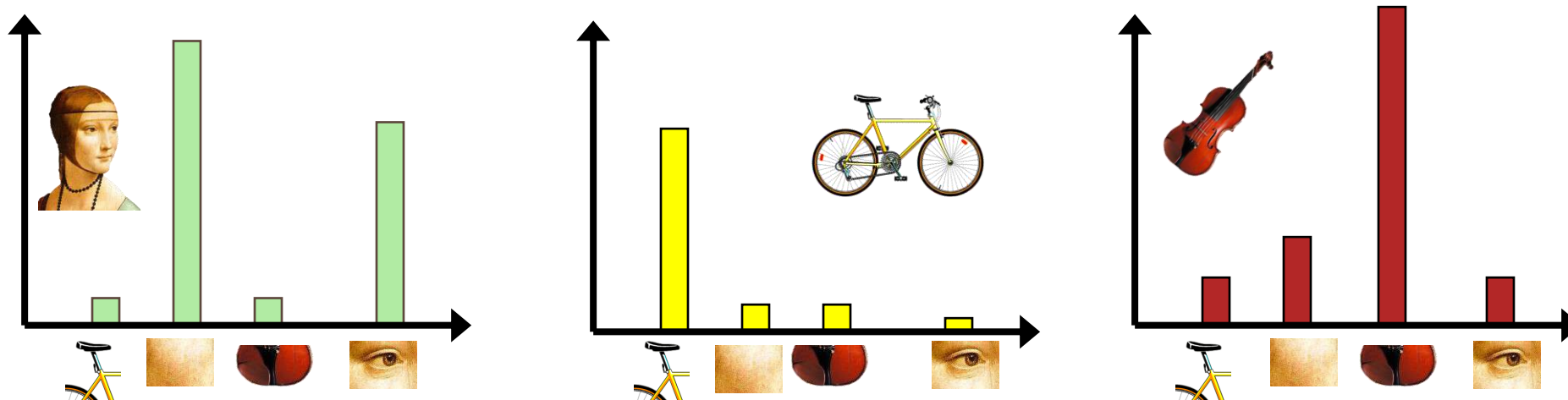
1. Extract features
2. Learn “visual vocabulary”
3. Quantize features using visual vocabulary





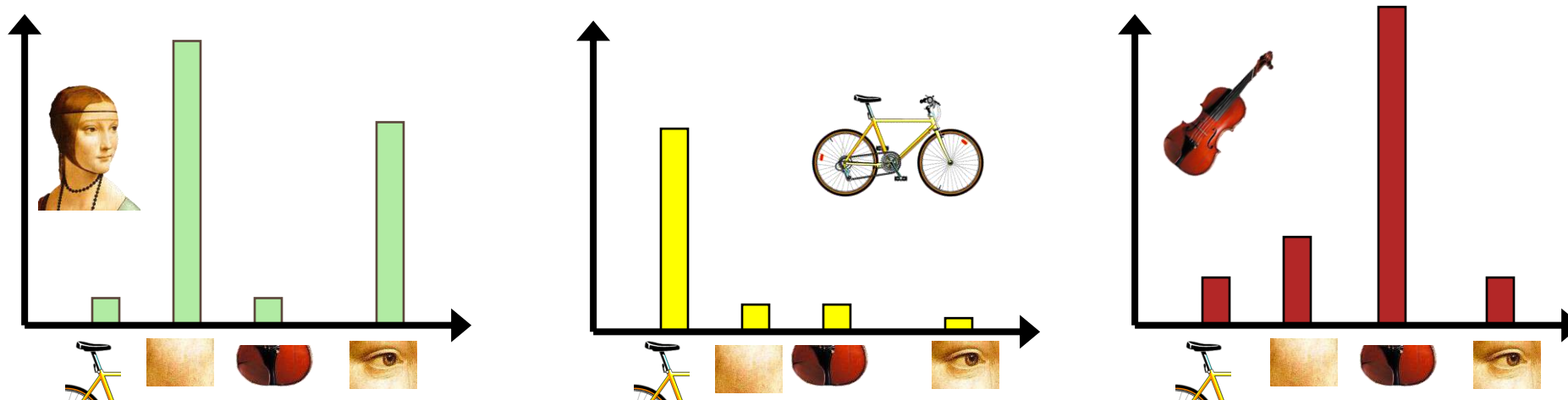
# Bag of features: outline

1. Extract features
2. Learn “visual vocabulary”
3. Quantize features using visual vocabulary
4. Represent images by frequencies of “visual words”



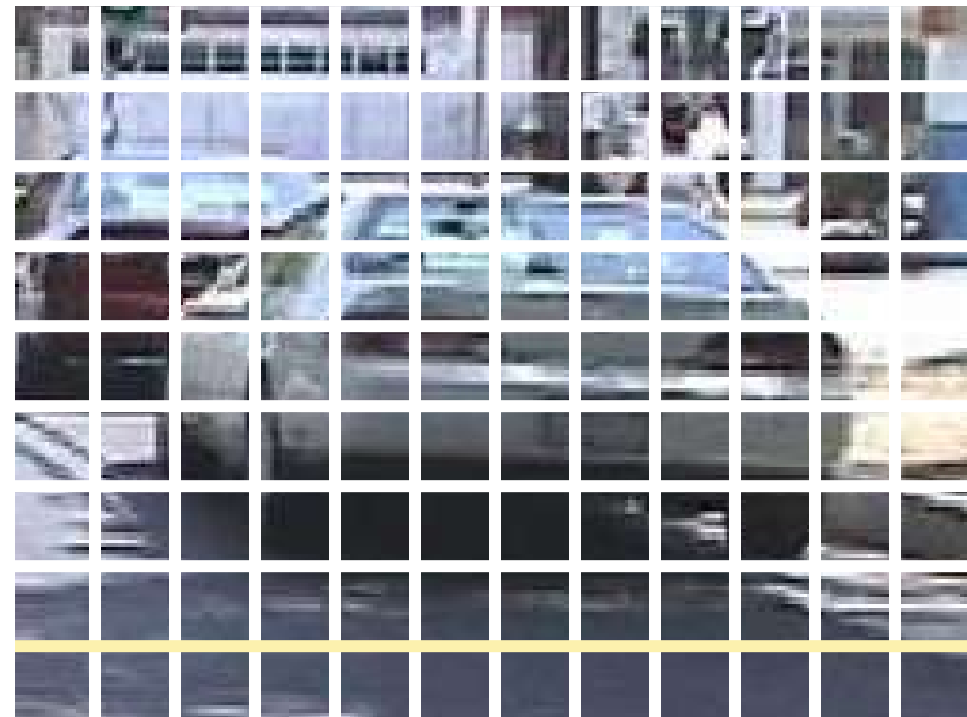
# Bag of features: outline

1. Extract features
2. Learn “visual vocabulary”
3. Quantize features using visual vocabulary
4. Represent images by frequencies of “visual words”



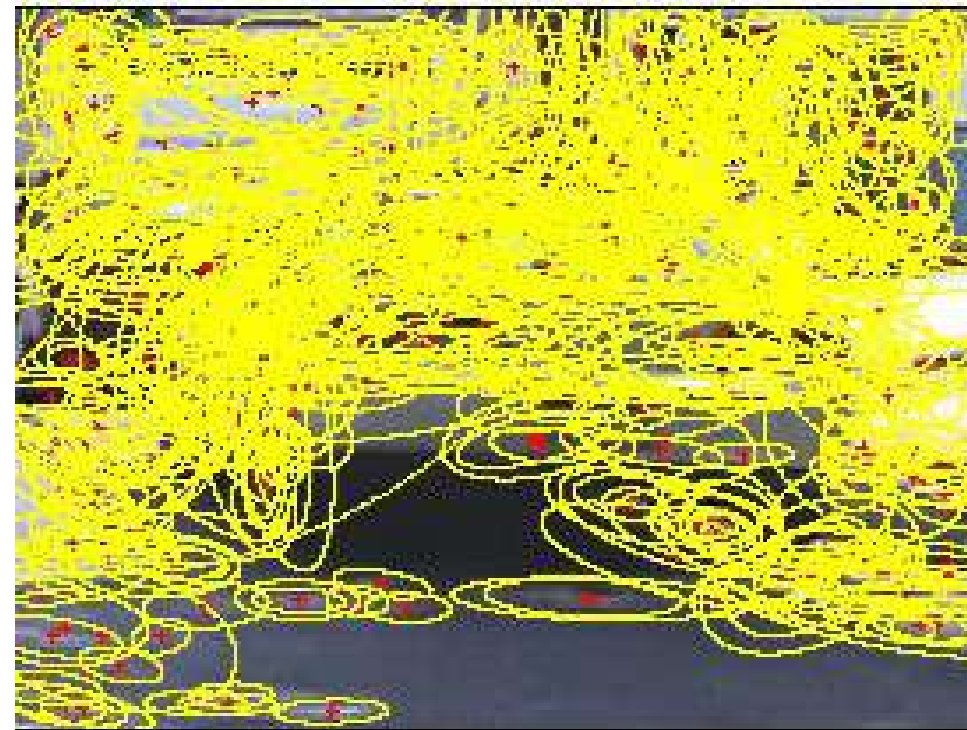
# 1. Feature extraction

- Regular grid
  - Vogel & Schiele, 2003
  - Fei-Fei & Perona, 2005



# 1. Feature extraction

- Regular grid
  - Vogel & Schiele, 2003
  - Fei-Fei & Perona, 2005
- Interest point detector
  - Csurka et al. 2004
  - Fei-Fei & Perona, 2005
  - Sivic et al. 2005



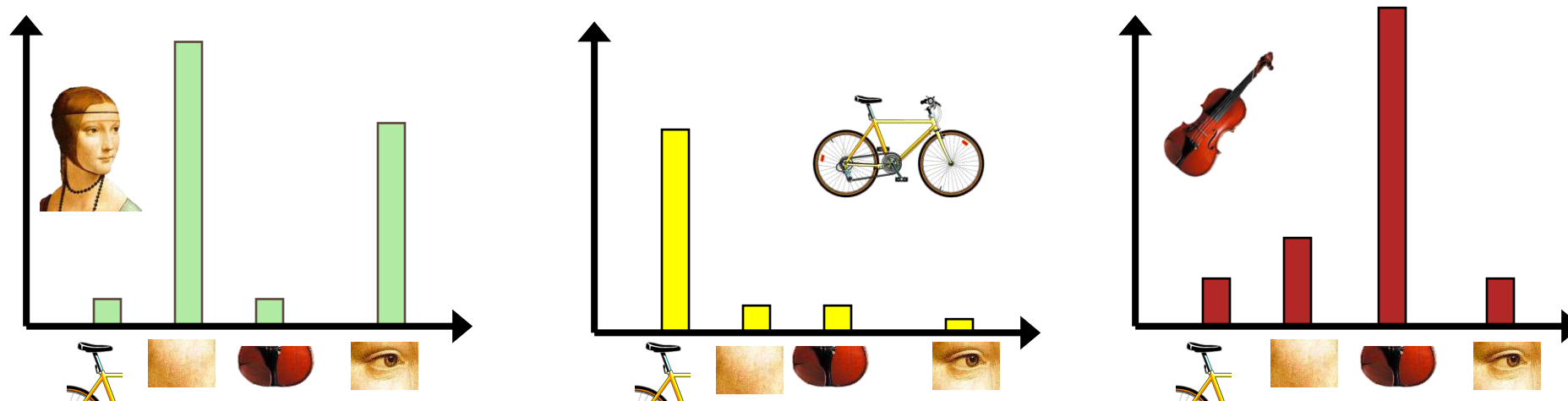


# 1. Feature extraction

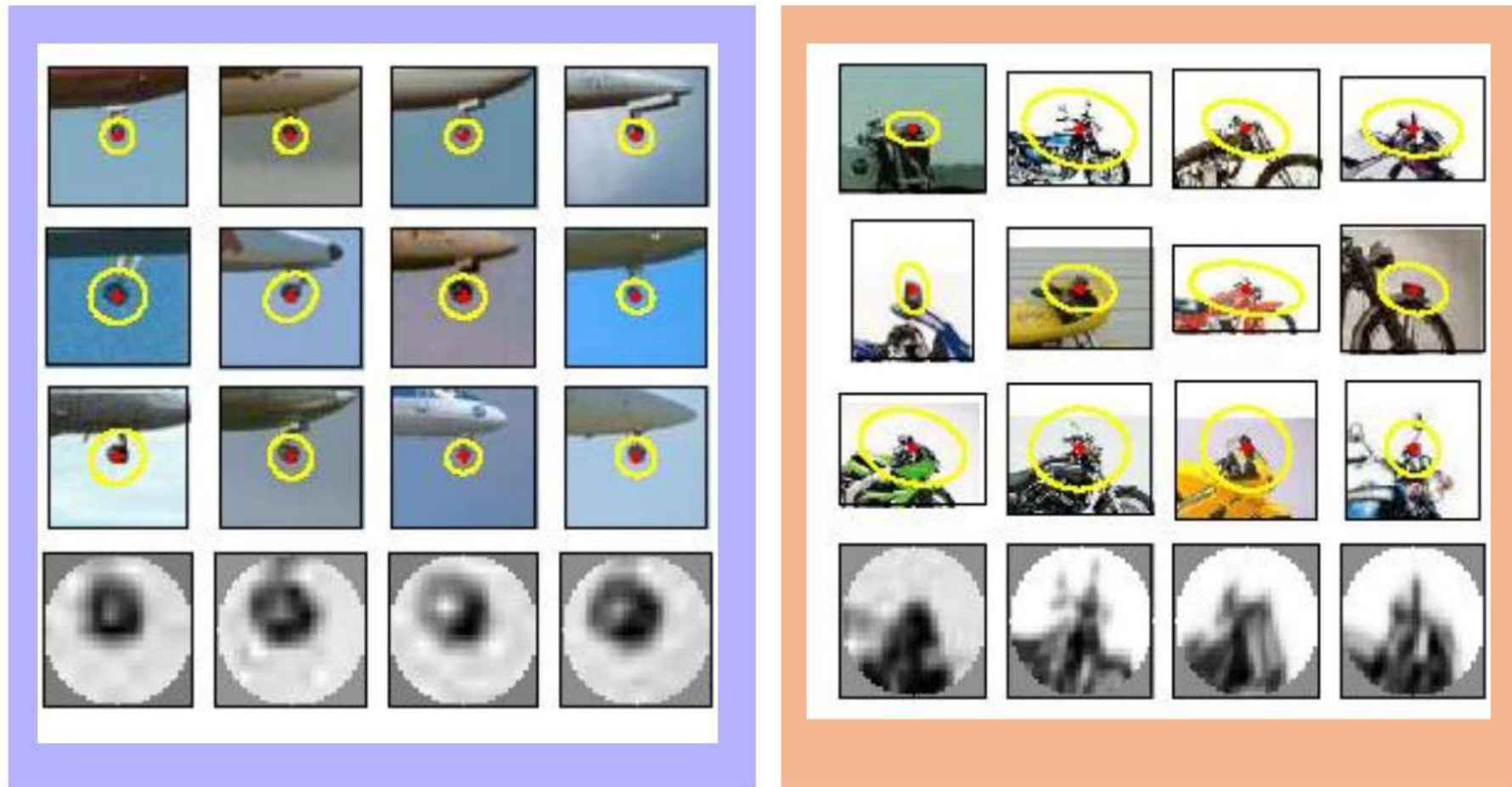
- Regular grid
  - Vogel & Schiele, 2003
  - Fei-Fei & Perona, 2005
- Interest point detector
  - Csurka et al. 2004
  - Fei-Fei & Perona, 2005
  - Sivic et al. 2005
- Other methods
  - Random sampling (Vidal-Naquet & Ullman, 2002)
  - Segmentation-based patches (Barnard et al. 2003)

# Bag of features: outline

1. Extract features
2. Learn “visual vocabulary”
3. Quantize features using visual vocabulary
4. Represent images by frequencies of “visual words”

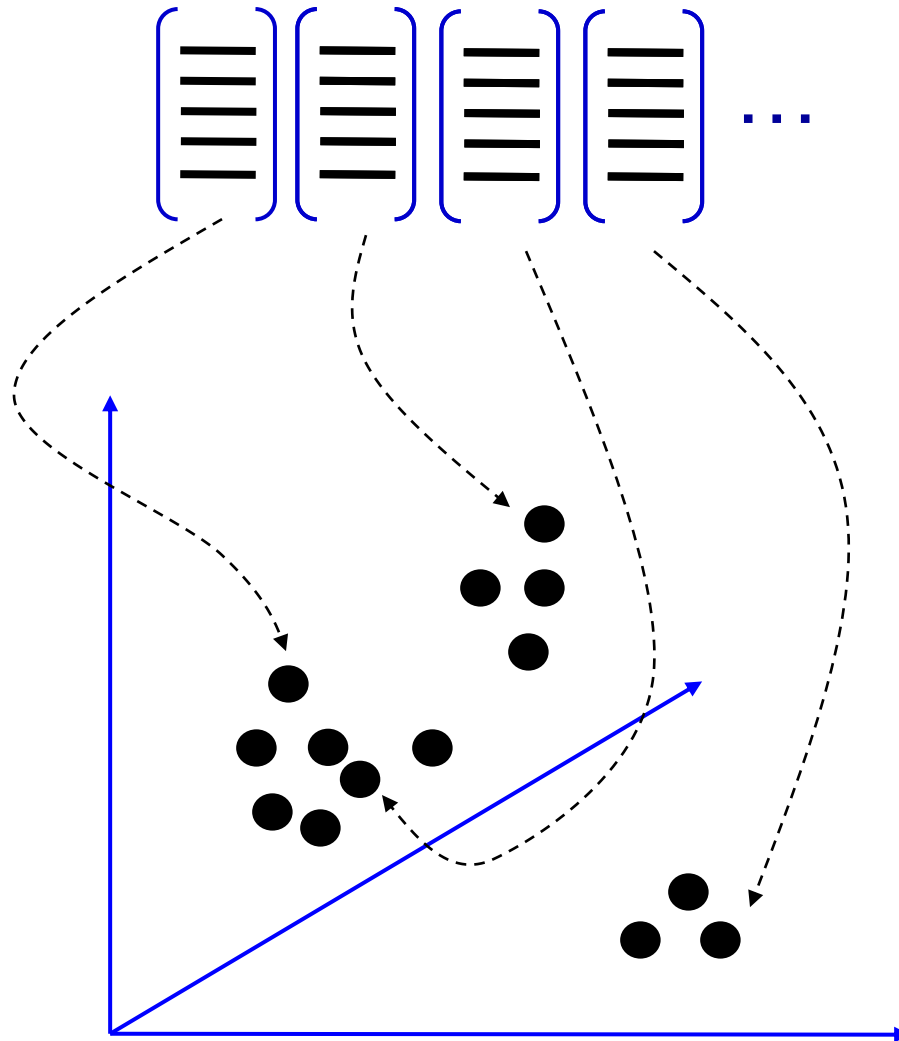


# Image patch examples of visual words



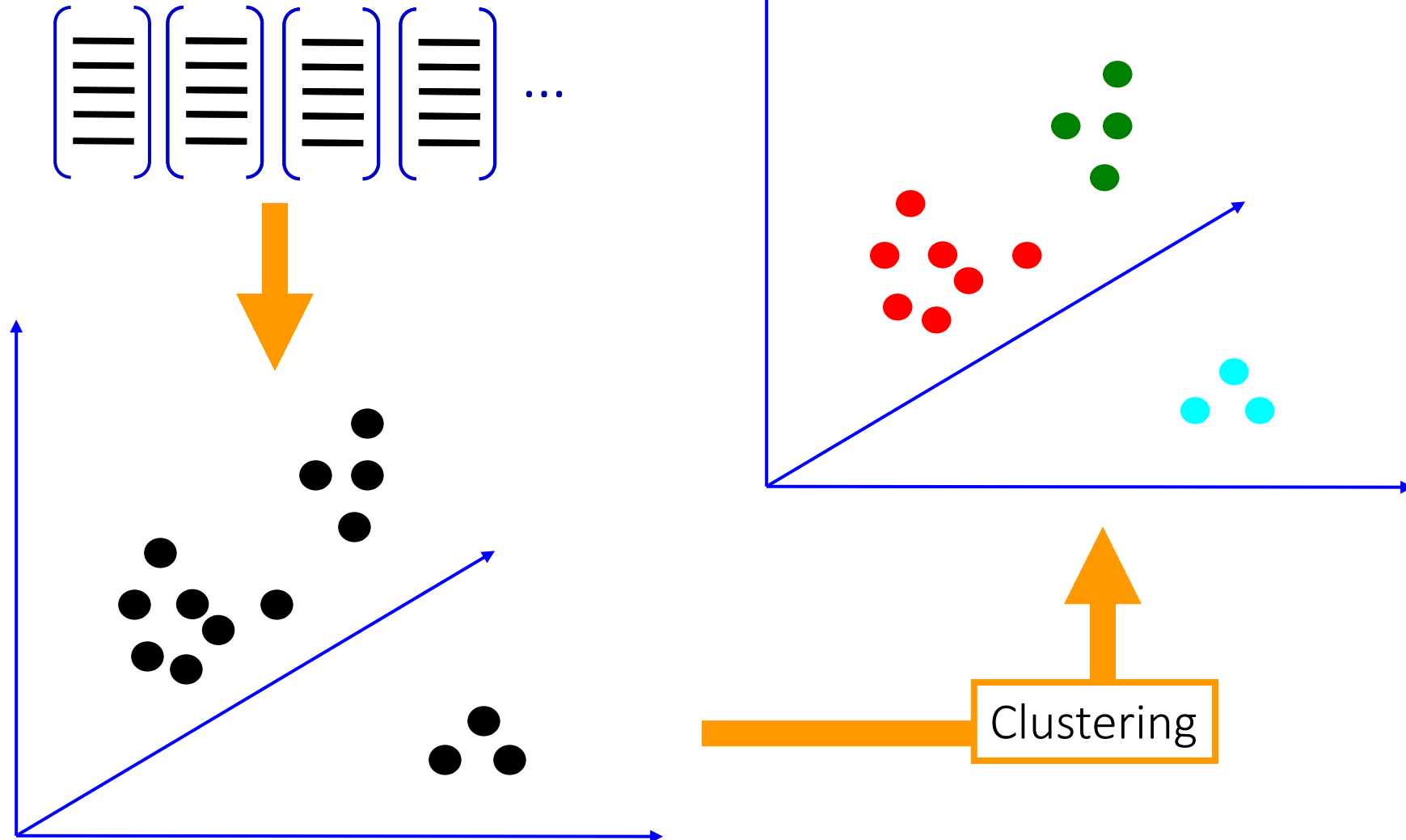
Sivic et al 2005

## 2. Learn the visual vocabulary



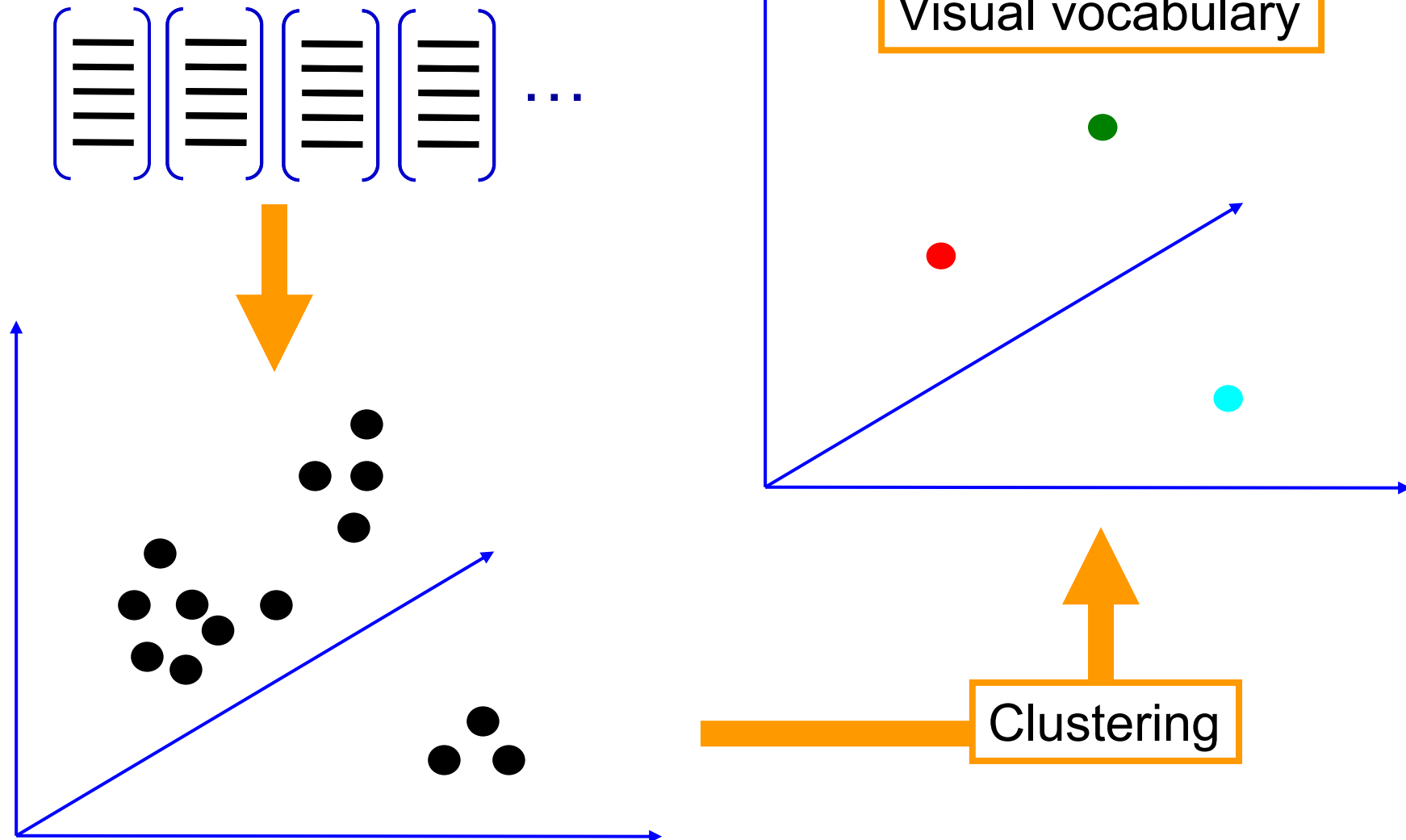


## 2. Learn the visual vocabulary



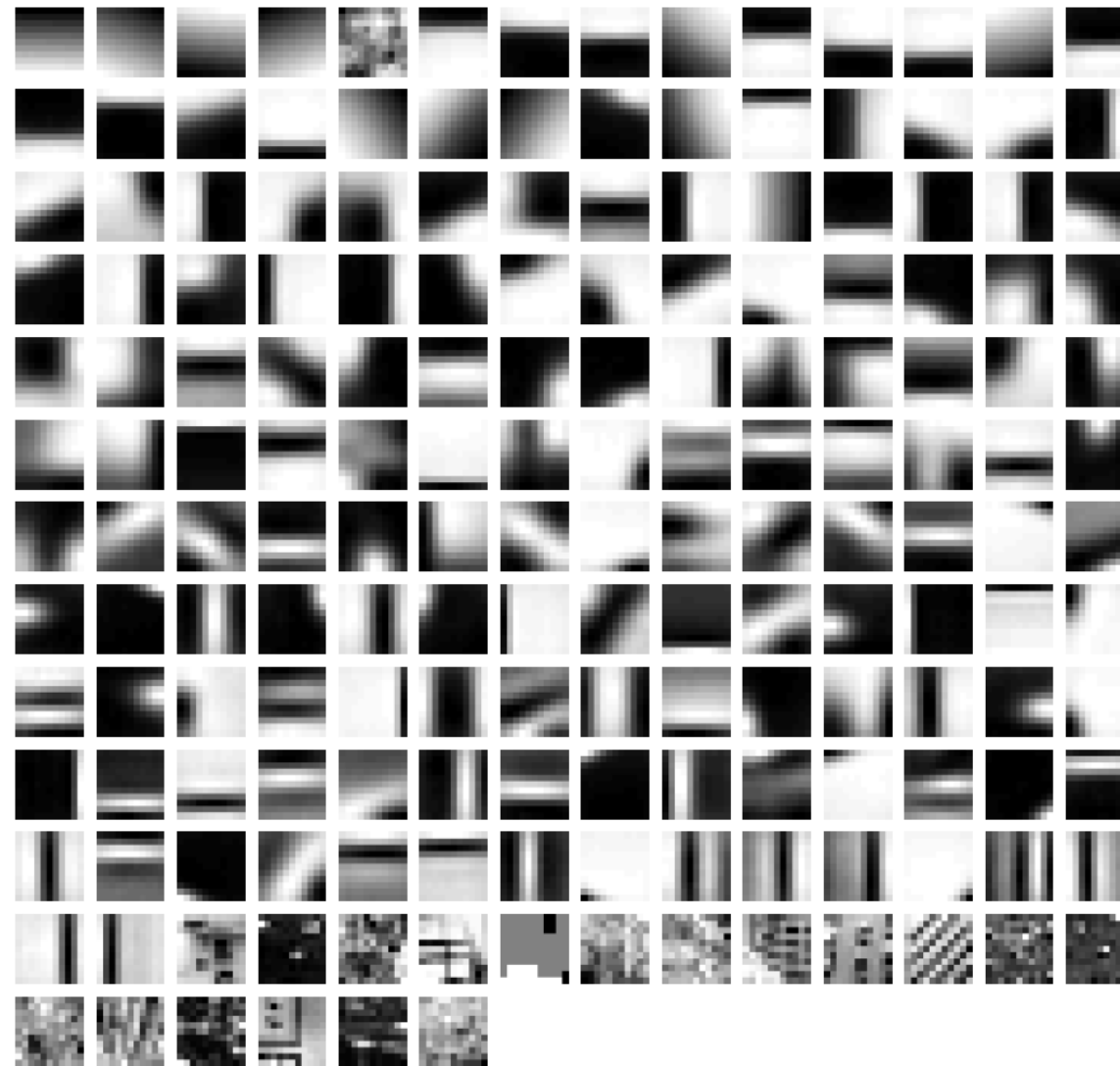
Slide credit: Josef Sivic

## 2. Learn the visual vocabulary



Slide credit: Josef Sivic

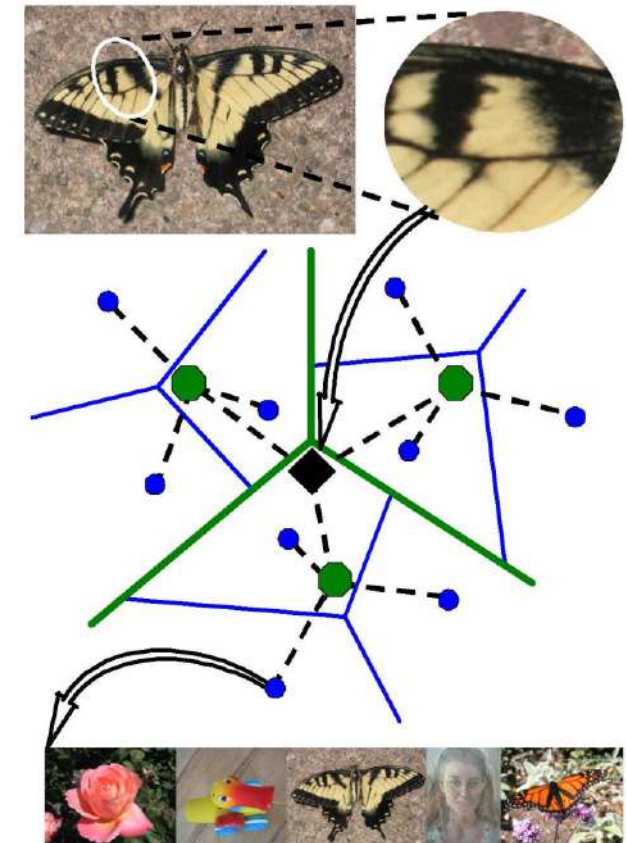
## Example visual vocabulary



Fei-Fei et al. 2005

# Visual vocabularies: issues

- How to choose vocabulary size?
  - Too small: visual words not representative of all patches
  - Too large: quantization artifacts, overfitting
- Computational efficiency
  - Solution: Vocabulary trees (Nister & Stewenius, 2006)





## Bag of features: outline

1. Extract features
2. Learn “visual vocabulary”
3. Quantize features using visual vocabulary
4. Represent images by frequencies of “visual words”

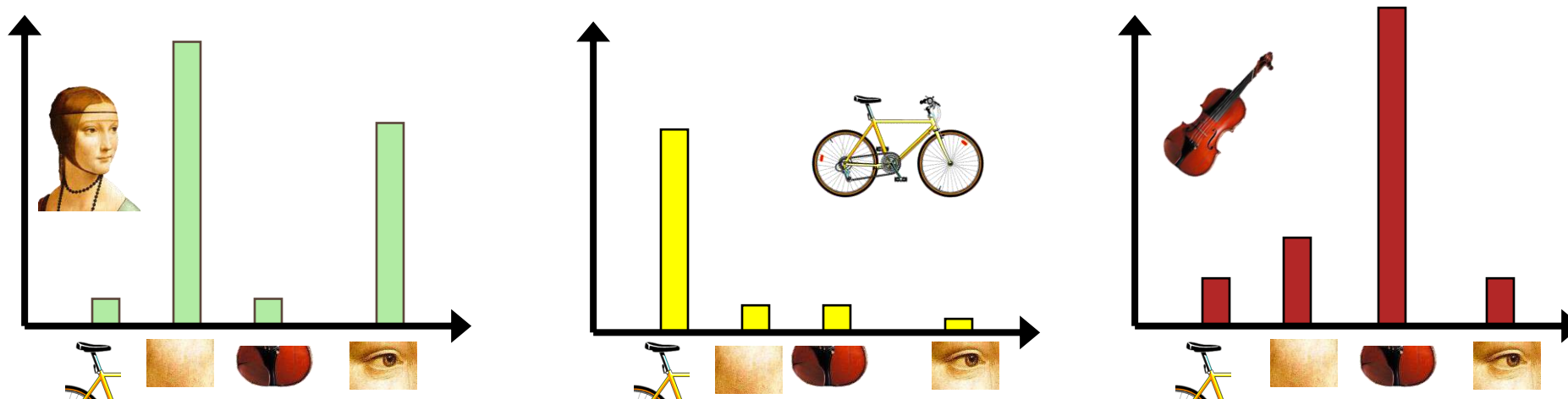


### 3. From clustering to vector quantization

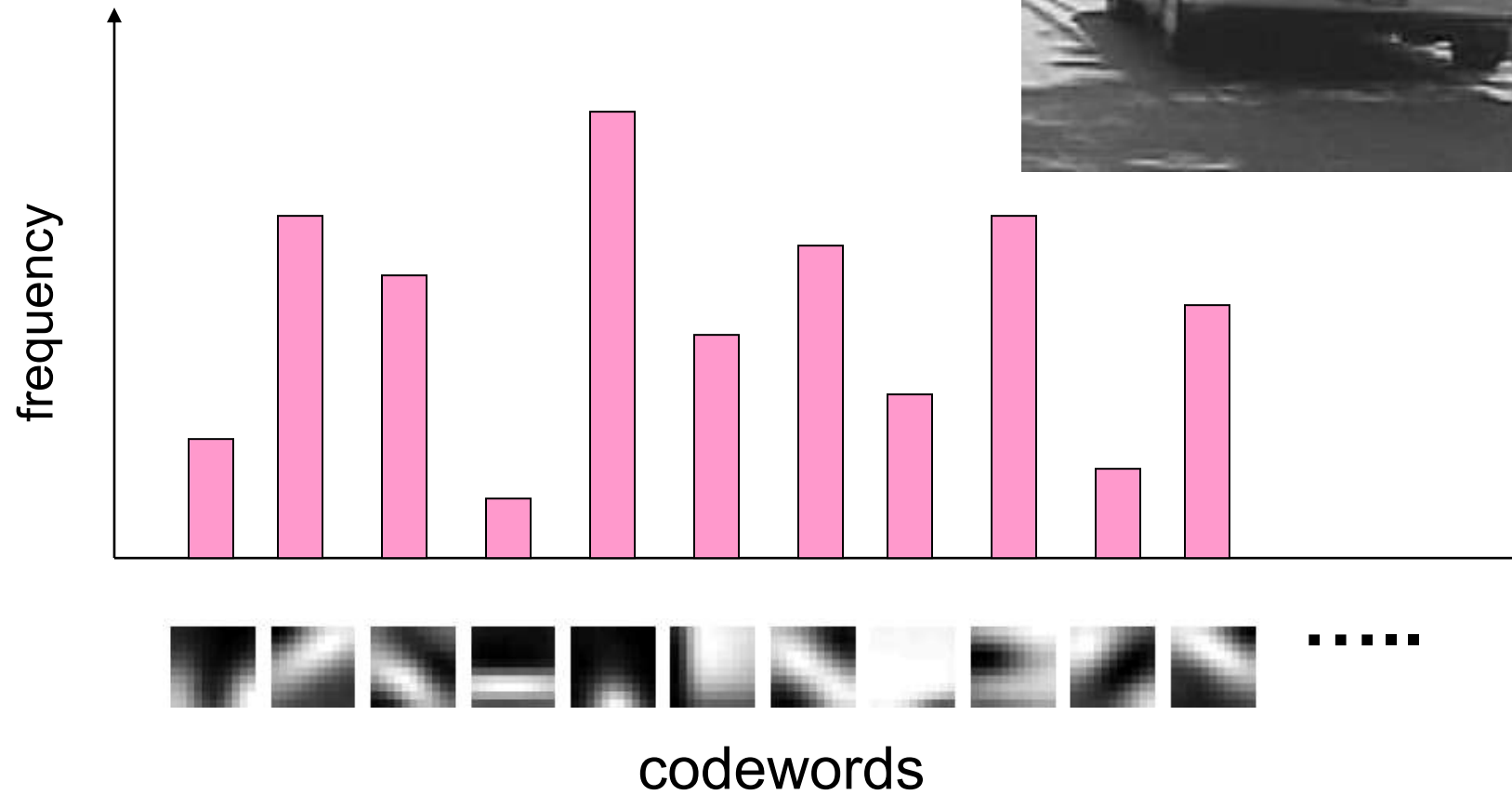
- Clustering is a common method for learning a visual vocabulary or codebook
  - Unsupervised learning process
  - Each cluster center produced by k-means becomes a codevector
  - Codebook can be learned on separate training set
  - Provided the training set is sufficiently representative, the codebook will be “universal”
- The codebook is used for quantizing features
  - A *vector quantizer* takes a feature vector and maps it to the index of the nearest codevector in a codebook
  - Codebook = visual vocabulary
  - Codevector = visual word

# Bag of features: outline

1. Extract features
2. Learn “visual vocabulary”
3. Quantize features using visual vocabulary
4. Represent images by frequencies of “visual words”



# 4. Image representation





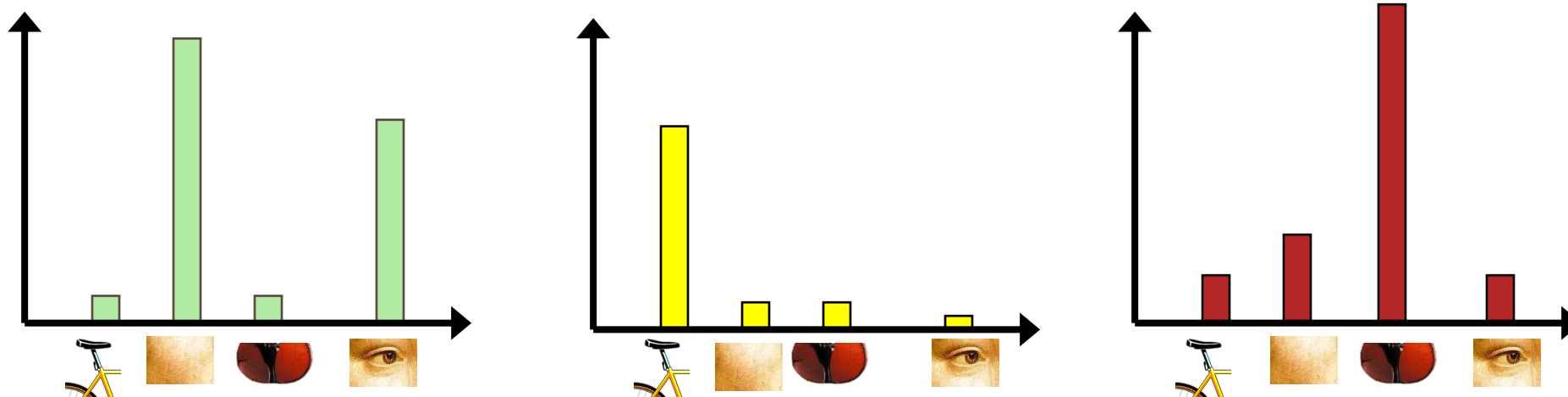
# Today's Agenda

- Visual bag of words (BoW)
  - Background
  - Algorithm
- Applications
  - Image search
- Spatial Pyramid Matching



# Image classification

- Given the bag-of-features representations of images from different classes, how do we learn a model for distinguishing them?



# Uses of BoW representation

- Treat as feature vector for standard classifier
  - e.g k-nearest neighbors, support vector machine
- Cluster BoW vectors over image collection
  - Discover visual themes

# Large-scale image search



11,400 images of game covers  
(Caltech games dataset)



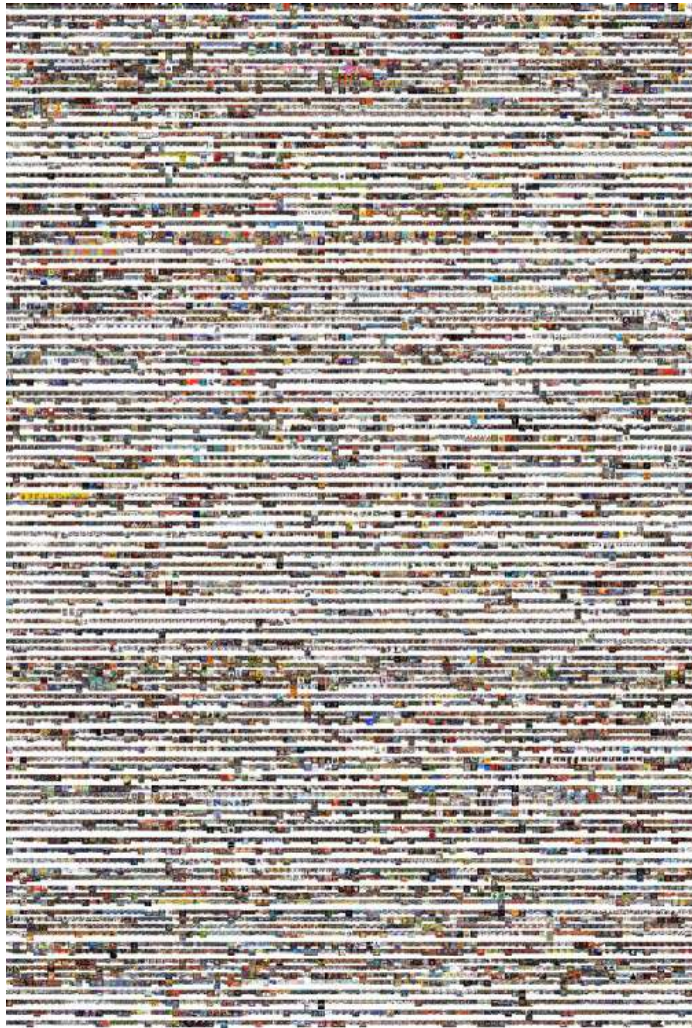
Bag-of-words models have been useful in matching an image to a large database of object *instances*



How do I find this image in the database?



# Large-scale image search



## Build the database:

- Extract features from the database images
- Learn a vocabulary using k-means (typical k: 100,000)
- Compute *weights* for each word
- Create an inverted file mapping words → images

# Weighting the words

- Just as with text, some visual words are more discriminative than others

***the, and, or*** vs. ***cow, AT&T, Cher***

- The bigger fraction of the documents a word appears in, the less useful it is for matching
  - e.g., a word that appears in *all* documents is not helping us

# TF-IDF weighting

- Instead of computing a regular histogram distance, we'll weight each word by its *inverse document frequency*
- Inverse Document Frequency (IDF) of word  $j$  =

$$\log \frac{\text{number of documents}}{\text{number of documents in which } j \text{ appears}}$$

# TF-IDF weighting

- Term Frequency (TF) of word  $j$  is the number of times it appears in the 'document', i.e., the image
- To compute the value of bin  $j$  in image  $I$ , compute TF-IDF:

*Term frequency of  $j$  in  $I$  × Inverse Document Frequency of  $j$*





# Inverted file

- Each image has  $\sim 1,000$  features
- We have  $\sim 100,000$  visual words
  - each histogram is extremely sparse (mostly zeros)
- Inverted file
  - mapping from words to 'documents', i.e., images

```
"a": {2}
"banana": {2}
"is": {0, 1, 2}
"it": {0, 1, 2}
"what": {0, 1}
```



# Inverted file

- Can quickly use the inverted file to compute similarity between a new image and all the images in the database
  - Only consider database images whose bins overlap the query image

# Large-scale image search

query image



top 6 results

- Cons:
  - performance degrades as the database grows

# Large-scale image search



- Cons:
  - performance degrades as the database grows



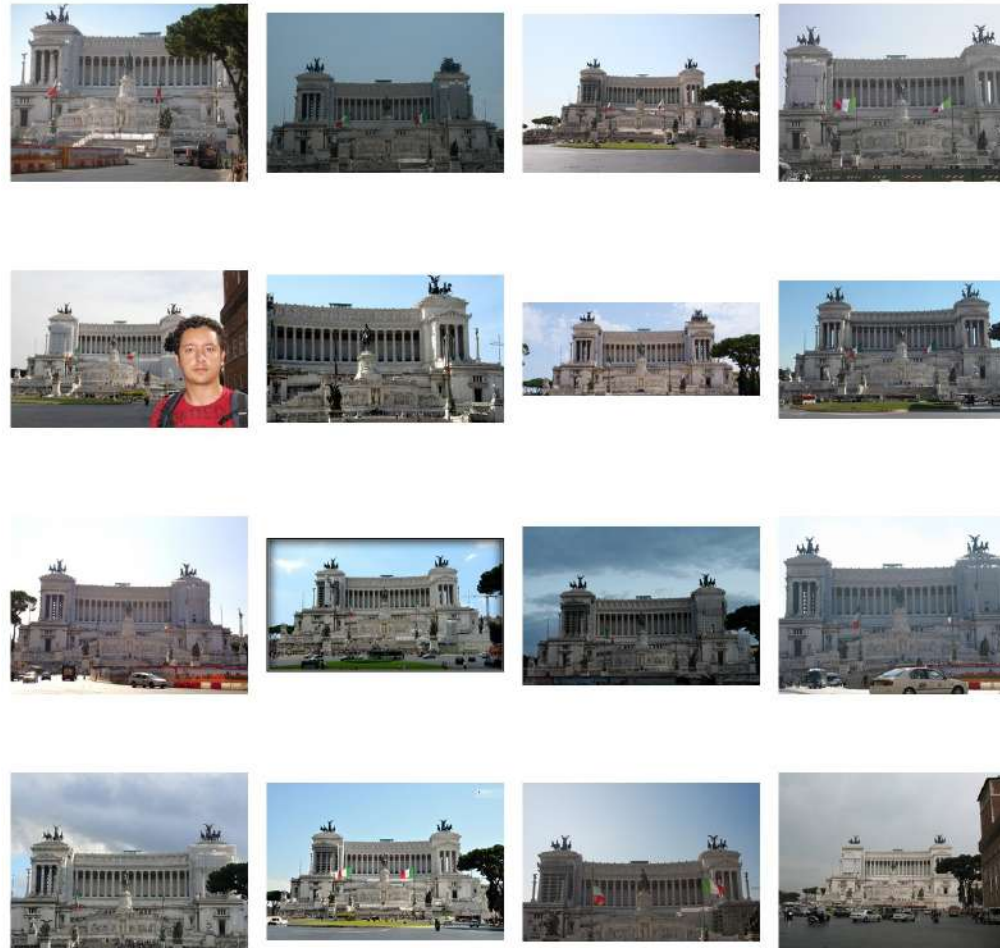
# Large-scale image search

- Pros:
  - Works well for CD covers, movie posters
  - Real-time performance possible

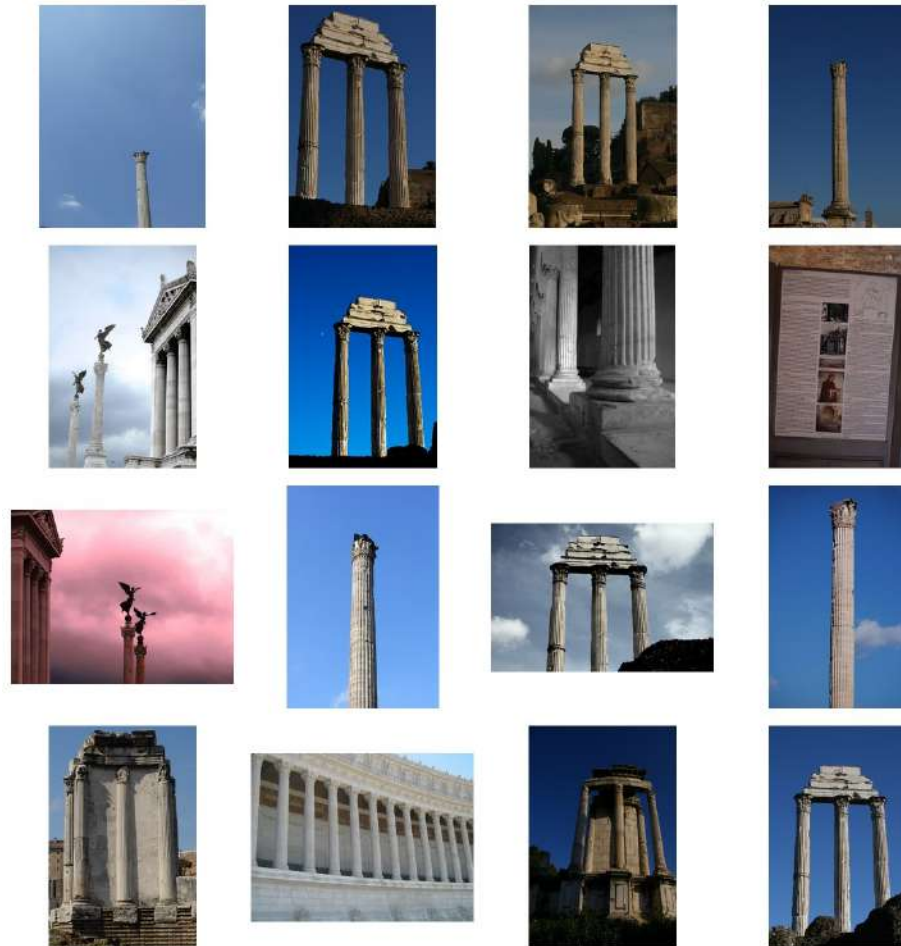


Real-time retrieval from a database of 40,000 CD covers  
Nister & Stewenius, **Scalable Recognition with a Vocabulary Tree**

## Example bag-of-words matches



## Example bag-of-words matches





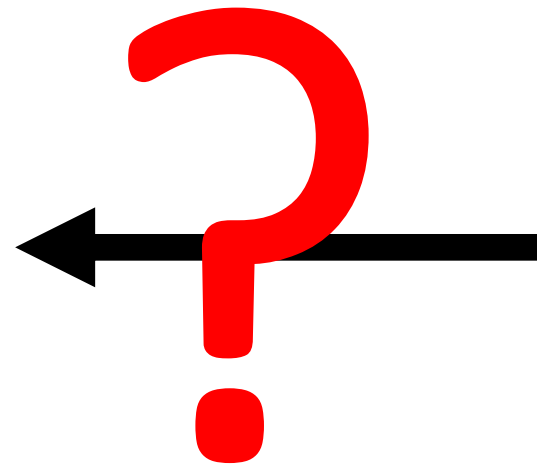
# Today's Agenda

- Visual bag of words (BoW)
  - Background
  - Algorithm
- Applications
  - Image search
  - Action recognition
- Spatial Pyramid Matching





## What about spatial info?



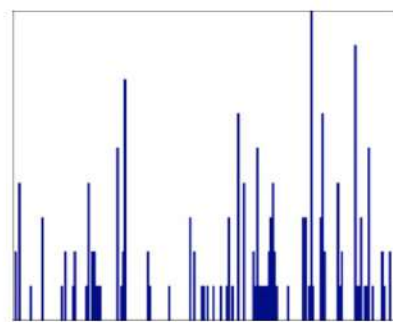
# Pyramids

- Very useful for representing images.
- Pyramid is built by using multiple copies of image.
- Each level in the pyramid is  $1/4$  of the size of previous level.
- The lowest level is of the highest resolution.
- The highest level is of the lowest resolution.

# Bag of words + pyramids



Locally orderless representation at several levels of spatial resolution

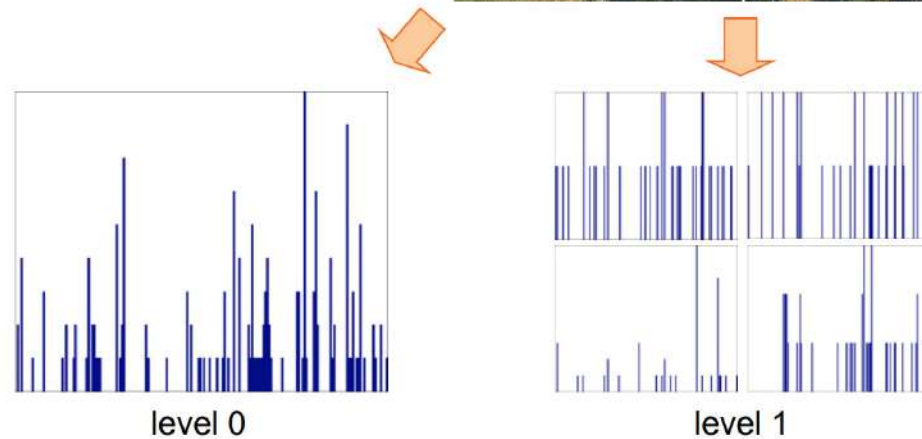


level 0

# Bag of words + pyramids

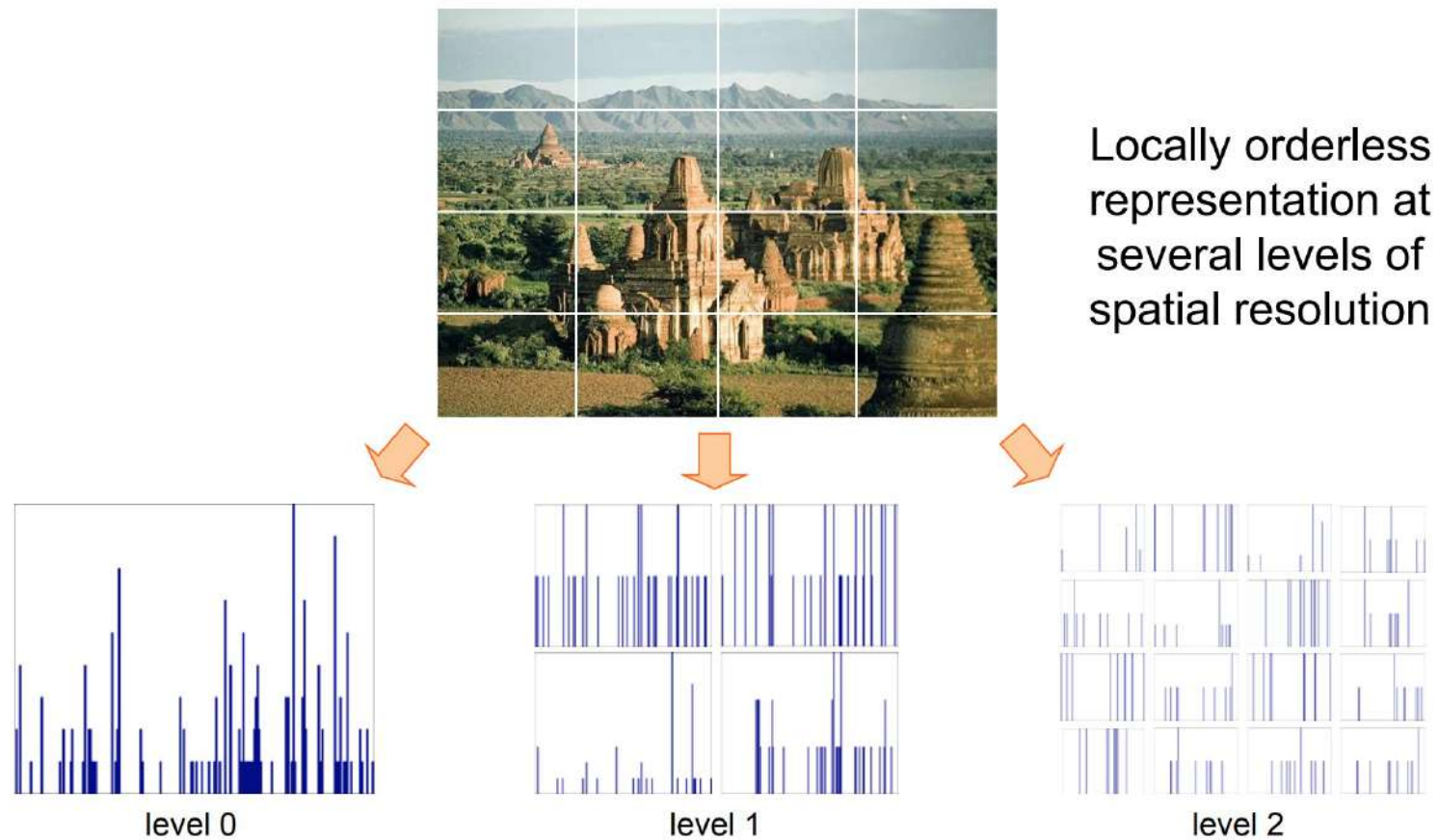


Locally orderless representation at several levels of spatial resolution

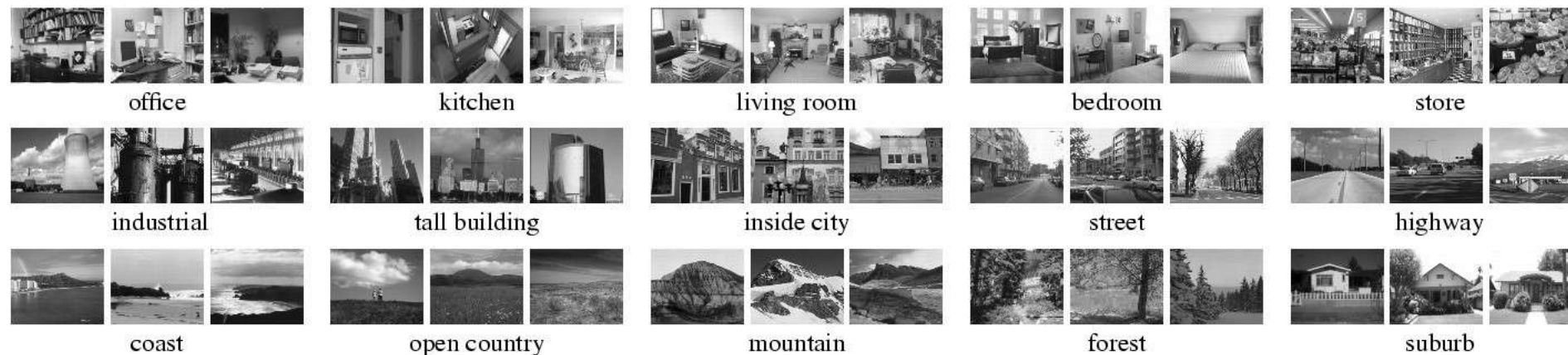




# Bag of words + pyramids



# Results: Scene category dataset



Multi-class classification results (100 training images per class)

Level	Weak features (vocabulary size: 16)		Strong features (vocabulary size: 200)	
	Single-level	Pyramid	Single-level	Pyramid
0 (1 × 1)	45.3 ±0.5		72.2 ±0.6	
1 (2 × 2)	53.6 ±0.3	56.2 ±0.6	77.9 ±0.6	79.0 ±0.5
2 (4 × 4)	61.7 ±0.6	64.7 ±0.7	79.4 ±0.3	<b>81.1 ±0.3</b>
3 (8 × 8)	63.3 ±0.8	<b>66.8 ±0.6</b>	77.2 ±0.4	80.7 ±0.3

Lazebnik, Schmid & Ponce (CVPR 2006)

Slide credit: Svetlana Lazebnik



## Results: Caltech101 dataset

[http://www.vision.caltech.edu/Image\\_Datasets/Caltech101/Caltech101.html](http://www.vision.caltech.edu/Image_Datasets/Caltech101/Caltech101.html)



Multi-class classification results (30 training images per class)

	Weak features (16)		Strong features (200)	
Level	Single-level	Pyramid	Single-level	Pyramid
0	15.5 ±0.9		41.2 ±1.2	
1	31.4 ±1.2	32.8 ±1.3	55.9 ±0.9	57.0 ±0.8
2	47.2 ±1.1	49.3 ±1.4	63.6 ±0.9	<b>64.6 ±0.8</b>
3	52.2 ±0.8	<b>54.0 ±1.1</b>	60.3 ±0.9	64.6 ±0.7

Lazebnik, Schmid & Ponce (CVPR 2006)

Slide credit: Svetlana Lazebnik

**MAI4CAREU**

Master programmes in Artificial  
Intelligence 4 Careers in Europe



# Thank you.







University of Cyprus – MSc Artificial Intelligence

# MAI644 – COMPUTER VISION

## Lecture 13: Object Detection

**Melinos Averkiou**

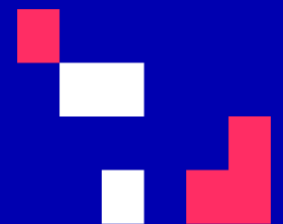
CYENS Centre of Excellence

University of Cyprus - Department of Computer Science

[m.averkiou@cyens.org.cy](mailto:m.averkiou@cyens.org.cy)



**CYENS**  
CENTRE OF EXCELLENCE



## Last time

- Visual bag of words (BoW)
  - Background
  - Algorithm
- Applications
  - Image search
- Spatial Pyramid Matching

# Today's Agenda

- Object detection
  - Task definition
  - Benchmarks
  - Evaluation
- A simple object detector

[material based on Niebles-Krishna]

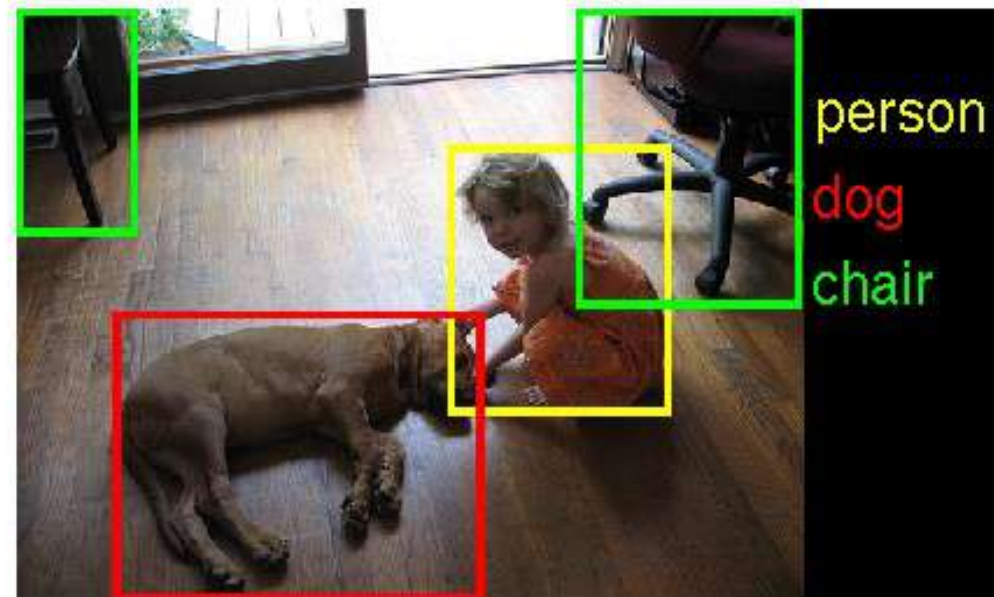
# Today's Agenda

- Object detection
  - Task definition
  - Benchmarks
  - Evaluation
- A simple object detector



# Object Detection

- **Problem:** Detecting and localizing generic objects from various categories, such as cars, people, etc.
- **Challenges:**
  - Illumination,
  - viewpoint,
  - deformations,
  - Intra-class variability





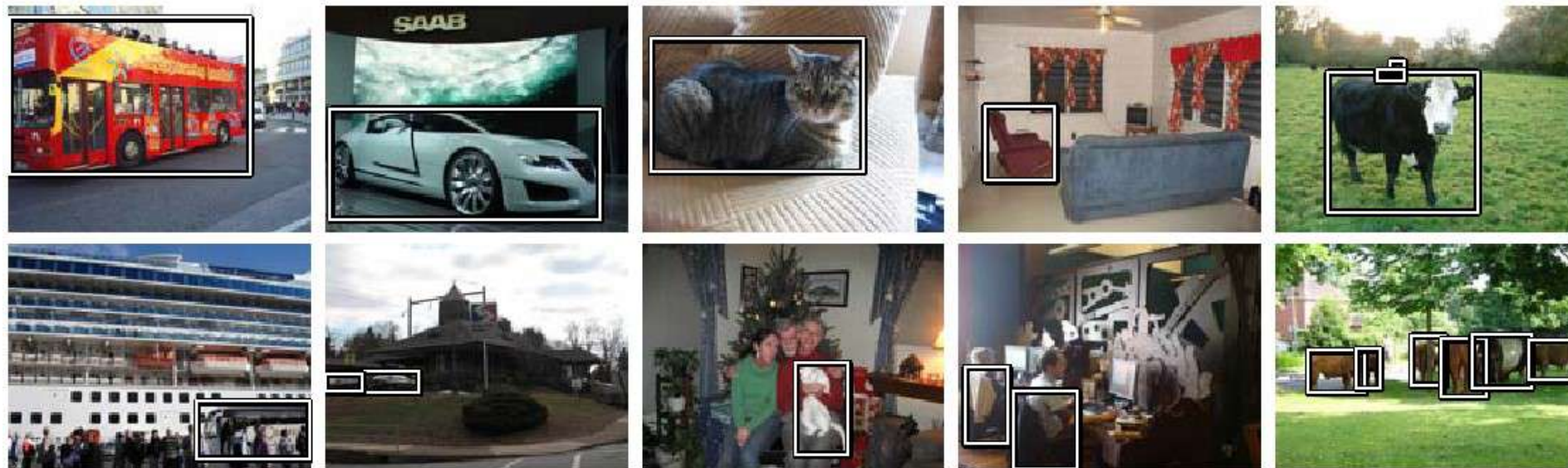
# Today's Agenda

- Object detection
  - Task definition
  - Benchmarks
  - Evaluation
- A simple object detector



# Object Detection Benchmarks

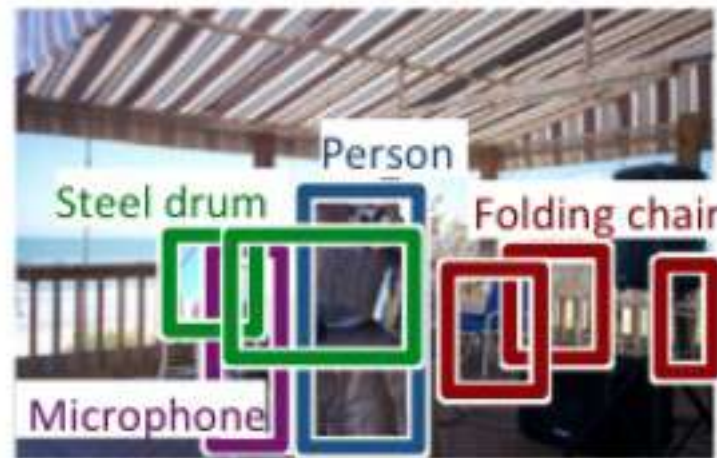
- PASCAL VOC Challenge



- 20 categories
- Annual classification, detection, segmentation challenges

# Object Detection Benchmarks

- PASCAL VOC Challenge
- ImageNet Large Scale Visual Recognition Challenge (ILSVR)
  - 200 Categories for detection





# Object Detection Benchmarks

- PASCAL VOC Challenge
- ImageNet Large Scale Visual Recognition Challenge (ILSVR)
- Common Objects in Context (COCO)
  - 80 Object categories



# Today's Agenda

- Object detection
  - Task definition
  - Benchmarks
  - Evaluation
- A simple object detector
- Deformable parts model
  - Overview
  - Method
  - Pipeline
  - Results and analysis

# How do we evaluate object detection?



— predictions  
 — ground truth

# How do we evaluate object detection?



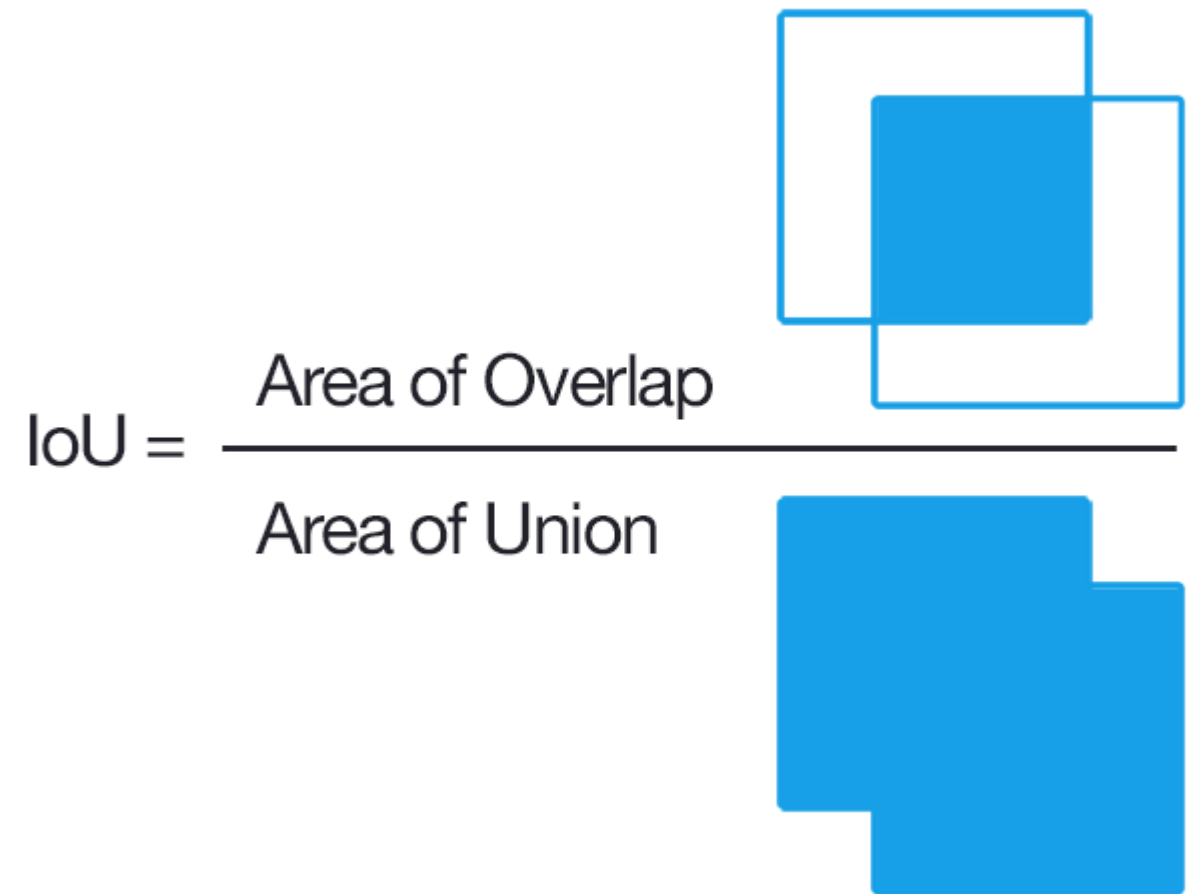
— predictions  
— ground truth

**True positive:**  
- The overlap of the prediction with the ground truth is **MORE** than 0.5



# How do we measure overlap ?

Intersection over Union (IoU)



# How do we evaluate object detection?



— predictions  
 — ground truth

**True positive:**

**False positive:**

- The overlap of the prediction with the ground truth is **LESS** than 0.5

# How do we evaluate object detection?



— predictions  
— ground truth

**True positive:**

**False positive:**

**False negative:**

- The objects that our model doesn't find

# How do we evaluate object detection?



— predictions  
— ground truth

**True positive:**

**False positive:**

**False negative:**

- The objects that our model doesn't find

What is a **True Negative**?



	<u>Predicted 1</u>	<u>Predicted 0</u>
<u>True 1</u>	true positive	false negative
<u>True 0</u>	false positive	true negative



	<u>Predicted 1</u>	<u>Predicted 0</u>
<u>True 1</u>	true positive	false negative
<u>True 0</u>	false positive	true negative

	<u>Predicted 1</u>	<u>Predicted 0</u>
<u>True 1</u>	TP	FN
<u>True 0</u>	FP	TN





	<u>Predicted 1</u>	<u>Predicted 0</u>
<u>True 1</u>	true positive	false negative
<u>True 0</u>	false positive	true negative

	<u>Predicted 1</u>	<u>Predicted 0</u>
<u>True 1</u>	TP	FN
<u>True 0</u>	FP	TN

	<u>Predicted 1</u>	<u>Predicted 0</u>
<u>True 1</u>	hits	misses
<u>True 0</u>	false alarms	correct rejections





	<u>Predicted 1</u>	<u>Predicted 0</u>
<u>True 1</u>	true positive	false negative
<u>True 0</u>	false positive	true negative

	<u>Predicted 1</u>	<u>Predicted 0</u>
<u>True 1</u>	TP	FN
<u>True 0</u>	FP	TN

	<u>Predicted 1</u>	<u>Predicted 0</u>
<u>True 1</u>	hits	misses
<u>True 0</u>	false alarms	correct rejections

$$precision = \frac{TP}{TP + FP}$$

$$recall = \frac{TP}{TP + FN}$$





# How do we evaluate object detection?



— predictions  
— ground truth

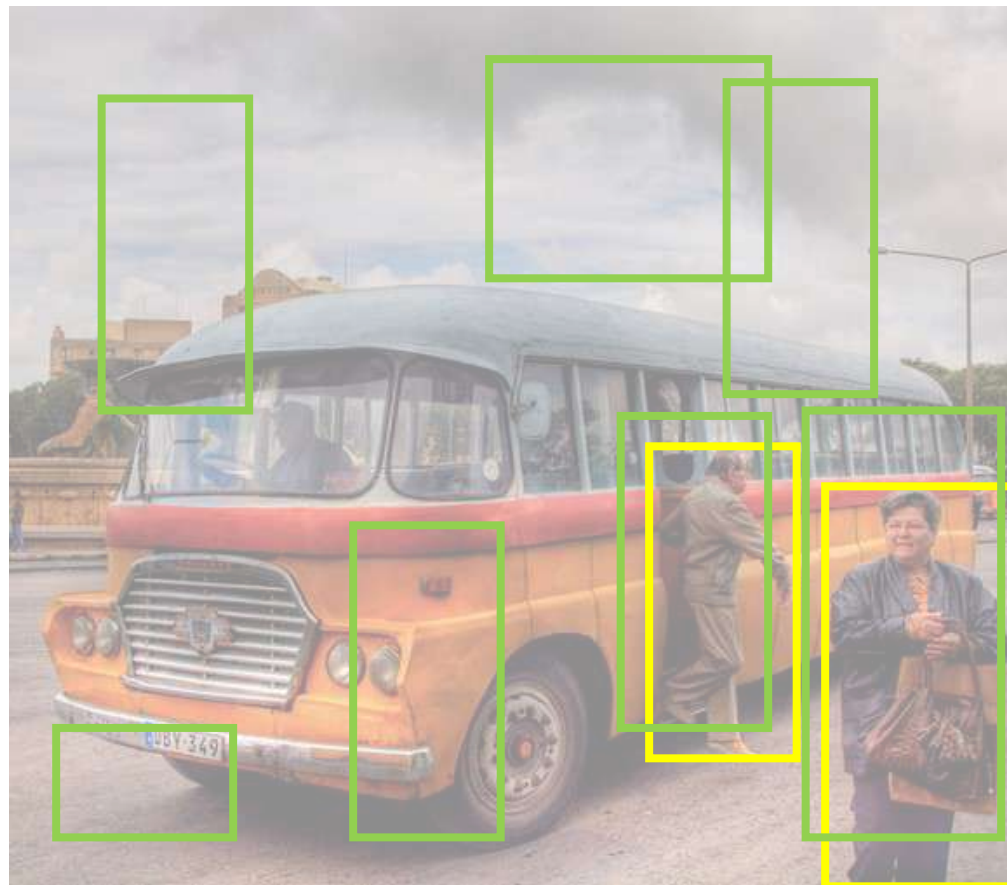
True positive: 1  
False positive: 2  
False negative: 1

So what is the  
- precision?  
- recall?

# Precision versus recall

- Precision:
  - how many of the object detections are correct?
- Recall:
  - how many of the ground truth objects can the model detect?

In reality, our model makes a lot of predictions with varying scores between 0 and 1



- predictions
- ground truth

Here are all the boxes that are predicted with score  $> 0$ .

This means that our:

- Recall is perfect!
- But our precision is BAD!

In reality, our model makes a lot of predictions with varying scores between 0 and 1



\_\_\_\_\_ predictions  
 \_\_\_\_\_ ground truth

There are no boxes that are predicted with **score = 1**.

- This means that our
- Precision is undefined!
  - And our recall is BAD!



# How do we evaluate object detection?

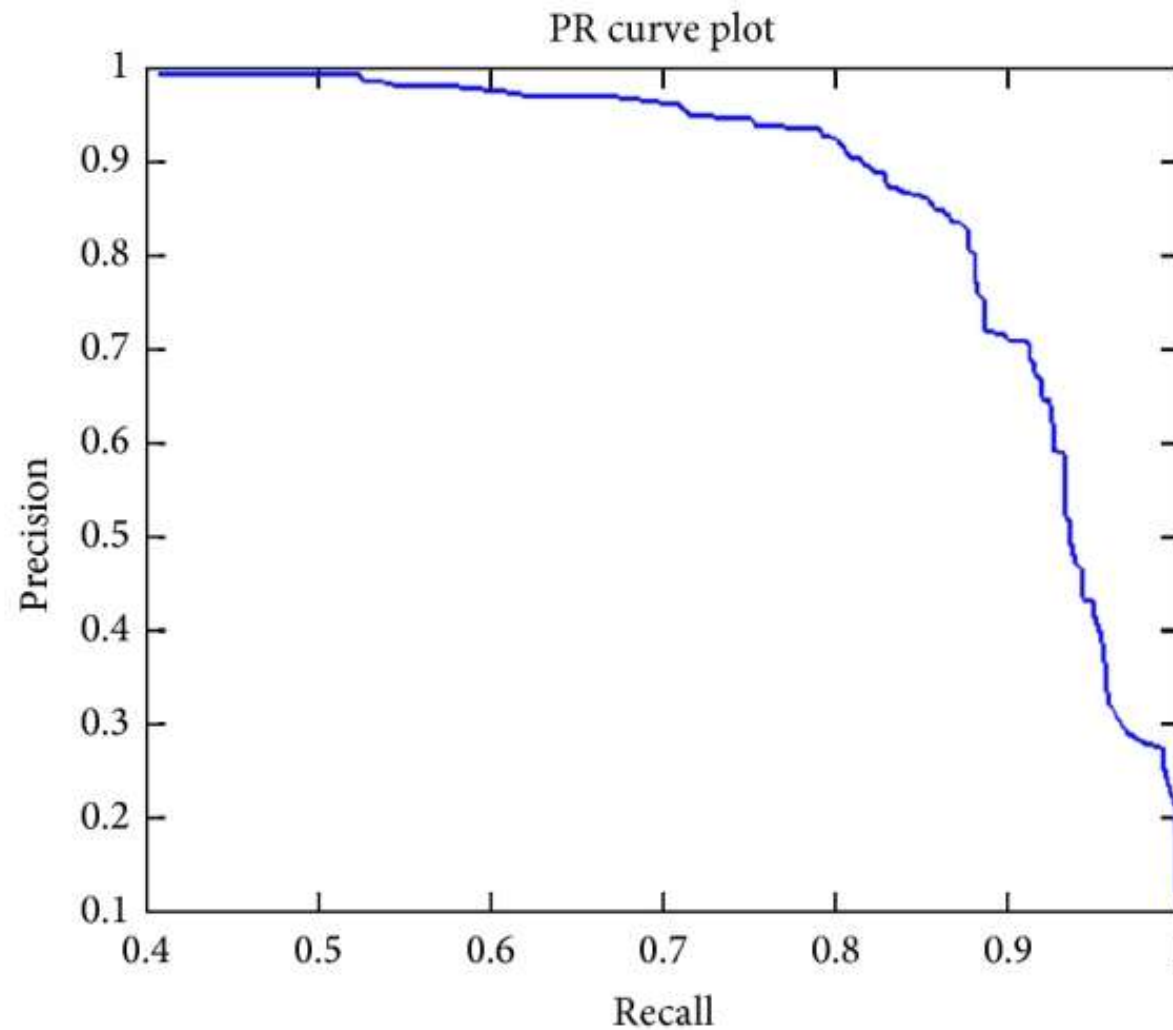


— predictions  
— ground truth

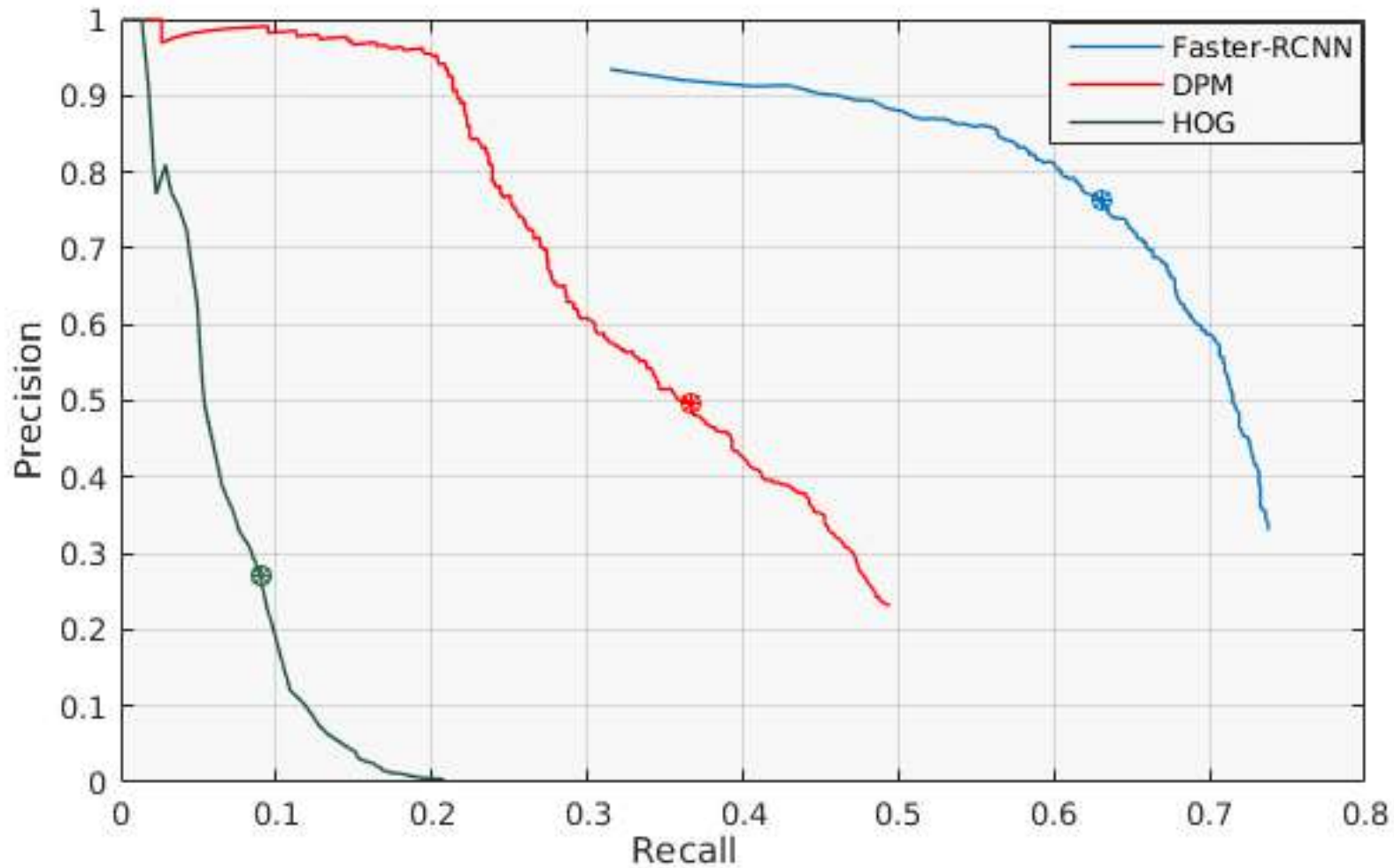
Here are all the boxes  
that are predicted with  
score  $> 0.5$

We are setting a  
threshold of 0.5

# Precision – recall curve (PR curve)



# Which model is the best?



## True Positives - Person

UoCTTI\_LSVN-MDPM



MIZZOU\_DEF-HOG-LBP



NECUIUC\_CLS-DTCT





## False Positives - Person

UoCTTI\_LSVN-MDPM



MIZZOU\_DEF-HOG-LBP



NECUIUC\_CLS-DTCT





## “Near Misses” - Person

UoCTTI\_LSYM-MDPM



MIZZOU\_DEF-HOG-LBP



NECUIUC\_CLS-DTCT





## True Positives - Bicycle

UoCTTI\_L SVM-MDPM



OXFORD\_MKL



NECUIUC\_CLS-DTCT





## False Positives - Bicycle

UoCTTI\_L SVM-MDPM



OXFORD\_MKL



NECUIUC\_CLS-DTCT





# Today's Agenda

- Object detection
  - Task definition
  - Benchmarks
  - Evaluation
- A simple object detector

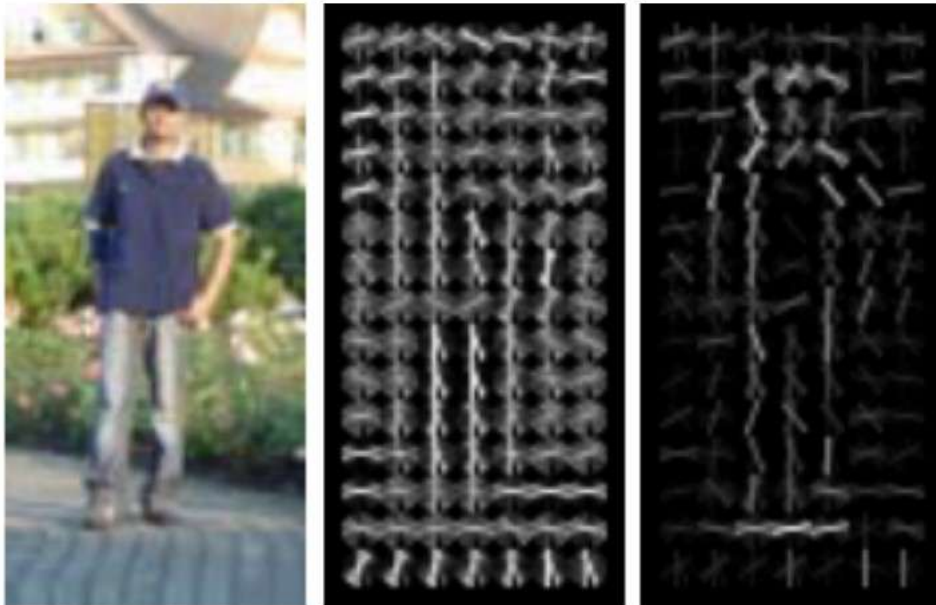
# Dalal-Triggs method



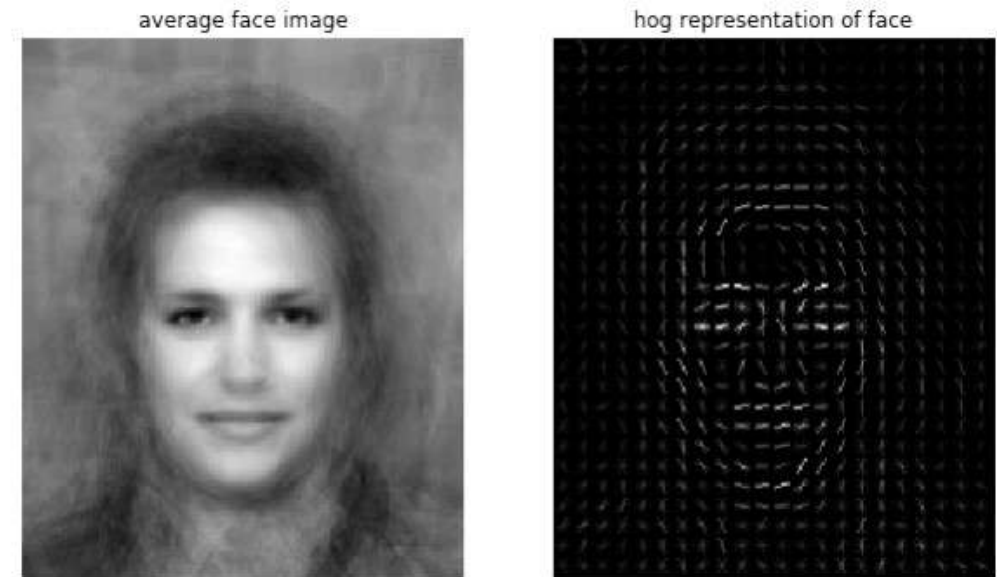
sliding window

# Recap – HoG features

Find a HoG template and use as filter



Train a linear SVM classifier on HoG



Take the average image and extract HoG

# Sliding window + HoG features



No person here

- Slide through the image and check if there is an object at every location
- Compare HOG feature template to HOG features from each location in the image.
  - Use dot product



# Sliding window + HoG features



YES!! Person match found

- Slide through the image and check if there is an object at every location
- Compare HOG feature template to HOG features from each location in the image.
- If a comparison produces a high score, output detection at the corresponding location



# Sliding window + HoG features



- But what if we were looking for buses?

No bus found





# Sliding window + HoG features



- But what if we were looking for buses?

No bus found





# Sliding window + HoG features



- We will never find the object if we don't choose our window size wisely!

No bus found



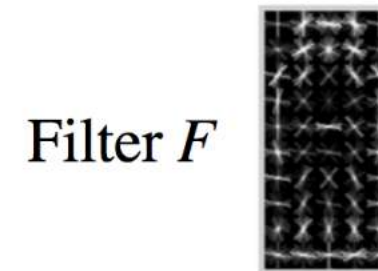
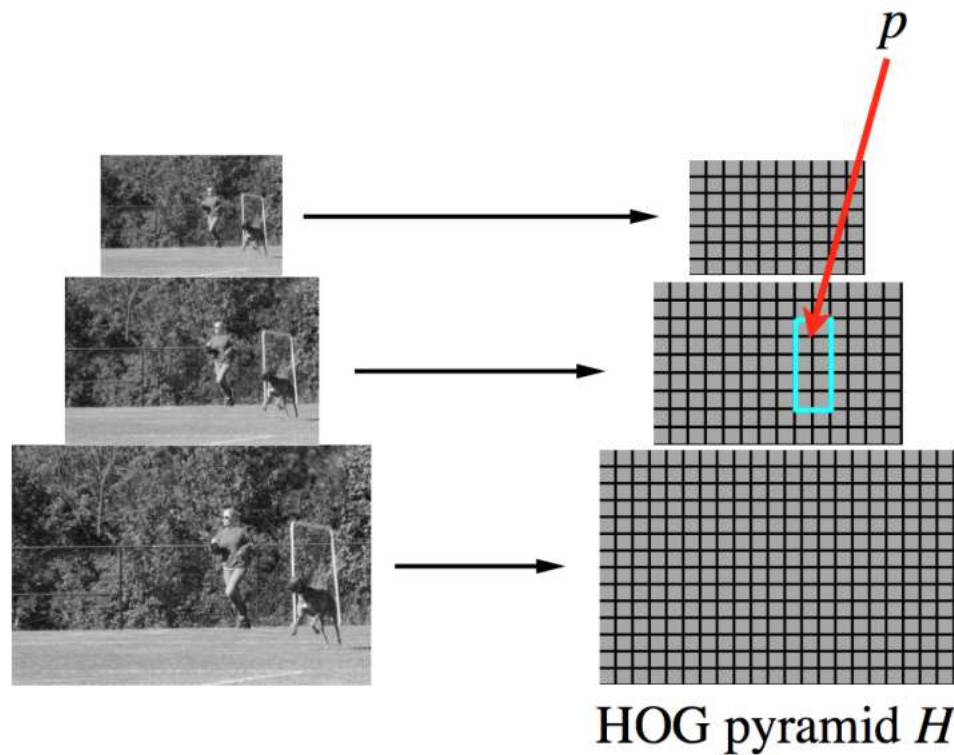


# Sliding window + HoG features



We need to do a **multi scale** sliding window search

# Create a feature pyramid



Score of  $F$  at position  $p$  is  
 $F \cdot \phi(p, H)$

$\phi(p, H)$  = concatenation of  
 HOG features from  
 subwindow specified by  $p$

**MAI4CAREU**

Master programmes in Artificial  
Intelligence 4 Careers in Europe



# Thank you.





University of Cyprus – MSc Artificial Intelligence

# MAI644 – COMPUTER VISION

## Lecture 14: Camera Models

**Melinos Averkiou**

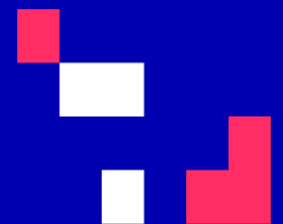
CYENS Centre of Excellence

University of Cyprus - Department of Computer Science

[m.averkiou@cyens.org.cy](mailto:m.averkiou@cyens.org.cy)



**CYENS**  
CENTRE OF EXCELLENCE





## Last time

- Object detection
  - Task definition
  - Benchmarks
  - Evaluation
- A simple object detector

# Today's Agenda

- Perspective projection
- Vanishing points
- Full camera model

[material based on Paris Kaimakis]



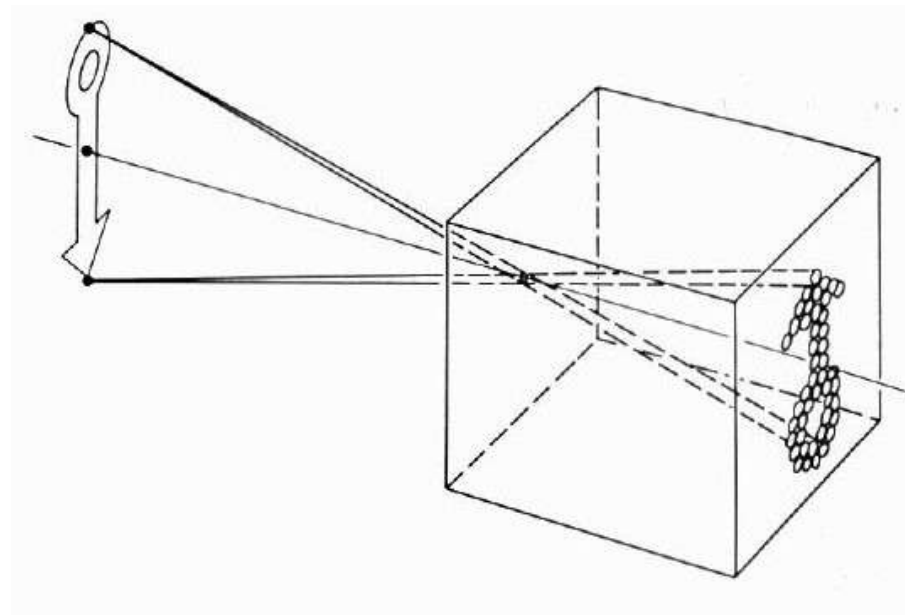
# Today's Agenda

- Perspective projection
- Vanishing points
- Full camera model



# Perspective Projection

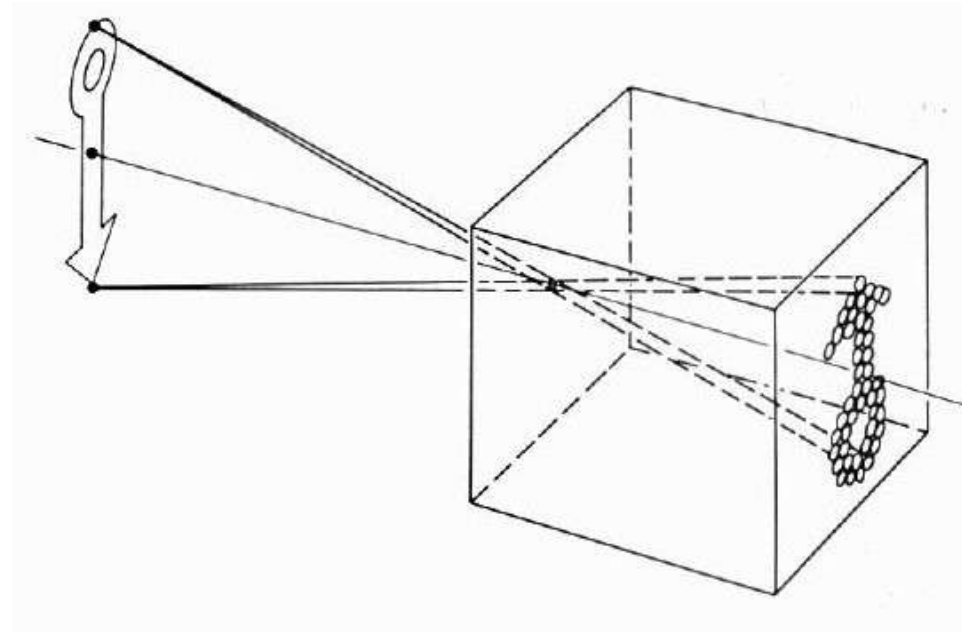
- The process of capturing 2D images of the 3D world is modelled by **perspective projection**
- We model perspective projection using the **pin-hole camera**





# Pinhole Camera model

- Captures **pencil of rays** – all rays through a single point
- The point is called **Center of Projection**
- The image is formed on the **Image Plane**

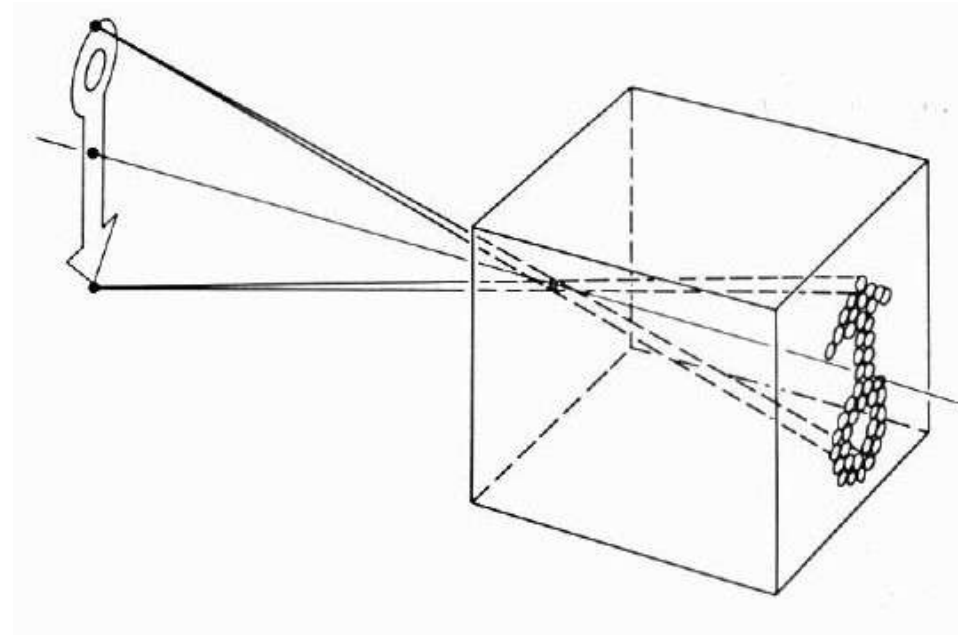


[Slide by Steve Seitz]

# Pinhole Camera - Camera Obscura

## The first camera

- Known to Aristotle
- How does the aperture size affect the image?

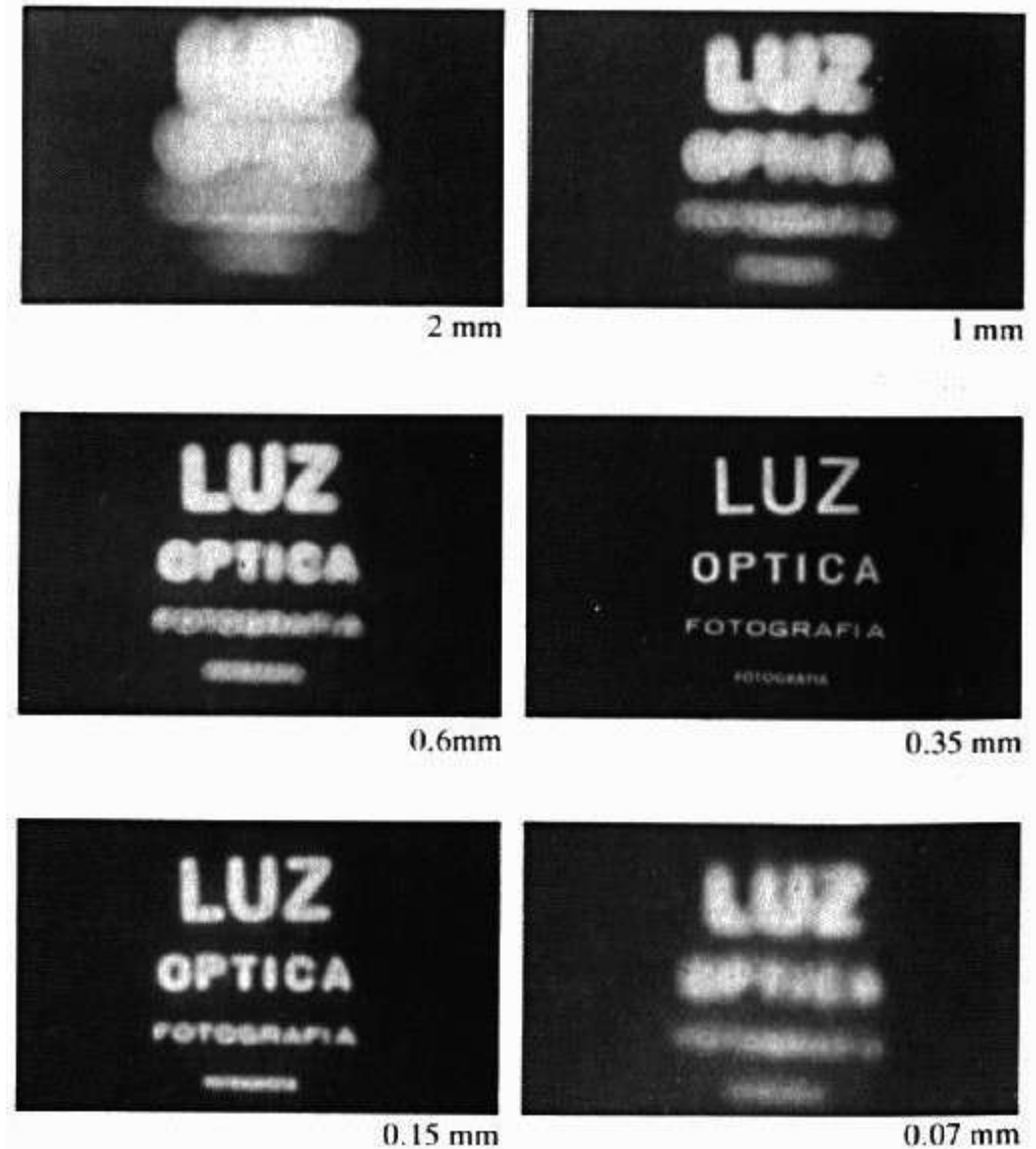
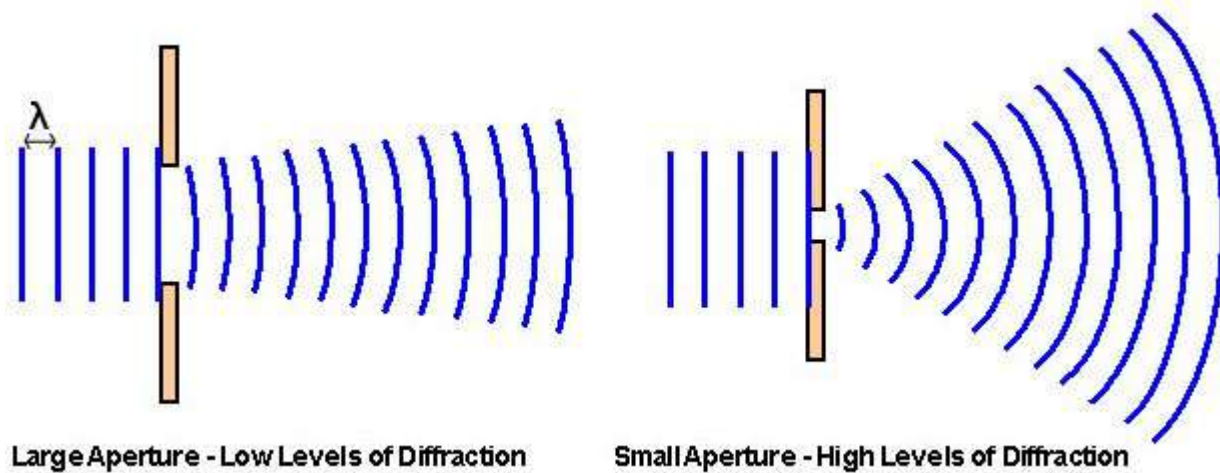


[Slide by Steve Seitz]

# Shrinking the aperture

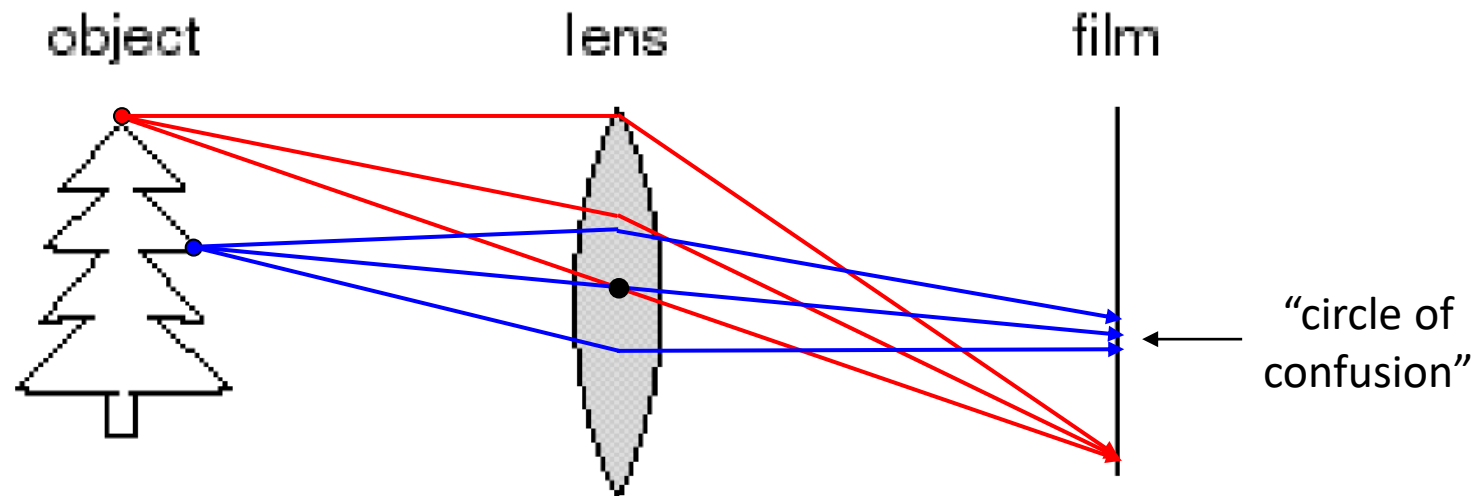
Why not make the aperture as small as possible?

- Less light gets through
- Diffraction effects



[Slide by Steve Seitz]

# Adding a lens



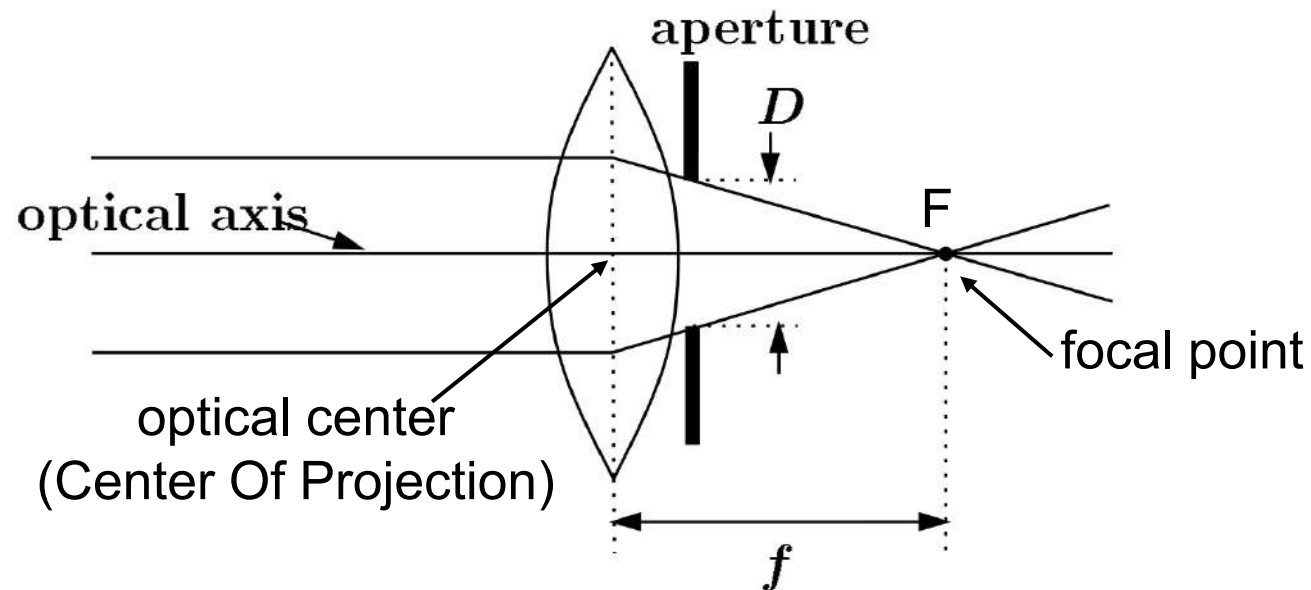
A lens focuses light onto the film

- Rays passing through the center are not deviated
- There is a specific distance at which objects are “in focus”
  - other points project to a “circle of confusion” in the image

[Slide by Steve Seitz]



# Lenses

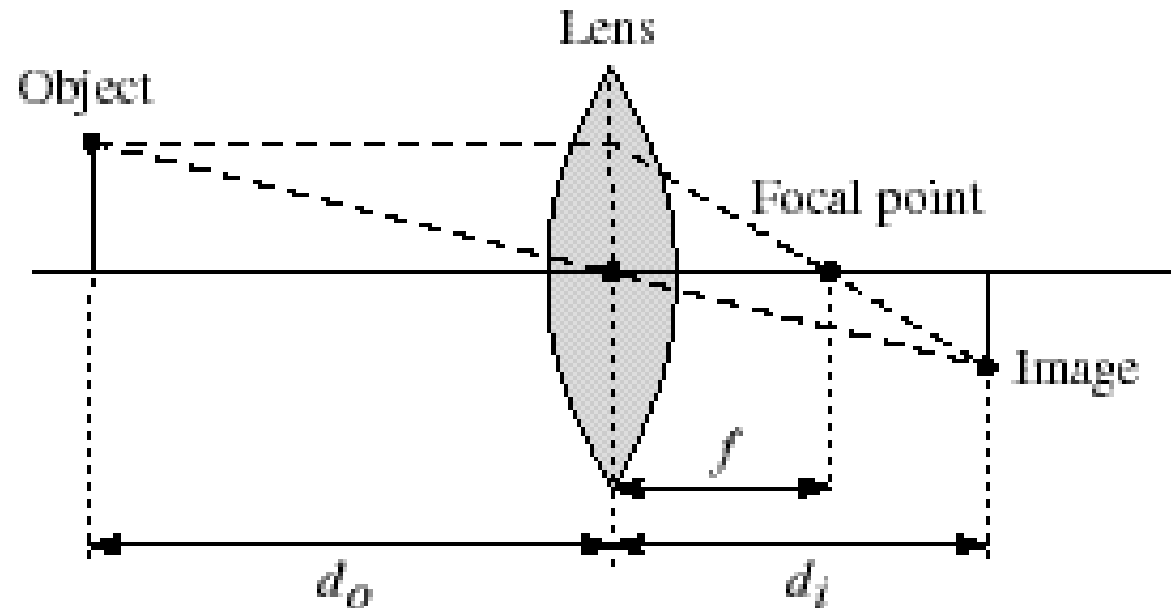


A lens focuses parallel rays onto a single **focal point**

- Focal point is on a plane located at a distance  $f$  (**focal length**) beyond the plane of the lens
  - $f$  is a function of the shape and index of refraction of the lens
- Aperture of diameter  $D$  restricts the range of rays
  - aperture may be on either side of the lens
- Lenses are typically spherical (easier to produce)

[Slide by Steve Seitz]

# Thin lenses



Thin lens equation: 
$$\frac{1}{d_o} + \frac{1}{d_i} = \frac{1}{f}$$

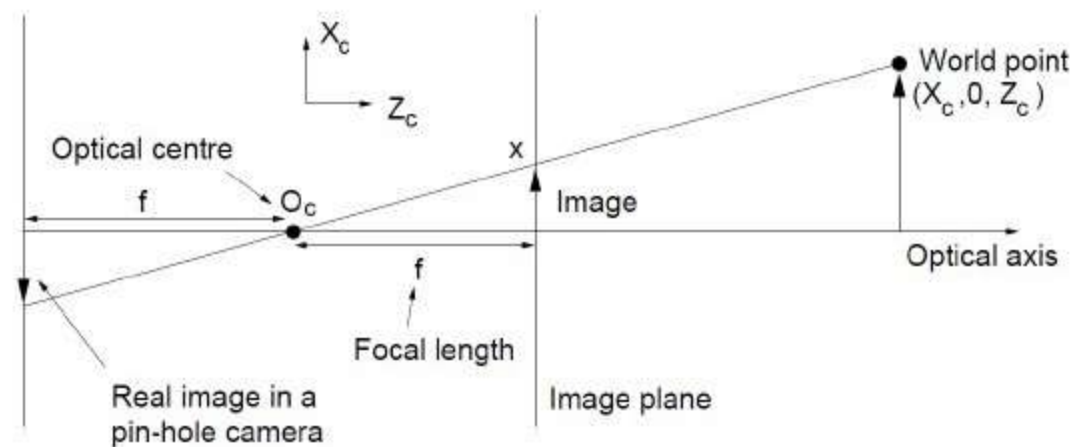
- Any object point satisfying this equation is in focus
- Thin lens applet (needs flash player):

[http://www.phy.ntnu.edu.tw/java/Lens/lens\\_e.html](http://www.phy.ntnu.edu.tw/java/Lens/lens_e.html) (by Fu-Kwun Hwang)

[Slide by Steve Seitz]

# Perspective Projection

- Objects closer to the camera appear larger, those far away from the camera appear smaller
- We can place the image plane in front on the camera to avoid flipping

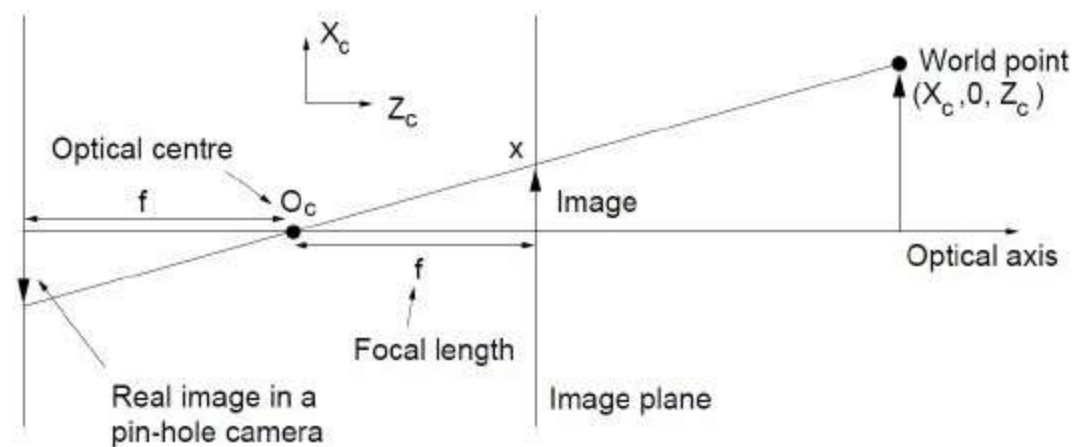


# Perspective Projection equation

- Using similar triangles we find that:

$$\frac{x}{f} = \frac{X_c}{Z_c} \iff x = \frac{f}{Z_c} X_c$$

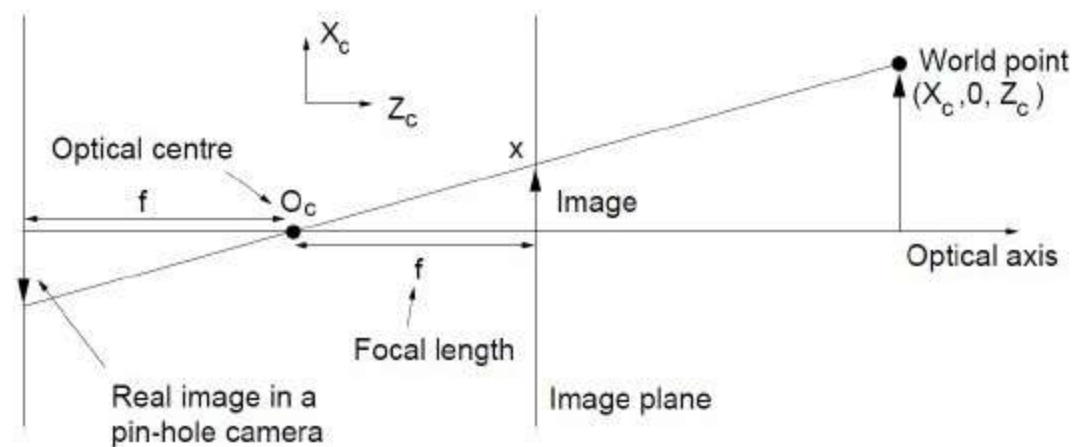
- Similarly





# Perspective Projection equation

- Therefore, under perspective projection every 3D point  $\mathbf{X}(X_C, Y_C, Z_C)$  projects to a 2D point  $\mathbf{x}(\frac{f}{Z_C} X_C, \frac{f}{Z_C} Y_C)$  on the camera's image plane

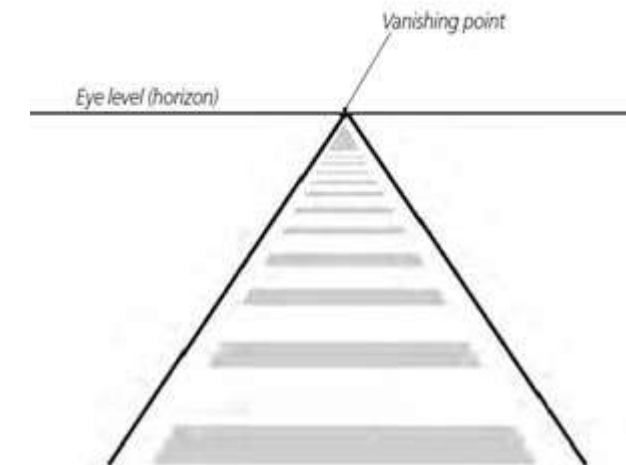
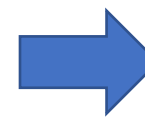
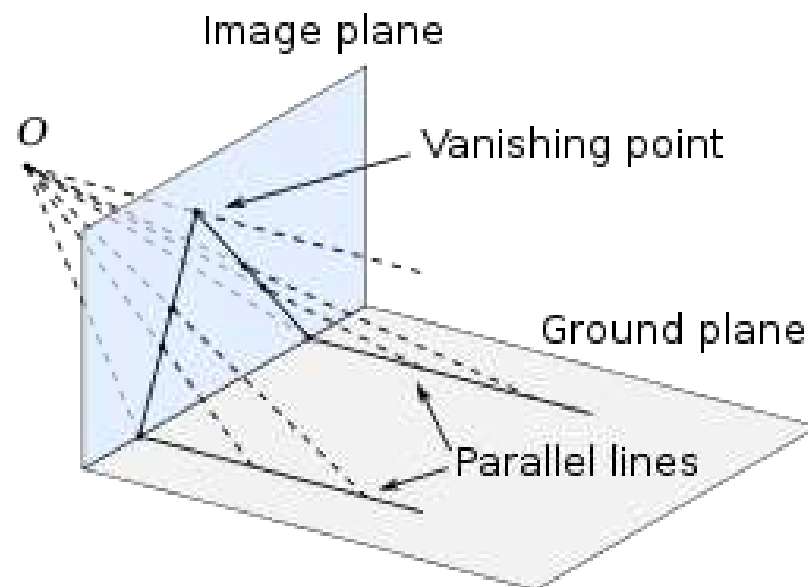


# Today's Agenda

- Perspective projection
- Vanishing points
- Full camera model

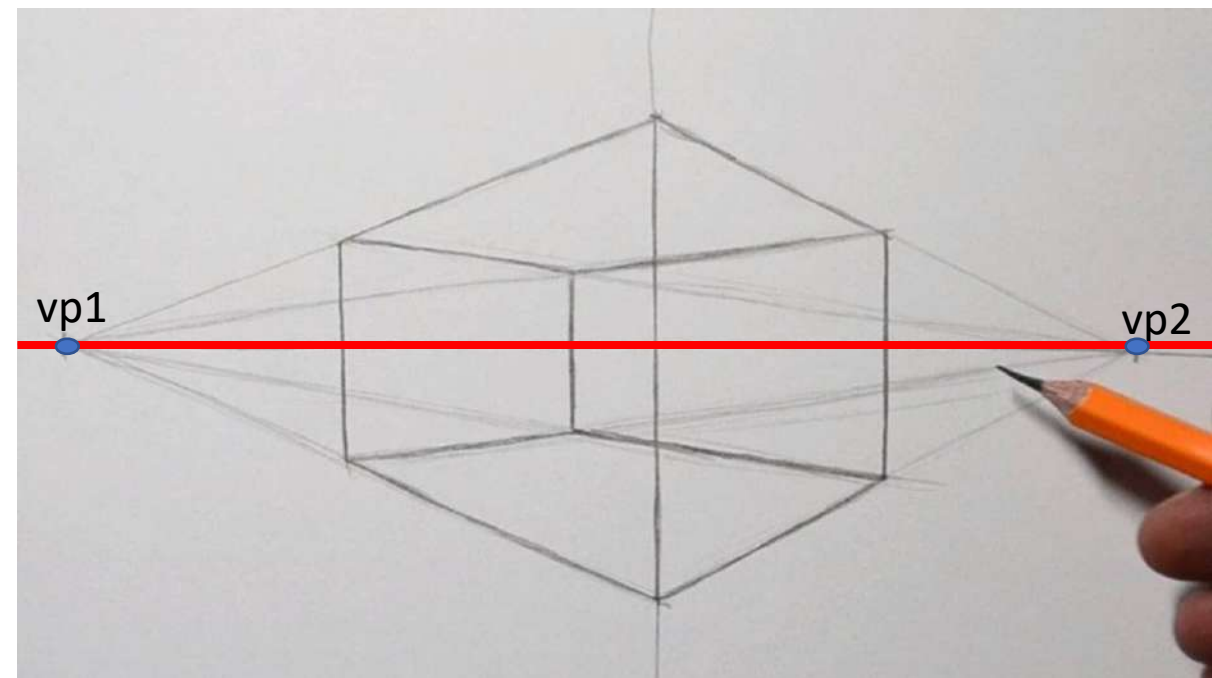
# Vanishing points

- Vanishing points are points in the image where parallel lines appear to meet.
- Each set of parallel lines in the world (that are not parallel to the image plane) will have a different vanishing point in the image.



# Vanishing points

- In the same way, parallel planes in the world meet in a line in the image, often called a **horizon line**.
- Any set of parallel lines lying on these planes in the 3D world will have a vanishing point on the horizon line.





# Vanishing points

Heavily used in Renaissance paintings to introduce 'depth' into the scene



*Pietro Perugino, Christ Handing the Keys to St. Peter, 1481-1482. Image via [Wikimedia Commons](https://commons.wikimedia.org/wiki/File:Pietro_Perugino_-_Christ_Handing_the_Keys_to_Saint_Peter_-_WGA01499.jpg).*

# Vanishing point example

- Let's find the image location of the vanishing point for a line in the world

$$X_c = \mathbf{a} + \lambda \mathbf{b} = (a_x + \lambda b_x, a_y + \lambda b_y, a_z + \lambda b_z)$$

$$\Rightarrow \mathbf{x} = \frac{f}{Z_c} (X_c, Y_c) = f \left[ \frac{a_x + \lambda b_x}{a_z + \lambda b_z}, \frac{a_y + \lambda b_y}{a_z + \lambda b_z} \right]$$

- As  $\lambda$  goes to infinity we move down the line and  $\mathbf{x}$  converges to the vanishing point

$$\mathbf{x}_{vp} = \lim_{\lambda \rightarrow \infty} \mathbf{x} = f \left[ \frac{b_x}{b_z}, \frac{b_y}{b_z} \right]$$

# Vanishing point example

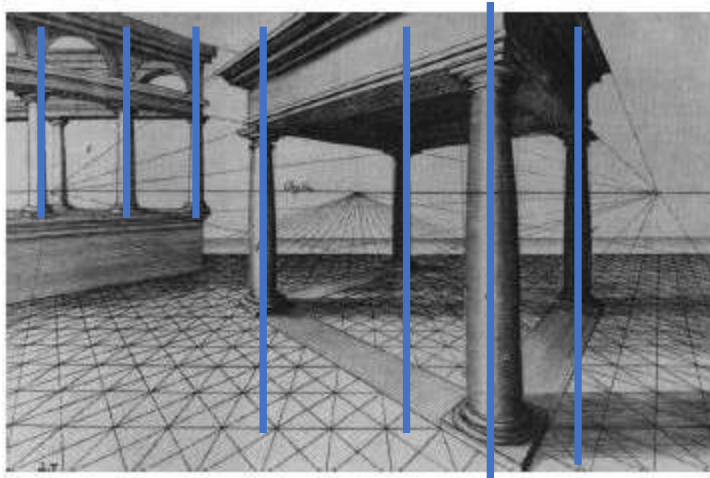
- As  $\lambda$  goes to infinity we move down the line and  $x$  converges to the vanishing point

$$x_{vp} = \lim_{\lambda \rightarrow \infty} x = f \left[ \frac{b_x}{b_z}, \frac{b_y}{b_z} \right]$$

- The vanishing point depends only on the line's orientation (the vector  $b$  in the line equation), and not its position.
- When  $b_z=0$ , the line is parallel to the image plane, and the vanishing point is at infinity. Remember:
  - Each set of parallel lines in the world (**that are not parallel to the image plane**) will have a different vanishing point in the image.

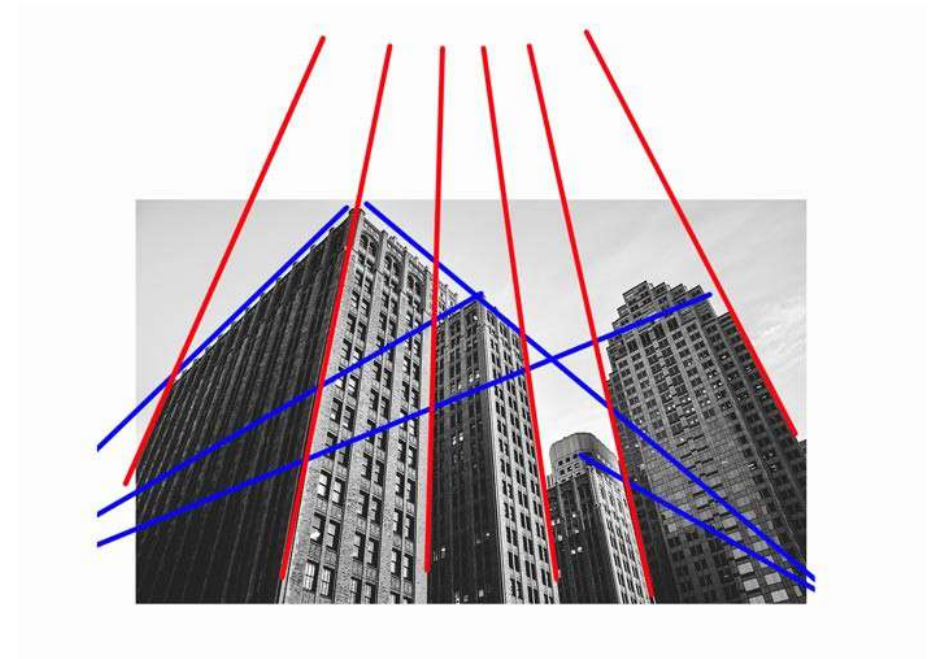


# Vanishing point example



$b_z=0$

VS



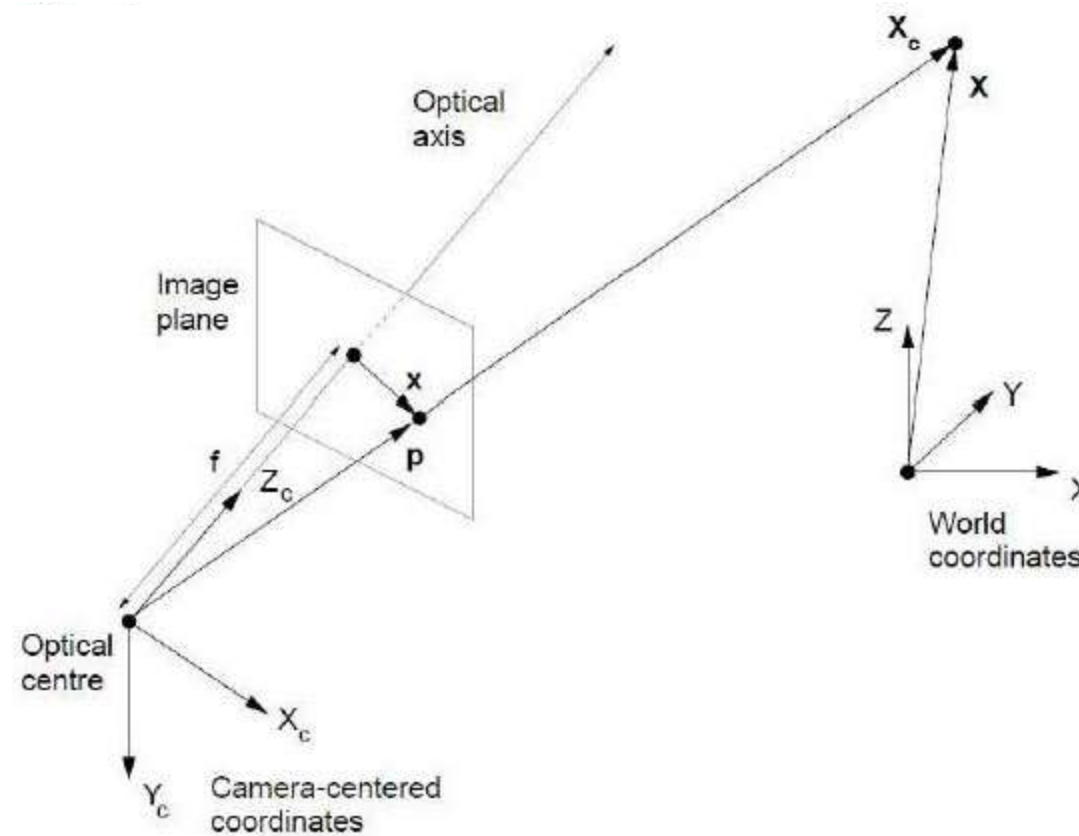


# Today's Agenda

- Perspective projection
- Vanishing points
- Full camera model

# Full camera model

A full camera model describes the mapping from 3D world to 2D pixel coordinates



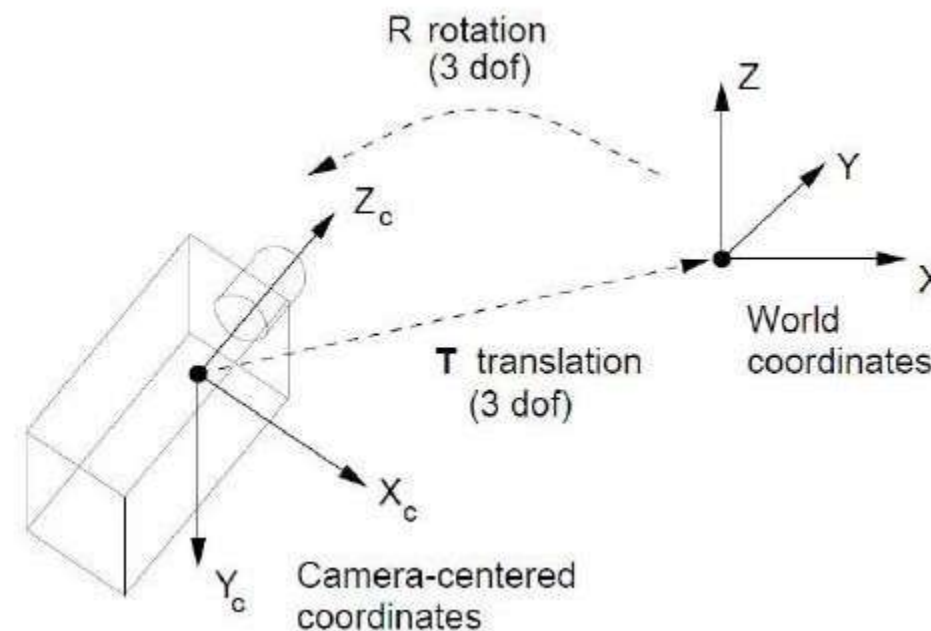
# Full camera model

It consists of three transformations:

1. The **Euclidean (Rigid) transformation** between the camera and the world, i.e., the translation and rotation of the camera with respect the world origin – takes points from 3D world coordinates to 3D camera coordinates
2. The **perspective projection** onto the camera plane – takes points from 3D camera coordinates to 2D image coordinates
3. **CCD imaging**, i.e., the geometry of the CCD array (the size and shape of the pixels) and its position with respect to the optical axis – takes points from 2D image coordinates to 2D pixel coordinates

# Full camera model – Euclidean transformation

We attach a coordinate system  $\mathbf{X}(X, Y, Z)$  to the world and another coordinate system to the camera  $\mathbf{X}_C(X_C, Y_C, Z_C)$



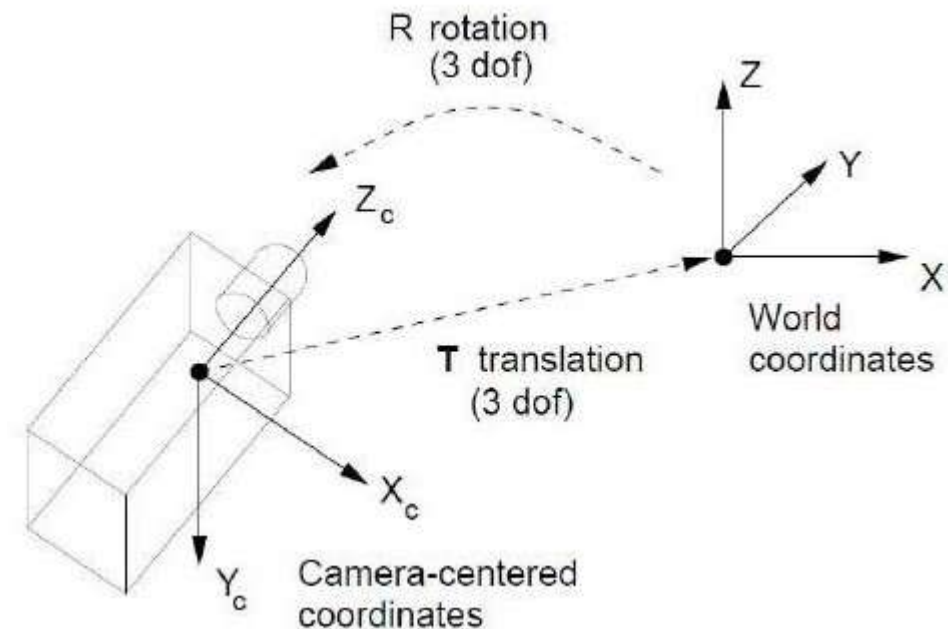


# Full camera model – Euclidean transformation

The Euclidean transformation can be described by a rotation matrix R and a translation vector T

$$X_c = RX + T$$

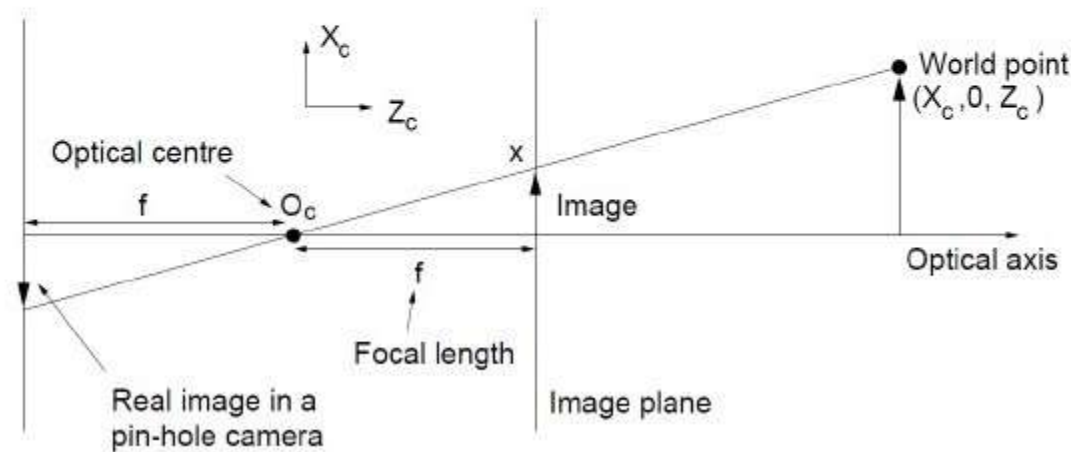
$$\begin{bmatrix} X_c \\ Y_c \\ Z_c \end{bmatrix} = \begin{bmatrix} r_{11} & r_{12} & r_{13} \\ r_{21} & r_{22} & r_{23} \\ r_{31} & r_{32} & r_{33} \end{bmatrix} \begin{bmatrix} X \\ Y \\ Z \end{bmatrix} + \begin{bmatrix} T_x \\ T_y \\ T_z \end{bmatrix}$$



# Full camera model – Perspective projection

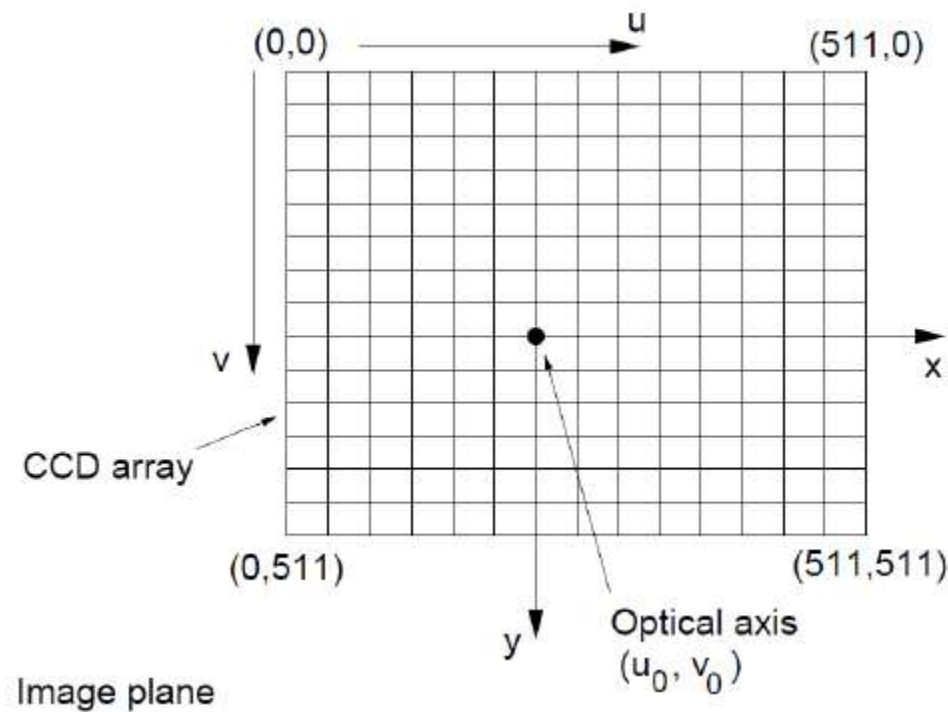
Perspective projection is modelled by

$$\left. \begin{aligned} x &= \frac{f}{Z_c} X_c \\ y &= \frac{f}{Z_c} Y_c \end{aligned} \right\}$$



# Full camera model – CCD Imaging

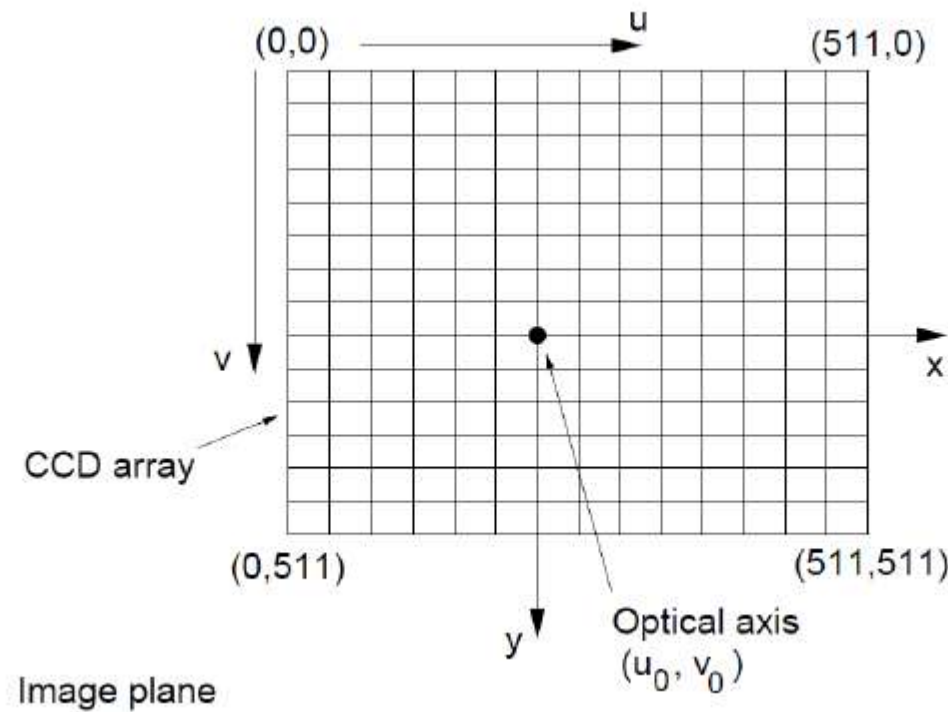
To model CCD Imaging we define pixel coordinates  $\mathbf{w}(u,v)$  in addition to the image coordinates  $\mathbf{x}(x,y)$



# Full camera model – CCD Imaging

$w(u,v)$  and  $x(x,y)$  are related as follows

$$\left. \begin{aligned} u &= u_0 + k_u x \\ v &= v_0 + k_v y \end{aligned} \right\}$$





# Full camera model – Putting it all together

So here is the full camera model combining all three transformations

$$\left. \begin{aligned} u &= u_0 + k_u x \\ v &= v_0 + k_v y \end{aligned} \right\} \Rightarrow \left. \begin{aligned} u &= u_0 + k_u \frac{f}{Z_c} X_c \\ v &= v_0 + k_v \frac{f}{Z_c} Y_c \end{aligned} \right\}$$

$$\Rightarrow \left. \begin{aligned} u &= u_0 + k_u f \frac{r_{11}X + r_{12}Y + r_{13}Z + T_x}{r_{31}X + r_{32}Y + r_{33}Z + T_z} \\ v &= v_0 + k_v f \frac{r_{21}X + r_{22}Y + r_{23}Z + T_y}{r_{31}X + r_{32}Y + r_{33}Z + T_z} \end{aligned} \right\}$$

**MAI4CAREU**

Master programmes in Artificial  
Intelligence 4 Careers in Europe



# Thank you.





University of Cyprus – MSc Artificial Intelligence

# MAI644 – COMPUTER VISION

## Lecture 15: Camera Calibration

**Melinos Averkiou**

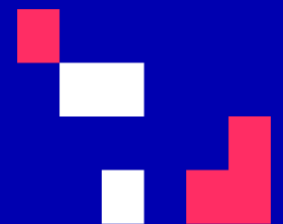
CYENS Centre of Excellence

University of Cyprus - Department of Computer Science

[m.averkiou@cyens.org.cy](mailto:m.averkiou@cyens.org.cy)



**CYENS**  
CENTRE OF EXCELLENCE



## Last time

- Perspective projection
- Vanishing points
- Full camera model



# Today's Agenda

- Full camera model in matrix form
- Camera calibration
- Calibration – Projective camera model
- Calibration – Affine camera model

[material based on Paris Kaimakis]



# Today's Agenda

- Full camera model in matrix form
- Camera calibration
- Calibration - Projective camera model
- Calibration – Affine camera model



# Remember - Full camera model

It consists of three transformations:

1. The **Euclidean (Rigid) transformation** between the camera and the world, i.e., the translation and rotation of the camera with respect the world origin – takes points from 3D world coordinates to 3D camera coordinates
2. The **perspective projection** onto the camera plane – takes points from 3D camera coordinates to 2D image coordinates
3. **CCD imaging**, i.e., the geometry of the CCD array (the size and shape of the pixels) and its position with respect to the optical axis – takes points from 2D image coordinates to 2D pixel coordinates

# Remember - Full camera model – Putting it all together

So here is the full camera model combining all three transformations

$$\left. \begin{aligned} u &= u_0 + k_u x \\ v &= v_0 + k_v y \end{aligned} \right\} \Rightarrow \left. \begin{aligned} u &= u_0 + k_u \frac{f}{Z_c} X_c \\ v &= v_0 + k_v \frac{f}{Z_c} Y_c \end{aligned} \right\}$$

$$\Rightarrow \left. \begin{aligned} u &= u_0 + k_u f \frac{r_{11}X + r_{12}Y + r_{13}Z + T_x}{r_{31}X + r_{32}Y + r_{33}Z + T_z} \\ v &= v_0 + k_v f \frac{r_{21}X + r_{22}Y + r_{23}Z + T_y}{r_{31}X + r_{32}Y + r_{33}Z + T_z} \end{aligned} \right\}$$



# Full camera model in matrix form - Euclidean

- How can we express all three transformations in the form of matrices ?
  - Homogeneous coordinates
- Let's start from the Euclidean (rigid) transformation connecting the world and camera coordinate systems.
- This is composed of a rotation (3DOF) expressing the camera pose with respect to the world coordinate frame, and a translation (3DOF) expressing the camera location with respect to the world origin.

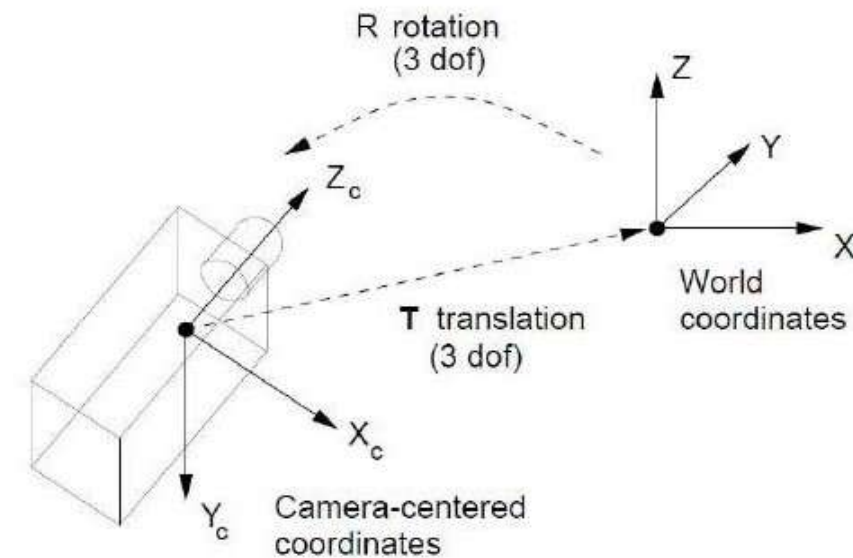
# Full camera model in matrix form - Euclidean

- We want to express a ‘world’ 3D point  $\mathbf{X}(X, Y, Z)$  as a 3D point  $\mathbf{X}_C(X_C, Y_C, Z_C)$  in camera coordinates
- In homogeneous coordinates we have  $\tilde{\mathbf{X}}(\lambda X, \lambda Y, \lambda Z, \lambda)$  and  $\tilde{\mathbf{X}}_C(\lambda X_C, \lambda Y_C, \lambda Z_C, \lambda)$  –  $\lambda$  can be set to 1 as it has no effect on  $\mathbf{X}_C$ , the cartesian equivalent of  $\tilde{\mathbf{X}}_C$
- The Euclidean transformation can now be expressed as:

$$\begin{bmatrix} \lambda X_C \\ \lambda Y_C \\ \lambda Z_C \\ \lambda \end{bmatrix} = \begin{bmatrix} r_{11} & r_{12} & r_{13} & T_x \\ r_{21} & r_{22} & r_{23} & T_y \\ r_{31} & r_{32} & r_{33} & T_z \\ 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} \lambda X \\ \lambda Y \\ \lambda Z \\ \lambda \end{bmatrix}$$

- Or equivalently:

$$\tilde{\mathbf{X}}_C = \mathbf{P}_e \tilde{\mathbf{X}} \quad \text{where} \quad \mathbf{P}_e = \left[ \begin{array}{ccc|c} & \mathbf{R} & & \mathbf{T} \\ \hline 0 & 0 & 0 & 1 \end{array} \right]$$



# Full camera model in matrix form - Perspective

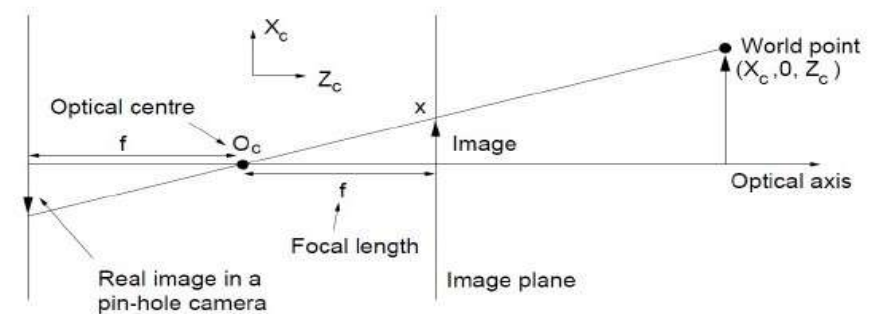
- We want to compute the projection of a 3D point  $\mathbf{X}_c (X_c, Y_c, Z_c)$  expressed in camera coordinates onto a 2D point  $\mathbf{x}(x,y)$  lying on the image plane
- In homogeneous coordinates we have  $\tilde{\mathbf{X}}_c (\lambda X_c, \lambda Y_c, \lambda Z_c, \lambda)$  and  $\tilde{\mathbf{x}}(sx, sy, s)$  – again  $\lambda$  and  $s$  can be set to 1
- Perspective projection can now be expressed as:

$$\begin{bmatrix} sx \\ sy \\ s \end{bmatrix} = \begin{bmatrix} f & 0 & 0 & 0 \\ 0 & f & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} \lambda X_c \\ \lambda Y_c \\ \lambda Z_c \\ \lambda \end{bmatrix}$$

$$\left. \begin{aligned} x &= \frac{f}{Z_c} X_c \\ y &= \frac{f}{Z_c} Y_c \end{aligned} \right\}$$

- Or equivalently:

$$\tilde{\mathbf{x}} = \mathbf{P}_p \tilde{\mathbf{X}}_c$$

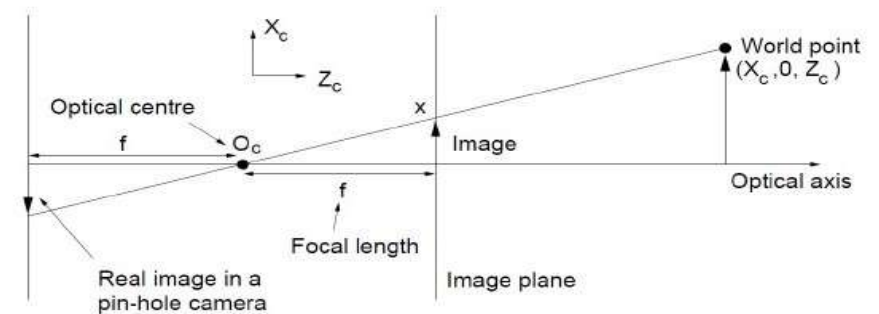


# Full camera model in matrix form - Perspective

- We can confirm that the expression in the previous slide is equivalent to the perspective equation by recovering  $\mathbf{x}$ , the cartesian equivalent of  $\tilde{\mathbf{x}}$

$$\tilde{\mathbf{x}} = \begin{bmatrix} sx \\ sy \\ s \end{bmatrix} = \begin{bmatrix} f\lambda X_c \\ f\lambda Y_c \\ \lambda Z_c \end{bmatrix} \Rightarrow \mathbf{x} = \begin{bmatrix} sx/s \\ sy/s \end{bmatrix} = \begin{bmatrix} f\lambda X_c/\lambda Z_c \\ f\lambda Y_c/\lambda Z_c \end{bmatrix} = \begin{bmatrix} fX_c/Z_c \\ fY_c/Z_c \end{bmatrix}$$

$$\left. \begin{aligned} x &= \frac{f}{Z_c} X_c \\ y &= \frac{f}{Z_c} Y_c \end{aligned} \right\}$$





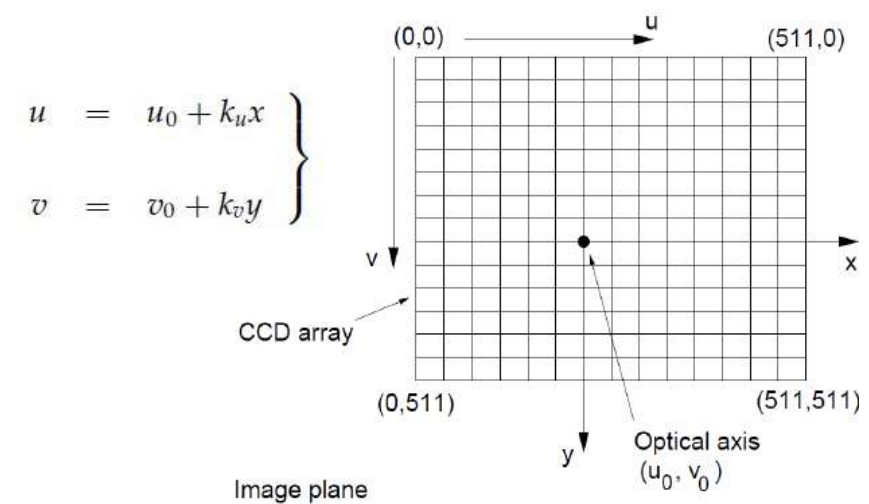
# Full camera model in matrix form – CCD Imaging

- Last, we want to express 2D point  $\mathbf{x}(x,y)$  lying on the image plane in pixel coordinates  $\mathbf{w}(u,v)$
- In homogeneous coordinates we have  $\tilde{\mathbf{x}}(sx, sy, s)$  and  $\tilde{\mathbf{w}}(su, sv, s)$  – again  $s$  can be set to 1
- This is a translation and scaling which can be expressed as:

$$\begin{bmatrix} su \\ sv \\ s \end{bmatrix} = \begin{bmatrix} k_u & 0 & u_0 \\ 0 & k_v & v_0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} sx \\ sy \\ s \end{bmatrix}$$

- Or equivalently:

$$\tilde{\mathbf{w}} = \mathbf{P}_c \tilde{\mathbf{x}}$$



# Full camera model in matrix form – all together

We can now express the whole process, from  $\tilde{X}$  to  $\tilde{w}$  as a single transformation  $P_{ps}$

$$\tilde{w} = P_{ps} \tilde{X}$$

where  $P_{ps} = P_c P_p P_e$

$$= \begin{bmatrix} k_u & 0 & u_0 \\ 0 & k_v & v_0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} f & 0 & 0 & 0 \\ 0 & f & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix} \left[ \begin{array}{ccc|c} & & & \mathbf{T} \\ & \mathbf{R} & & \\ \hline 0 & 0 & 0 & 1 \end{array} \right]$$

# Full camera model in matrix form – all together

- The transformation  $P_{ps}$  is not a general 3x4 matrix, because it has a special structure composed of  $P_e$ ,  $P_p$ , and  $P_c$ .
- We can simplify  $P_{ps}$  into an upper-triangular matrix  $\mathbf{K}$  composed of  $P_e$  and  $P_p$ , and a matrix representing the Euclidean transformation.

$$P_{ps} = \mathbf{K} [\mathbf{R} | \mathbf{T}]$$

$$= \begin{bmatrix} m_u & 0 & u_0 & 0 \\ 0 & m_v & v_0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} r_{11} & r_{12} & r_{13} & T_x \\ r_{21} & r_{22} & r_{23} & T_y \\ r_{31} & r_{32} & r_{33} & T_z \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

$$m_u = k_u f$$

$$m_v = k_v f$$

# Full camera model in matrix form – all together

- The matrix  $\mathbf{K}$  is called the camera calibration matrix and it contains all the viewing parameters (**intrinsic**s) coming from inside the camera.
- The Euclidean transformation matrix contains all the viewing parameters (**extrinsic**s) that are not controlled by the camera lens or sensor.





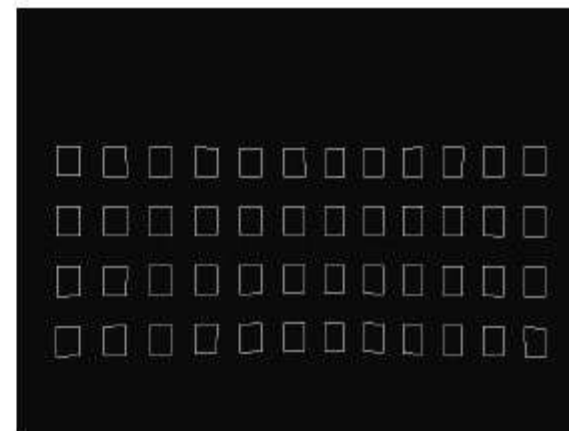
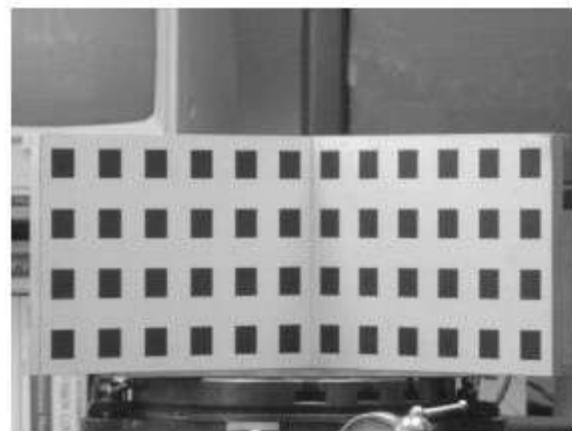
# Today's Agenda

- Full camera model in matrix form
- Camera calibration
- Calibration - Projective camera model
- Calibration – Affine camera model



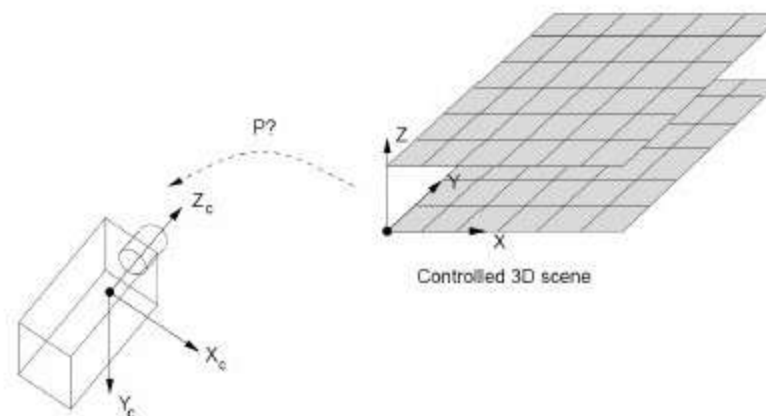
# Camera calibration

- Camera calibration is the name given to the process of discovering the parameters inside our camera model, i.e., the values inside the camera projection matrix.
- This is done by using an image of a controlled scene.
- We may want to use a scene with some sort of regular pattern.



# Camera calibration

- Camera calibration is the name given to the process of discovering the parameters inside our camera model, i.e., the values inside the camera projection matrix. This is done by using an image of a controlled scene.
- We may want to use a scene with some sort of regular pattern





# Today's Agenda

- Full camera model in matrix form
- Camera calibration
- Calibration - Projective camera model
- Calibration – Affine camera model





# Camera calibration – the projective camera

Remember that the **perspective** camera projection matrix  $P_{ps}$  is not a general 3x4 matrix, as it has a special structure composed of  $P_e$ ,  $P_p$ , and  $P_c$ .

$$\begin{aligned} \tilde{\mathbf{w}} &= \mathbf{P}_{ps} \tilde{\mathbf{X}} \\ \text{where } \mathbf{P}_{ps} &= \mathbf{P}_c \mathbf{P}_p \mathbf{P}_e \\ &= \begin{bmatrix} k_u & 0 & u_0 \\ 0 & k_v & v_0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} f & 0 & 0 & 0 \\ 0 & f & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix} \left[ \begin{array}{ccc|c} & & & \mathbf{T} \\ & \mathbf{R} & & \\ \hline 0 & 0 & 0 & 1 \end{array} \right] \end{aligned}$$

# Camera calibration – the projective camera

- Calibrating the perspective camera model may therefore be difficult. We can instead use the projective camera model which is described by a general 3x4 matrix.

$$\mathbf{P} = \begin{bmatrix} p_{11} & p_{12} & p_{13} & p_{14} \\ p_{21} & p_{22} & p_{23} & p_{24} \\ p_{31} & p_{32} & p_{33} & p_{34} \end{bmatrix}$$

- The projective camera has **11 degrees of freedom**, since the overall scale of  $\mathbf{P}$  does not matter when using homogeneous coordinates.
- It is more convenient to deal with a projective camera than a perspective one, since we don't have to worry about placing nonlinear constraints on the elements of  $\mathbf{P}$

# Camera calibration – the projective camera

- Using a projective camera, the whole imaging process is described by

$$\tilde{\mathbf{w}} = \mathbf{P}\tilde{\mathbf{X}}$$
$$\Rightarrow \begin{bmatrix} su \\ sv \\ s \end{bmatrix} = \begin{bmatrix} p_{11} & p_{12} & p_{13} & p_{14} \\ p_{21} & p_{22} & p_{23} & p_{24} \\ p_{31} & p_{32} & p_{33} & p_{34} \end{bmatrix} \begin{bmatrix} X \\ Y \\ Z \\ 1 \end{bmatrix}$$

- We must estimate 11 parameters (the overall scale does not matter, so let's set  $p_{34} = 1$ )

# Camera calibration – the projective camera

- Because the scene is controlled, we **know** the location of the world points  $X_i$
- We can also find the location of their projections  $w_i$  on the image, e.g., using corner detection.
- Each point we observe gives us two equations

$$u_i = \frac{su_i}{s} = \frac{p_{11}X_i + p_{12}Y_i + p_{13}Z_i + p_{14}}{p_{31}X_i + p_{32}Y_i + p_{33}Z_i + 1}$$

$$v_i = \frac{sv_i}{s} = \frac{p_{21}X_i + p_{22}Y_i + p_{23}Z_i + p_{24}}{p_{31}X_i + p_{32}Y_i + p_{33}Z_i + 1}$$



# Camera calibration – the projective camera

Rearranging them gives us two linear equations in the unknown parameters of the projection matrix

$$\begin{bmatrix} X_i & Y_i & Z_i & 1 & 0 & 0 & 0 & 0 & u_i X_i & u_i Y_i & u_i Z_i \\ 0 & 0 & 0 & 0 & X_i & Y_i & Z_i & 1 & v_i X_i & v_i Y_i & v_i Z_i \end{bmatrix} \begin{bmatrix} p_{11} \\ p_{12} \\ p_{13} \\ p_{14} \\ p_{21} \\ p_{22} \\ p_{23} \\ p_{24} \\ p_{31} \\ p_{32} \\ p_{33} \end{bmatrix} = \begin{bmatrix} -u_i \\ -v_i \end{bmatrix}$$

# Camera calibration – the projective camera

- As there are 11 unknowns, we need at least 6 points to calibrate this camera model to get enough equations (rows) in the linear system  $Ax=b$ .
- Each observed points adds two equations (rows) to the matrix  $A$
- We can solve the system of equations using linear least squares (pseudo-inverse of  $A$ )

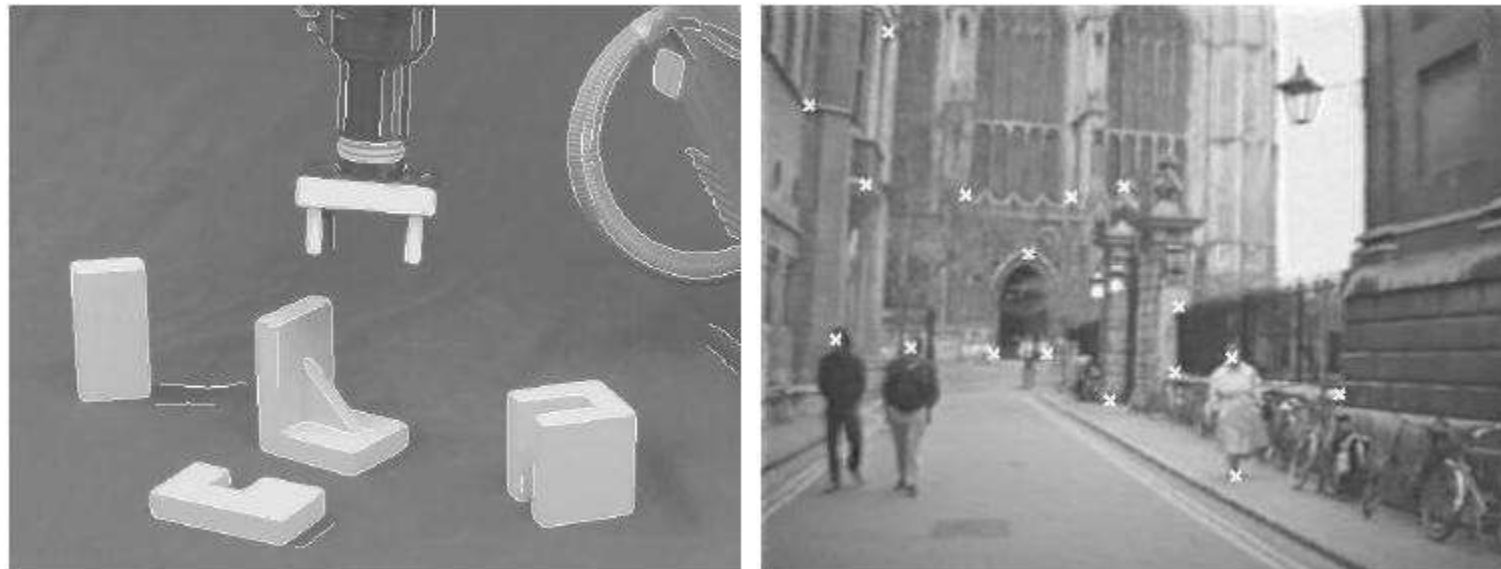
$$\begin{aligned} Ax &= b \\ \Rightarrow x &= (A^T A)^{-1} A^T b \end{aligned}$$

# Today's Agenda

- Full camera model in matrix form
- Camera calibration
- Calibration – Projective camera model
- Calibration – Affine camera model

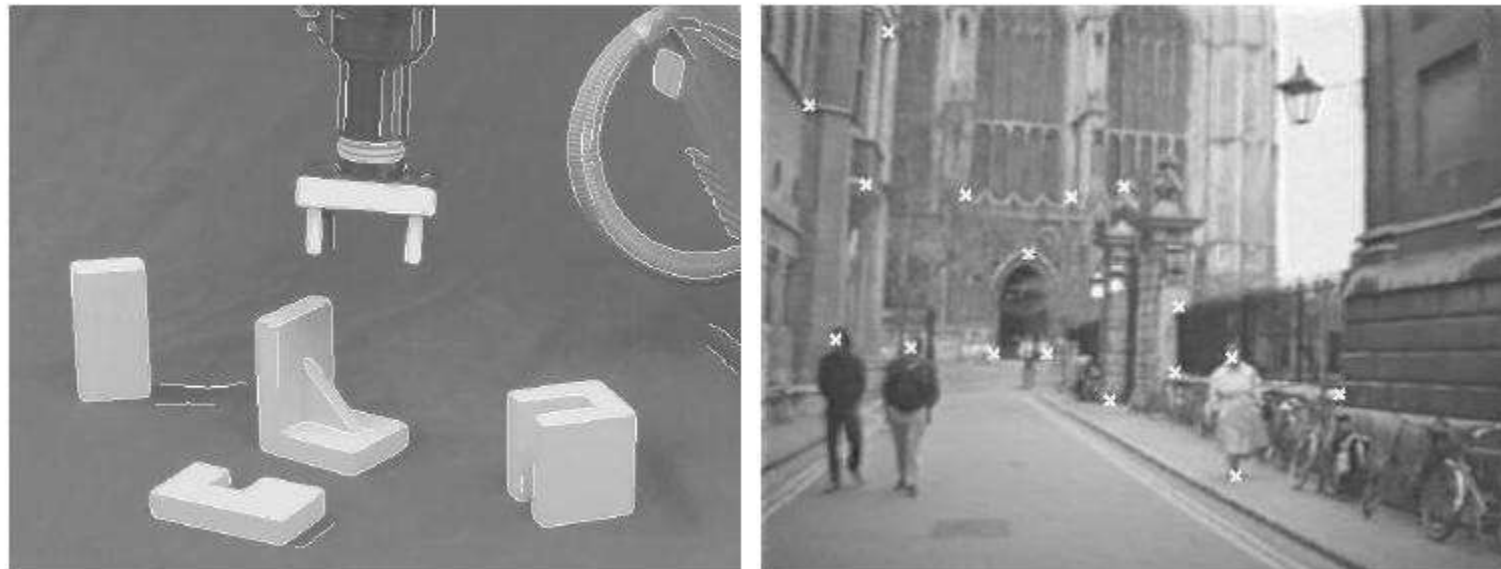
# Parallel Projection

When the depth difference  $\Delta Z_C$  of objects in the scene is small compared to their distance  $Z_C$  to the camera, the resulting image is not described well by perspective projection.



# Parallel Projection

For example, parallel lines on the left image remain parallel after projection. In this case, the 3D to 2D projection is better described by **parallel projection**.

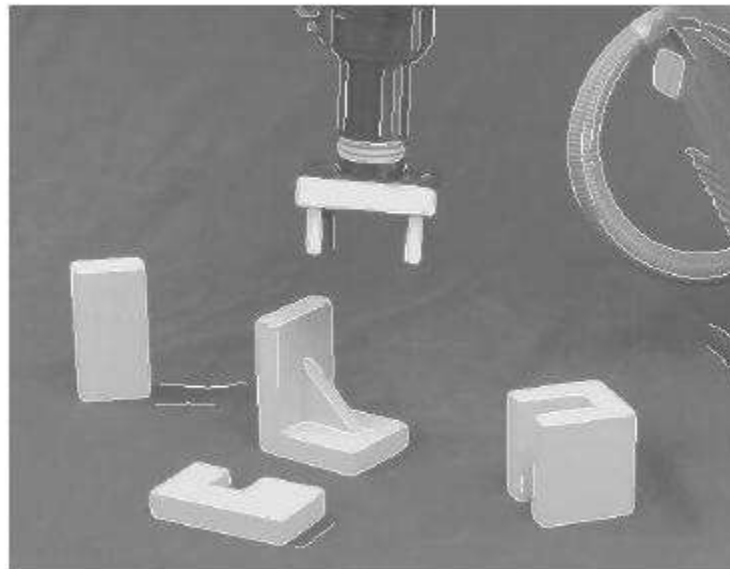






# Parallel Projection

A camera which creates images like the one on the left is known as a **weak-perspective camera**.



# Parallel Projection

What changes in our imaging transformations is the perspective projection matrix. Remember  $P_p$ :

$$\begin{bmatrix} sx \\ sy \\ s \end{bmatrix} = \begin{bmatrix} f & 0 & 0 & 0 \\ 0 & f & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} X_c \\ Y_c \\ Z_c \\ 1 \end{bmatrix}$$

# Parallel Projection

What changes in our imaging transformations is the perspective projection matrix.

If the depth variation in the scene is small, then  $Z_c \approx Z_c^{avg}$  and the perspective projection matrix  $P_p$  can be changed to the parallel projection matrix  $P_{pll}$

$$\begin{bmatrix} sx \\ sy \\ s \end{bmatrix} = \begin{bmatrix} f & 0 & 0 & 0 \\ 0 & f & 0 & 0 \\ 0 & 0 & 0 & Z_c^{avg} \end{bmatrix} \begin{bmatrix} X_c \\ Y_c \\ Z_c \\ 1 \end{bmatrix}$$

$$\tilde{x} = P_{pll} \tilde{X}_c$$

# Weak-Perspective Camera Model

This leads to the **weak perspective** camera model:

$$\tilde{\mathbf{w}} = \mathbf{P}_{wp} \tilde{\mathbf{X}}$$

where

$$\begin{aligned} \mathbf{P}_{wp} &= \mathbf{P}_c \mathbf{P}_{pll} \mathbf{P}_e \\ &= \begin{bmatrix} k_u & 0 & u_0 \\ 0 & k_v & v_0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} f & 0 & 0 & 0 \\ 0 & f & 0 & 0 \\ 0 & 0 & 0 & Z_c^{\text{avg}} \end{bmatrix} \left[ \begin{array}{ccc|c} & & & \\ & \mathbf{R} & & \mathbf{T} \\ \hline 0 & 0 & 0 & 1 \end{array} \right] \end{aligned}$$

# Weak-Perspective Camera Model

$P_{wp}$  is the projection matrix for a weak-perspective camera. It is a 3x4 matrix with a special structure composed of  $P_e$ ,  $P_{pll}$ , and  $P_c$ .

$$\tilde{\mathbf{w}} = \mathbf{P}_{wp} \tilde{\mathbf{X}}$$

where

$$\begin{aligned} \mathbf{P}_{wp} &= \mathbf{P}_c \mathbf{P}_{pll} \mathbf{P}_e \\ &= \begin{bmatrix} k_u & 0 & u_0 \\ 0 & k_v & v_0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} f & 0 & 0 & 0 \\ 0 & f & 0 & 0 \\ 0 & 0 & 0 & Z_c^{\text{avg}} \end{bmatrix} \begin{bmatrix} \mathbf{R} & | & \mathbf{T} \\ \hline 0 & 0 & 0 & 1 \end{bmatrix} \end{aligned}$$



# Camera calibration – the affine camera

- The special structure of the weak-perspective camera model makes it difficult to calibrate. This is why the affine camera model is often used instead:

$$\mathbf{P}_{aff} = \begin{bmatrix} p_{11} & p_{12} & p_{13} & p_{14} \\ p_{21} & p_{22} & p_{23} & p_{24} \\ 0 & 0 & 0 & p_{34} \end{bmatrix}$$

- $P_{aff}$  is the projection matrix for the affine camera. It has 8 degrees of freedom (remember the overall scale does not matter so we can set  $p_{34}$  to 1).
- It can be calibrated in the same way as the projective camera. There are eight degrees of freedom, so we need a minimum of 4 points.

**MAI4CAREU**

Master programmes in Artificial  
Intelligence 4 Careers in Europe



**CYENS**  
CENTRE OF EXCELLENCE



# Thank you.





University of Cyprus – MSc Artificial Intelligence

# MAI644 – COMPUTER VISION

## Lecture 16: Stereo Vision

**Melinos Averkiou**

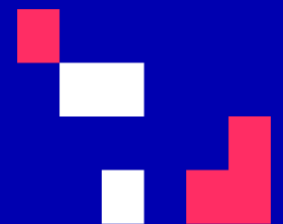
CYENS Centre of Excellence

University of Cyprus - Department of Computer Science

[m.averkiou@cyens.org.cy](mailto:m.averkiou@cyens.org.cy)



**CYENS**  
CENTRE OF EXCELLENCE



# Last time

- Full camera model in matrix form
- Camera calibration
- Calibration – Projective camera model
- Calibration – Affine camera model

# Today's Agenda

- Recovery of world position
- Triangulation
- Epipolar Geometry

[material based on Paris Kaimakis]





# Today's Agenda

- Recovery of world position
- Triangulation
- Epipolar Geometry



# Recovery of world position

- Previously we saw that the imaging process can be described as a transformation in homogeneous coordinates.
- If we can invert this transformation, the world coordinates of each pixel in the image can be computed.
- Recovering world coordinates of objects based on the projection on an image is known as **shape recovery** or **depth recovery**.

# Recovery of world position

- Is this possible using a single camera ?
- The camera needs to be calibrated, i.e. we know all its parameters
- Remember the projective camera model:

$$\tilde{\mathbf{w}} = \mathbf{P}\tilde{\mathbf{X}}$$
$$\Leftrightarrow \begin{bmatrix} su \\ sv \\ s \end{bmatrix} = \begin{bmatrix} p_{11} & p_{12} & p_{13} & p_{14} \\ p_{21} & p_{22} & p_{23} & p_{24} \\ p_{31} & p_{32} & p_{33} & p_{34} \end{bmatrix} \begin{bmatrix} X \\ Y \\ Z \\ 1 \end{bmatrix}$$

- Unfortunately the transformation described in  $\mathbf{P}$  is not invertible and the world point  $\mathbf{X}$  cannot be uniquely determined

# Recovery of world position

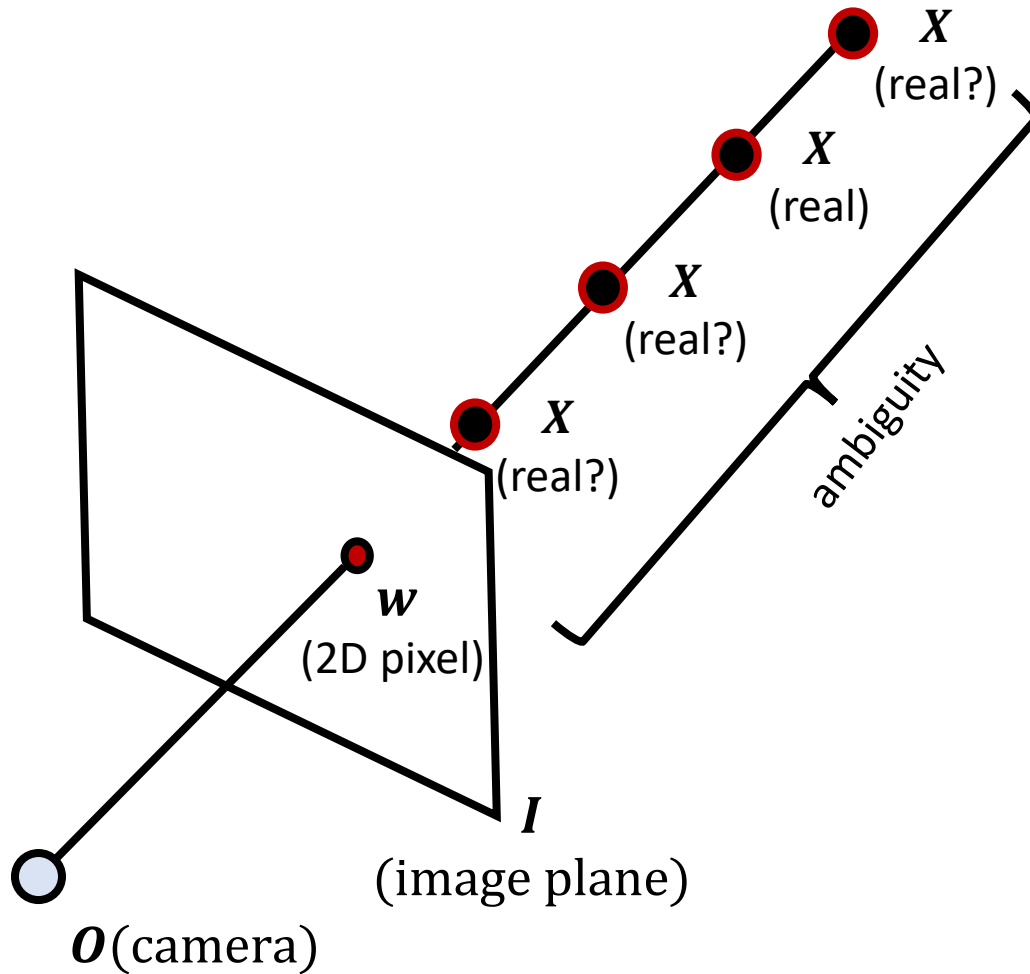
- Depth ambiguity



Courtesy slide S. Lazebnik

# Recovery of world position

- Each observed feature on the image gives 2 equations with 3 unknowns and therefore defines a line (a ray) of solutions for  $X$





# Recovery of world position

- Each observed feature on the image gives 2 equations with 3 unknowns and therefore defines a line (a ray) of solutions for  $\mathbf{X}$
- This system of equations is under-constrained.
- This can be seen by the size of  $\mathbf{P}$ . There are more columns than rows.

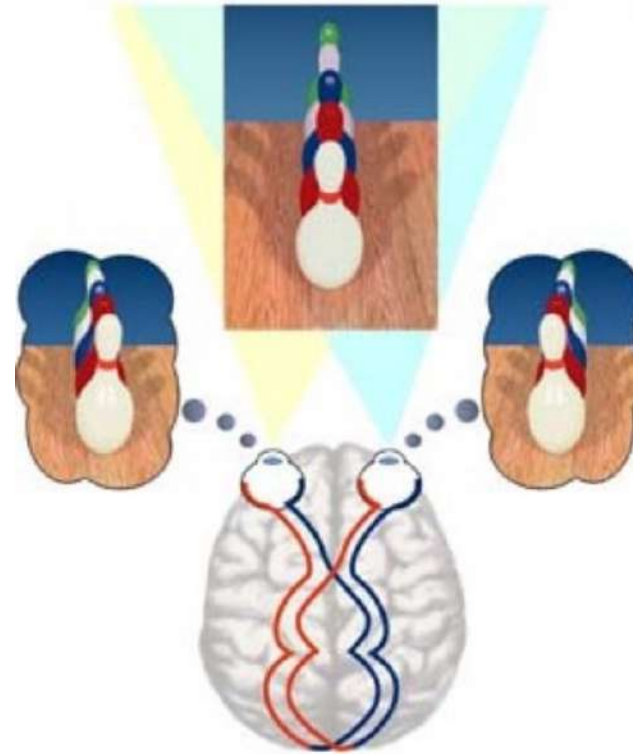
$$\tilde{\mathbf{w}} = \mathbf{P}\tilde{\mathbf{X}}$$
$$\Leftrightarrow \begin{bmatrix} su \\ sv \\ s \end{bmatrix} = \begin{bmatrix} p_{11} & p_{12} & p_{13} & p_{14} \\ p_{21} & p_{22} & p_{23} & p_{24} \\ p_{31} & p_{32} & p_{33} & p_{34} \end{bmatrix} \begin{bmatrix} X \\ Y \\ Z \\ 1 \end{bmatrix}$$

# Recovery of world position

- Under-constrained problems never have a unique solution.
- To uniquely recover  $\mathbf{X}$ , additional views must be used, so that the transformation between  $\mathbf{w}$  (pixel coordinates) and  $\mathbf{X}$  (world coordinates) is *forced* to become invertible.
- This is the subject of **stereo vision**.

# Recovery of world position

- Two eyes/cameras help.

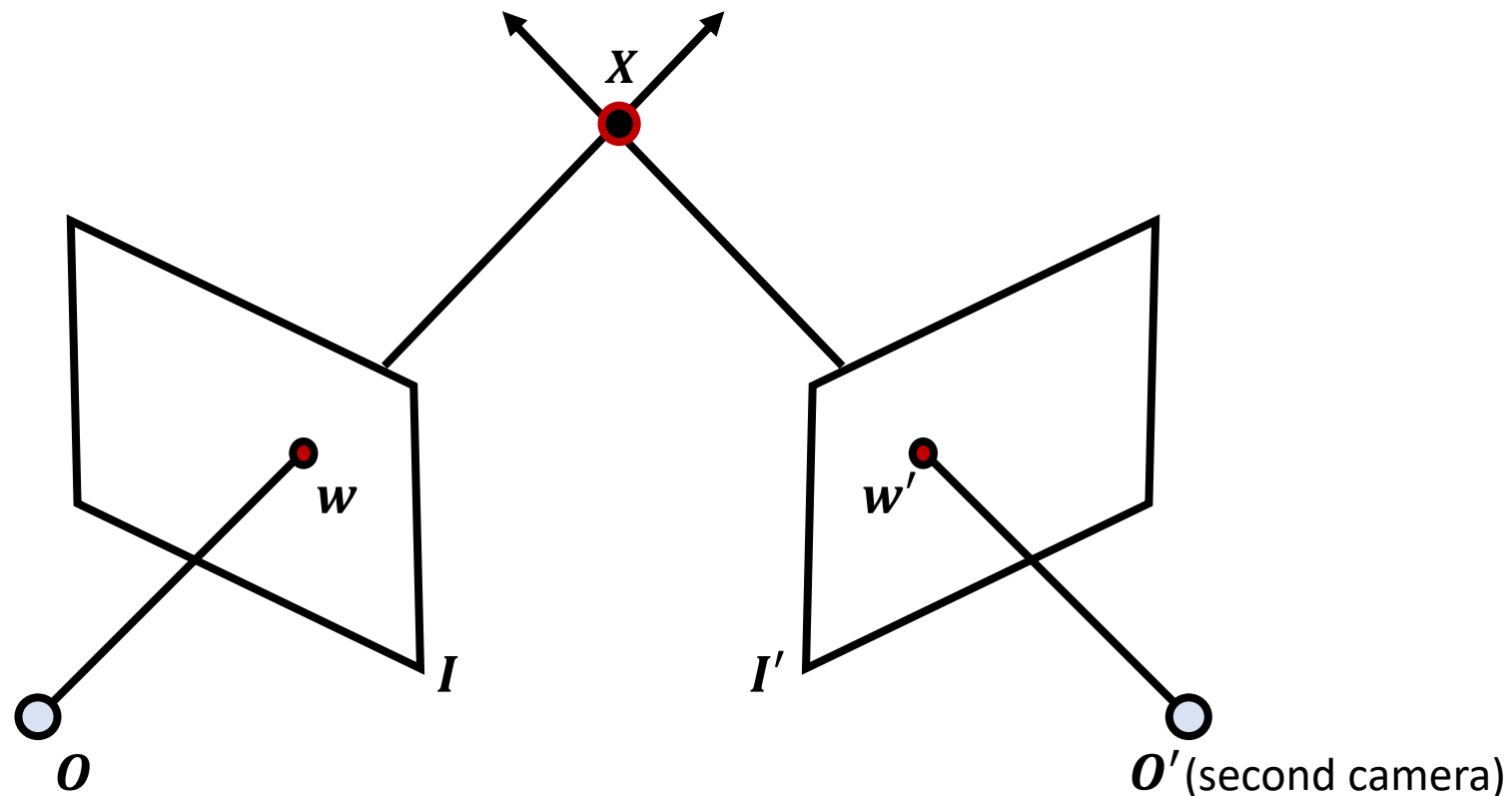


# Today's Agenda

- Recovery of world position
- Triangulation
- Epipolar Geometry

# Triangulation

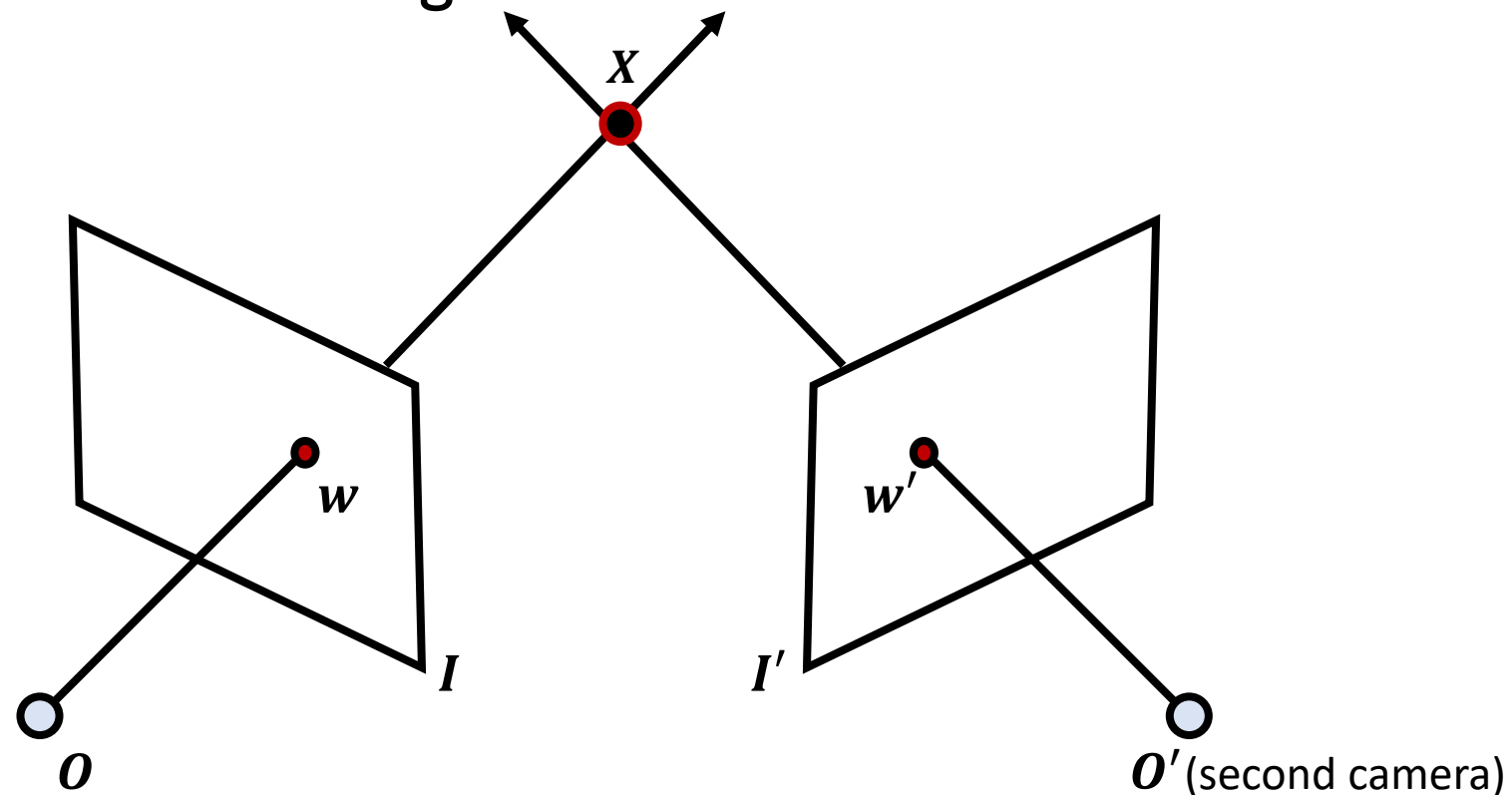
- In stereo vision at least two cameras are set up to view the  $3D$  scene.
- Each  $3D$  world location  $X$  projects to pixel  $w$  on camera 1 ( $O$ ) and to pixel  $w'$  on camera 2 ( $O'$ ).





# Triangulation

- If both cameras are calibrated, the 3D world location  $X$  projected on the pair of corresponding pixel locations  $w$  and  $w'$  can be estimated via a process known as **triangulation**.



# Triangulation

- Consider the projection of  $\mathbf{X}$  onto  $\mathbf{w}$ :

$$\tilde{\mathbf{w}} = \mathbf{P}\tilde{\mathbf{X}}$$

$$\Leftrightarrow \begin{bmatrix} su \\ sv \\ s \end{bmatrix} = \begin{bmatrix} p_{11} & p_{12} & p_{13} & p_{14} \\ p_{21} & p_{22} & p_{23} & p_{24} \\ p_{31} & p_{32} & p_{33} & p_{34} \end{bmatrix} \begin{bmatrix} X \\ Y \\ Z \\ 1 \end{bmatrix}$$

- Compute the equation for  $u$  and rearrange so our unknowns  $X, Y, Z$  are on the left:

$$u = \frac{su}{s} = \frac{p_{11}X + p_{12}Y + p_{13}Z + p_{14}}{p_{31}X + p_{32}Y + p_{33}Z + p_{34}}$$

$$\Rightarrow p_{11}X + p_{12}Y + p_{13}Z + p_{14} = p_{31}uX + p_{32}uY + p_{33}uZ + p_{34}u$$

$$\Rightarrow (p_{11} - p_{31}u)X + (p_{12} - p_{32}u)Y + (p_{13} - p_{33}u)Z = p_{34}u - p_{14}$$



# Triangulation

- Compute the equation for  $u$  and rearrange so our unknowns  $X, Y, Z$  are on the left:

$$u = \frac{su}{s} = \frac{p_{11}X + p_{12}Y + p_{13}Z + p_{14}}{p_{31}X + p_{32}Y + p_{33}Z + p_{34}}$$

$$\Rightarrow p_{11}X + p_{12}Y + p_{13}Z + p_{14} = p_{31}uX + p_{32}uY + p_{33}uZ + p_{34}u$$

$$\Rightarrow (p_{11} - p_{31}u)X + (p_{12} - p_{32}u)Y + (p_{13} - p_{33}u)Z = p_{34}u - p_{14}$$

- Put this in matrix form:

$$\begin{bmatrix} p_{11} - p_{31}u & p_{12} - p_{32}u & p_{13} - p_{33}u \end{bmatrix} \begin{bmatrix} X \\ Y \\ Z \end{bmatrix} = [p_{34}u - p_{14}]$$

$$\Leftrightarrow \mathbf{Ax} = \mathbf{b}$$



# Triangulation

- Put this in matrix form:

$$\begin{bmatrix} p_{11} - p_{31}u & p_{12} - p_{32}u & p_{13} - p_{33}u \end{bmatrix} \begin{bmatrix} X \\ Y \\ Z \end{bmatrix} = [p_{34}u - p_{14}]$$
$$\Leftrightarrow \mathbf{Ax} = \mathbf{b}$$

- $\mathbf{A}$  is a  $1 \times 3$  matrix
- The resulting system is under-constrained

# Triangulation

- Put this in matrix form:

$$\begin{aligned}
 & [p_{11} - p_{31}u \quad p_{12} - p_{32}u \quad p_{13} - p_{33}u] \begin{bmatrix} X \\ Y \\ Z \end{bmatrix} = [p_{34}u - p_{14}] \\
 & [p_{21} - p_{31}v \quad p_{22} - p_{32}v \quad p_{23} - p_{33}v] \begin{bmatrix} X \\ Y \\ Z \end{bmatrix} = [p_{34}v - p_{24}] \\
 & \Leftrightarrow \mathbf{Ax} = \mathbf{b}
 \end{aligned}$$

- By computing  $v$  in the same way as  $u$  and rearranging we can add a new row in  $\mathbf{A}$  and in  $\mathbf{b}$ , making it  $2 \times 3$



# Triangulation

- Put this in matrix form:

$$\begin{aligned}
 & [p'_{11} - p'_{31}u' \quad p'_{12} - p'_{32}u' \quad p'_{13} - p'_{33}u'] \begin{bmatrix} X \\ Y \\ Z \end{bmatrix} = [p'_{34}u' - p'_{14}] \\
 & [p'_{21} - p'_{31}v' \quad p'_{22} - p'_{32}v' \quad p'_{23} - p'_{33}v'] \begin{bmatrix} X \\ Y \\ Z \end{bmatrix} = [p'_{34}v' - p'_{24}] \\
 & \Leftrightarrow \mathbf{Ax} = \mathbf{b}
 \end{aligned}$$

- By also considering the projection of  $\mathbf{X}$  onto  $\mathbf{w}'$ , a new pair of rows will be added to  $\mathbf{A}$ , thus forcing it to be  $4 \times 3$ , i.e. over-constrained

# Triangulation

- Here is the resulting system

$$\begin{bmatrix} p_{11} - p_{31}u & p_{12} - p_{32}u & p_{13} - p_{33}u \\ p_{21} - p_{31}v & p_{22} - p_{32}v & p_{23} - p_{33}v \\ p'_{11} - p'_{31}u' & p'_{12} - p'_{32}u' & p'_{13} - p'_{33}u' \\ p'_{21} - p'_{31}v' & p'_{22} - p'_{32}v' & p'_{23} - p'_{33}v' \end{bmatrix} \begin{bmatrix} X \\ Y \\ Z \end{bmatrix} = \begin{bmatrix} p_{34}u - p_{14} \\ p_{34}v - p_{24} \\ p'_{34}u' - p'_{14} \\ p'_{34}v' - p'_{24} \end{bmatrix}$$

- Where  $p'_{11}$  etc. are the parameters inside the camera projection matrix for camera 2 ( $\mathbf{O}'$ ), and  $\mathbf{w}' = (u', v')$  are the pixel coordinates of  $\mathbf{X}$  projected on the image plane  $\mathbf{I}'$  of the second camera



# Triangulation

- Using an additional camera has forced  $\mathbf{A}$  to become over-constrained.

$$\begin{bmatrix} p_{11} - p_{31}u & p_{12} - p_{32}u & p_{13} - p_{33}u \\ p_{21} - p_{31}v & p_{22} - p_{32}v & p_{23} - p_{33}v \\ p'_{11} - p'_{31}u' & p'_{12} - p'_{32}u' & p'_{13} - p'_{33}u' \\ p'_{21} - p'_{31}v' & p'_{22} - p'_{32}v' & p'_{23} - p'_{33}v' \end{bmatrix} \begin{bmatrix} X \\ Y \\ Z \end{bmatrix} = \begin{bmatrix} p_{34}u - p_{14} \\ p_{34}v - p_{24} \\ p'_{34}u' - p'_{14} \\ p'_{34}v' - p'_{24} \end{bmatrix}$$

- Therefore we can now find a least-squares solution:

$$\begin{aligned} \mathbf{Ax} &= \mathbf{b} \\ \Leftrightarrow \mathbf{x} &= (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{b} \end{aligned}$$

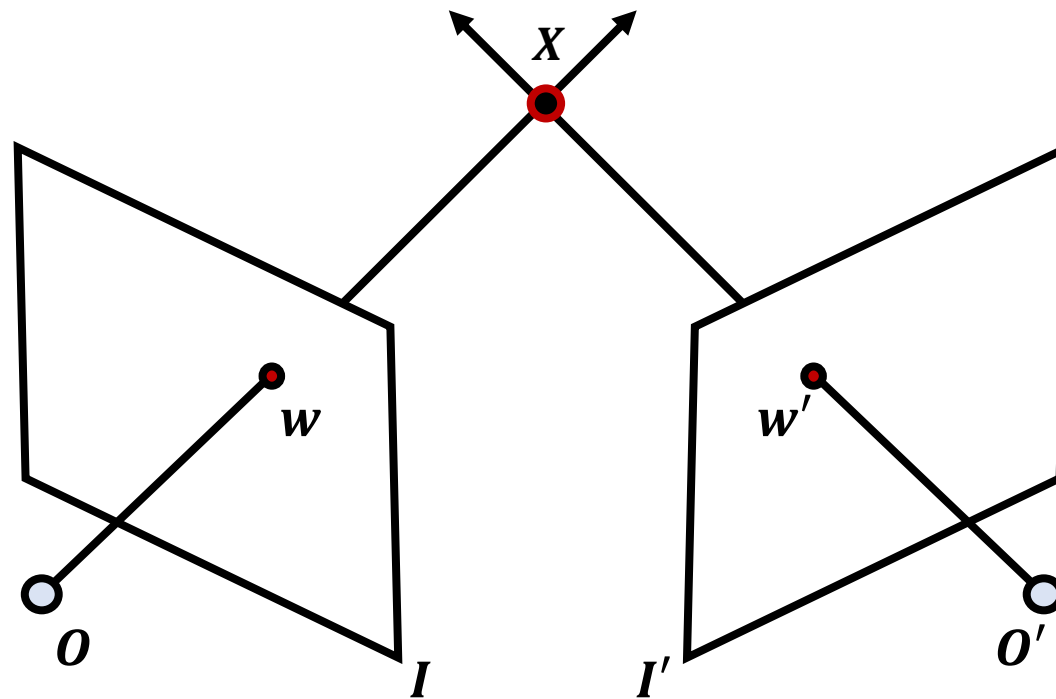


# Today's Agenda

- Recovery of world position
- Triangulation
- Epipolar Geometry

# Beyond triangulation

We have seen the simplest form of stereo vision: Given a pair of *calibrated* cameras observing a single feature at *corresponding* pixel locations  $\mathbf{w} \leftrightarrow \mathbf{w}'$ , the 3D position of the corresponding world location  $\mathbf{X}$  can be estimated via triangulation





# Beyond triangulation

How is the correspondence problem solved if there are several points  $\{\mathbf{w}_i\}_{i=1}^{N_1}$  in image 1 and several points  $\{\mathbf{w}'_j\}_{j=1}^{N_2}$  in image 2 ?

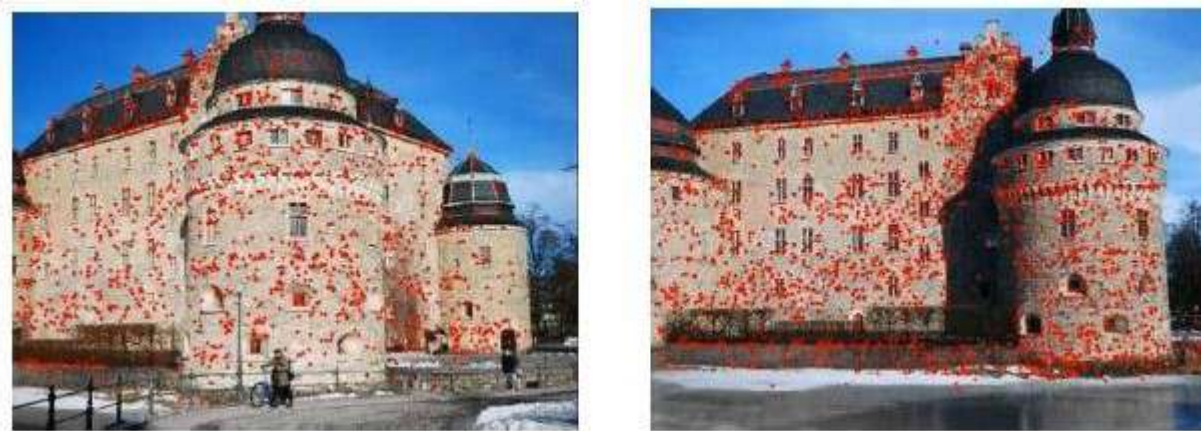


# Beyond triangulation

SIFT will give us a set of proposed correspondences  $\{w_i \leftrightarrow w'_j\}$  but there will be **many** outliers in these proposals

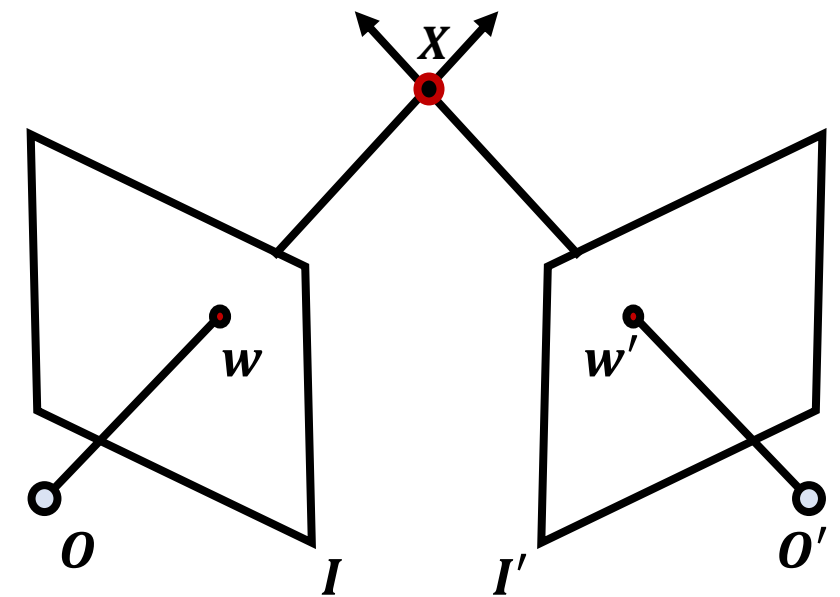
The question is then how can we remove outliers ?

Using the **epipolar constraint**



# Epipolar Geometry

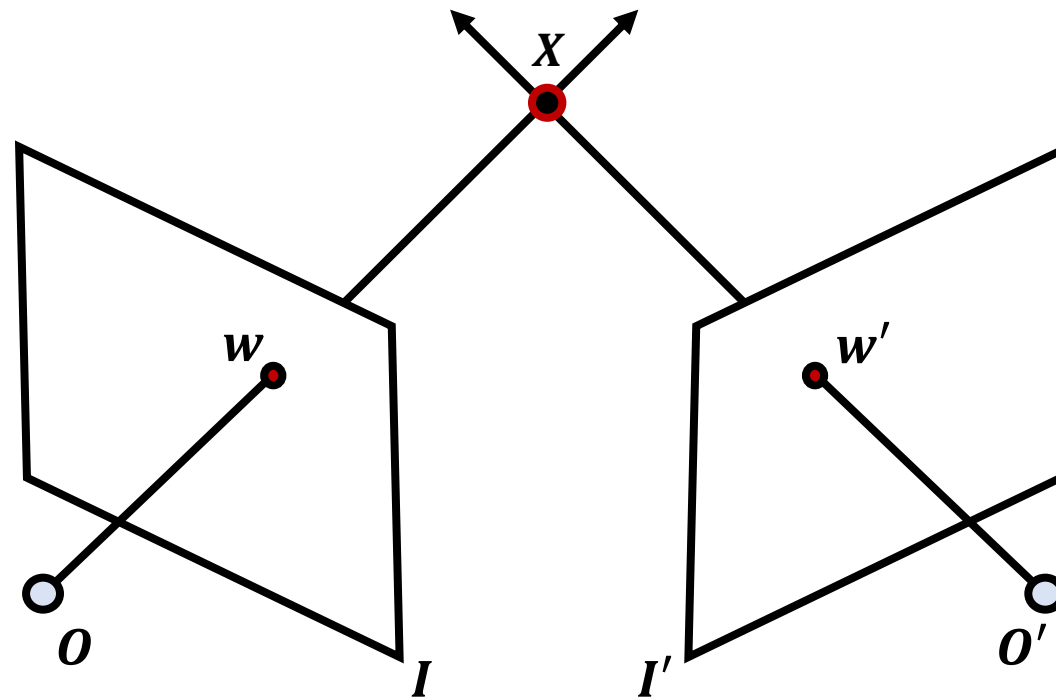
- To understand the **epipolar constraint**, we first need to understand the **geometry** that relates the
  - cameras
  - points in **3D** space
  - and their corresponding observations  $\{w_i \leftrightarrow w'_j\}$



- This type of geometry is referred to as the **epipolar geometry**

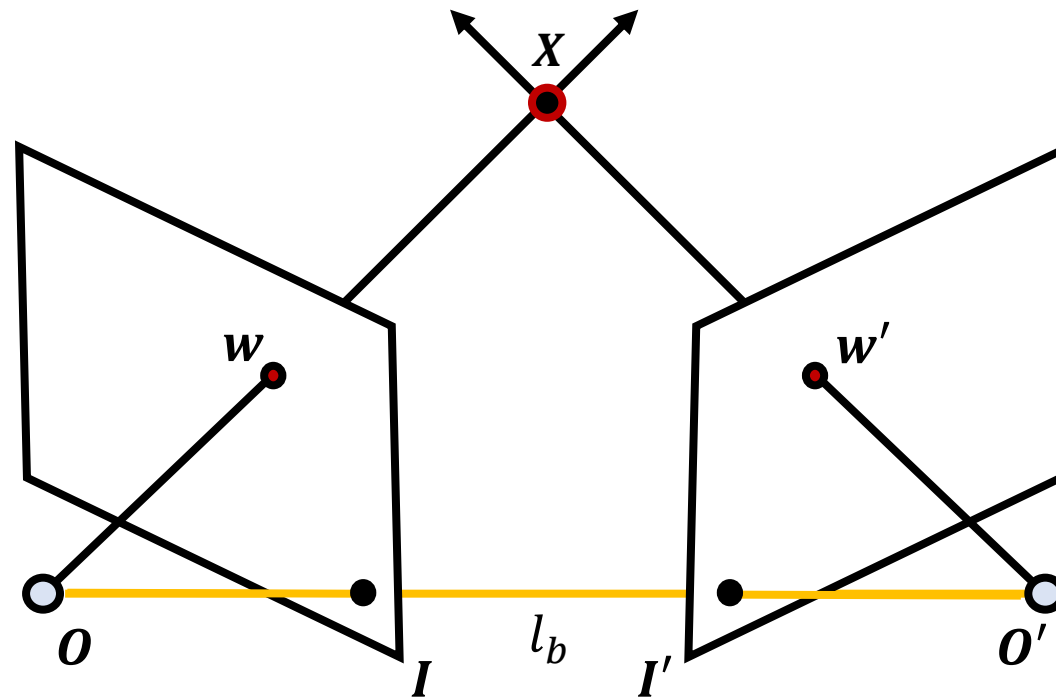
# Epipolar Geometry - fundamentals

Lets revisit our stereo pair



# Epipolar Geometry - fundamentals

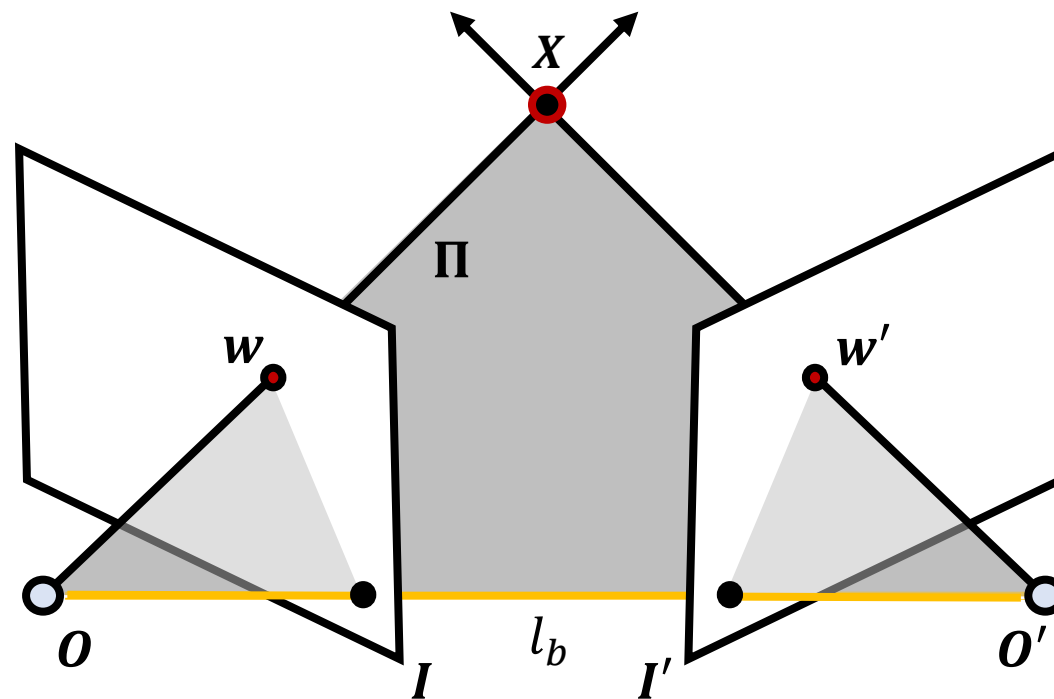
The **baseline**  $l_b$  is the line joining the two optical centers.





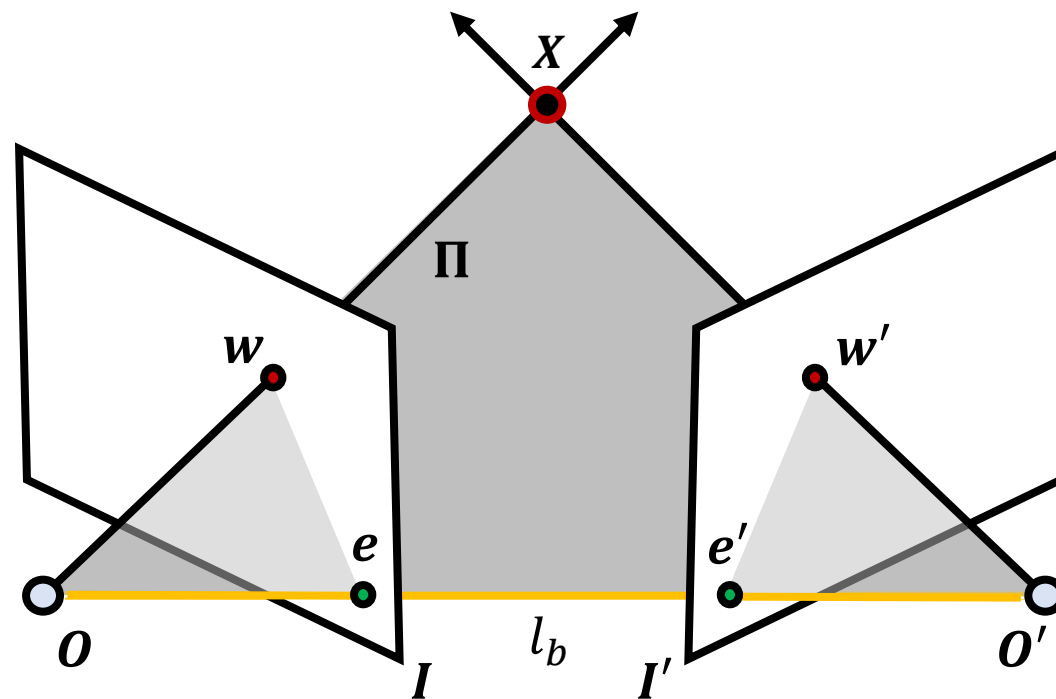
# Epipolar Geometry - fundamentals

The epipolar plane  $\Pi$  is the plane defined by the 3D point  $X$  and the optical centers of the cameras.



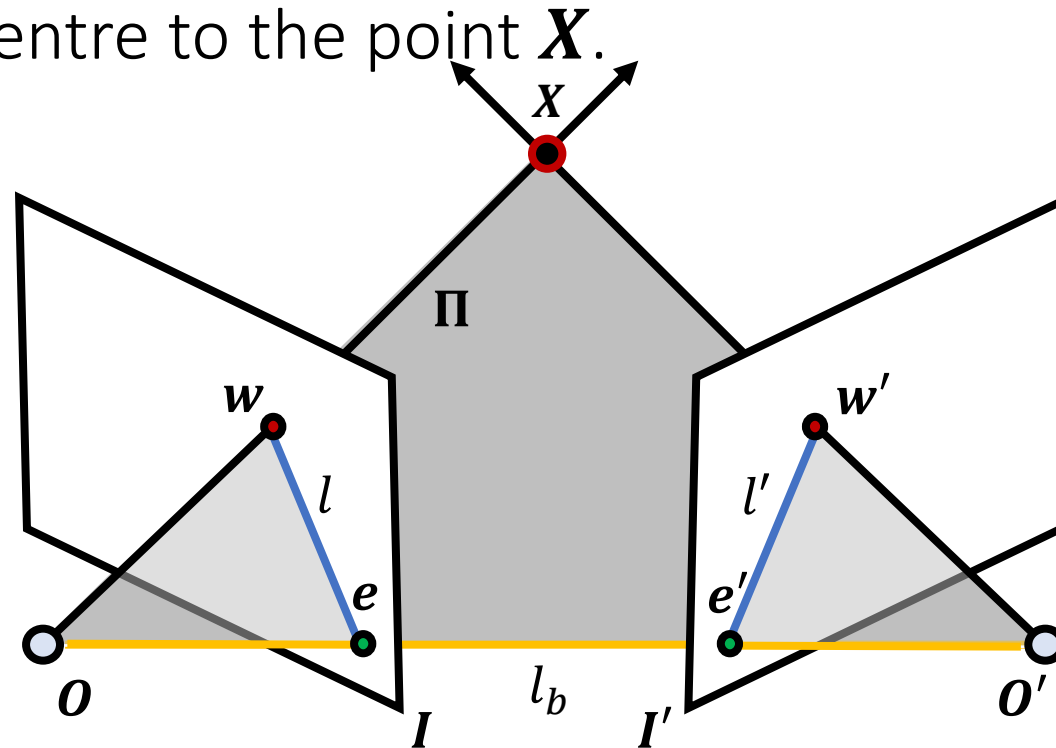
# Epipolar Geometry - fundamentals

An **epipole** is the point of intersection of the **baseline** with the image plane. There are two epipoles  $e$  and  $e'$ , one for each image.



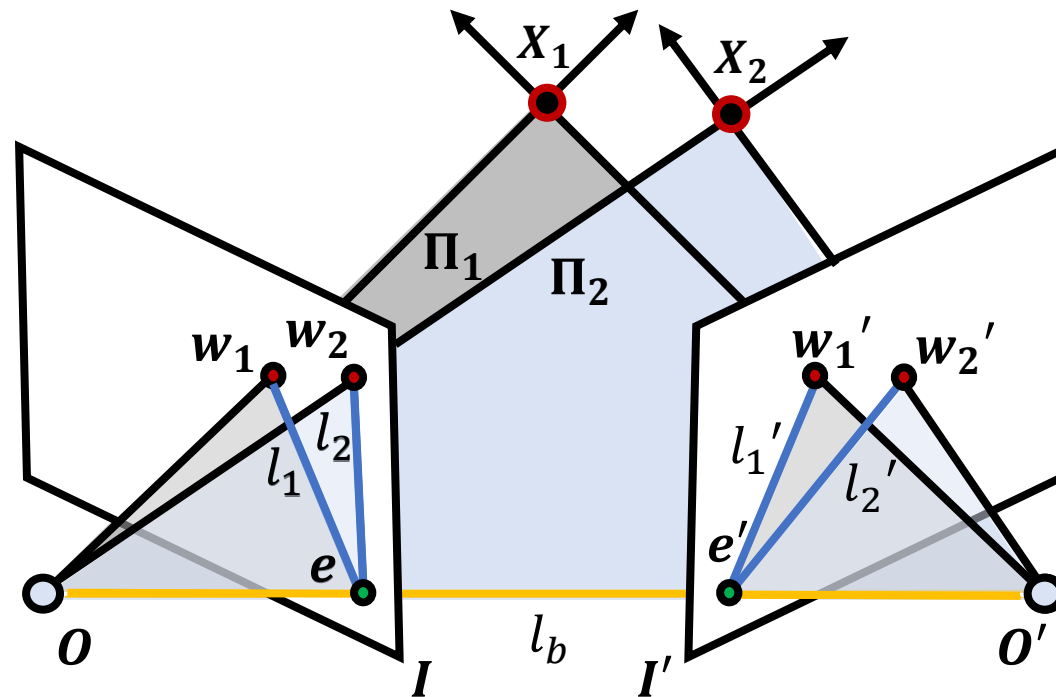
# Epipolar Geometry - fundamentals

An **epipolar line** is a line of intersection of the epipolar plane with an image plane. It is the image, in one camera, of the ray from the other camera's optical centre to the point  $X$ .



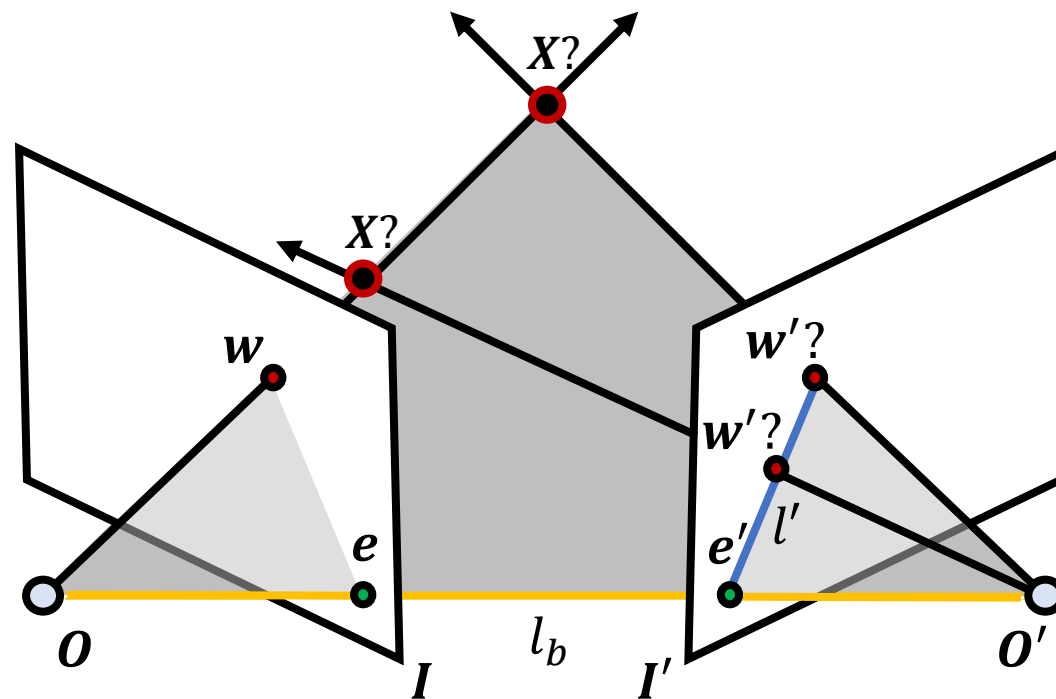
# Epipolar Geometry - fundamentals

For different world points  $X$ , the epipolar plane rotates about the baseline. All epipolar lines intersect at their corresponding epipole.



# Epipolar Geometry - fundamentals

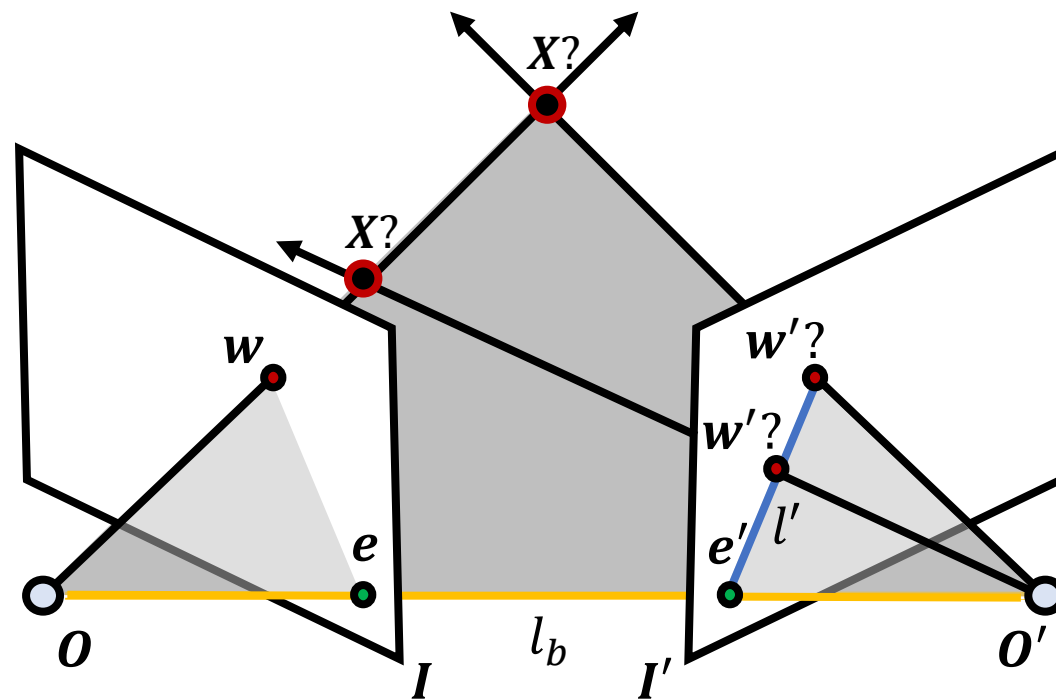
The **epipolar constraint** limits the search for correspondences, from the region of the whole image, to only the pixels spanned by the epipolar line.





# Epipolar Geometry - fundamentals

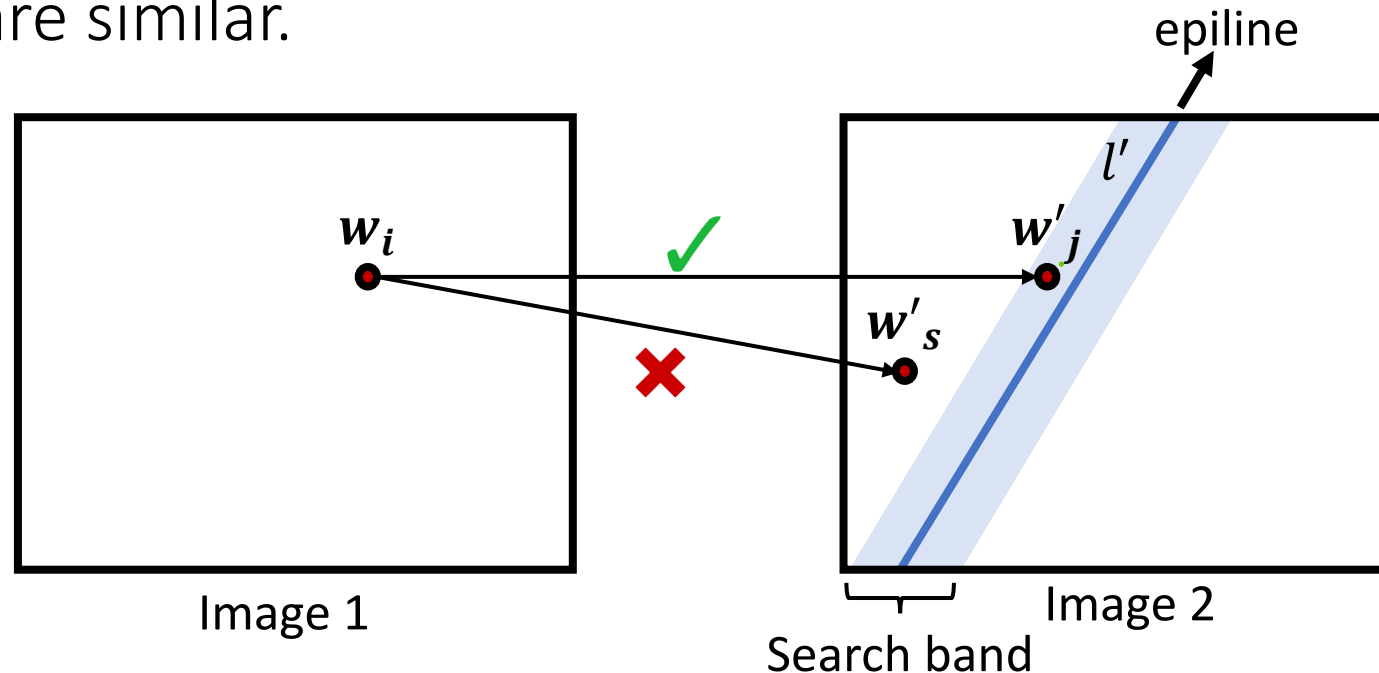
If a point feature  $w$  is observed in one image, then its location  $w'$  in the other image must lie on its corresponding epipolar line  $l'$





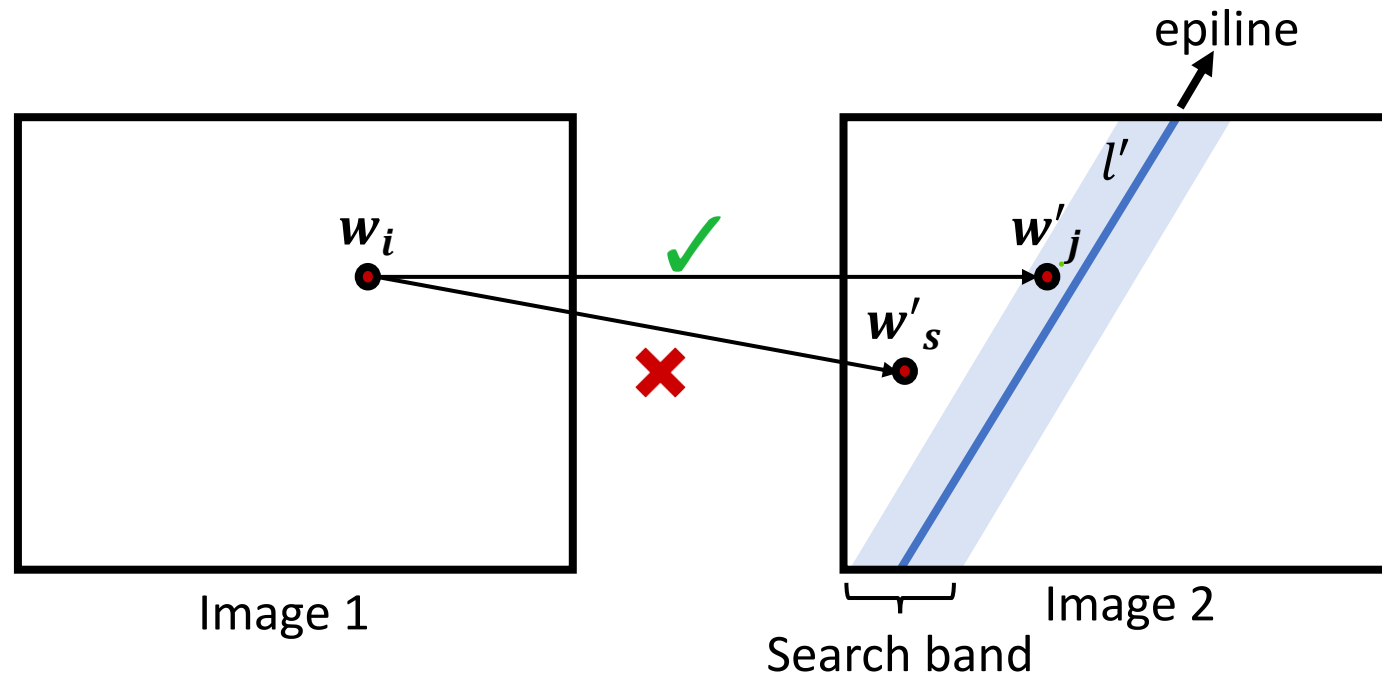
# Epipolar Geometry - fundamentals

So a simple algorithm for determining correspondences is to match each feature  $w_i$  from camera 1 to a feature  $w'_j$  from camera 2 which is close to the epipolar line of camera 2, provided that the SIFT descriptors for  $w_i$  and  $w'_j$  are similar.



# Epipolar Geometry - fundamentals

Therefore, we must calculate the equation of the epipolar lines.





# Essential Matrix

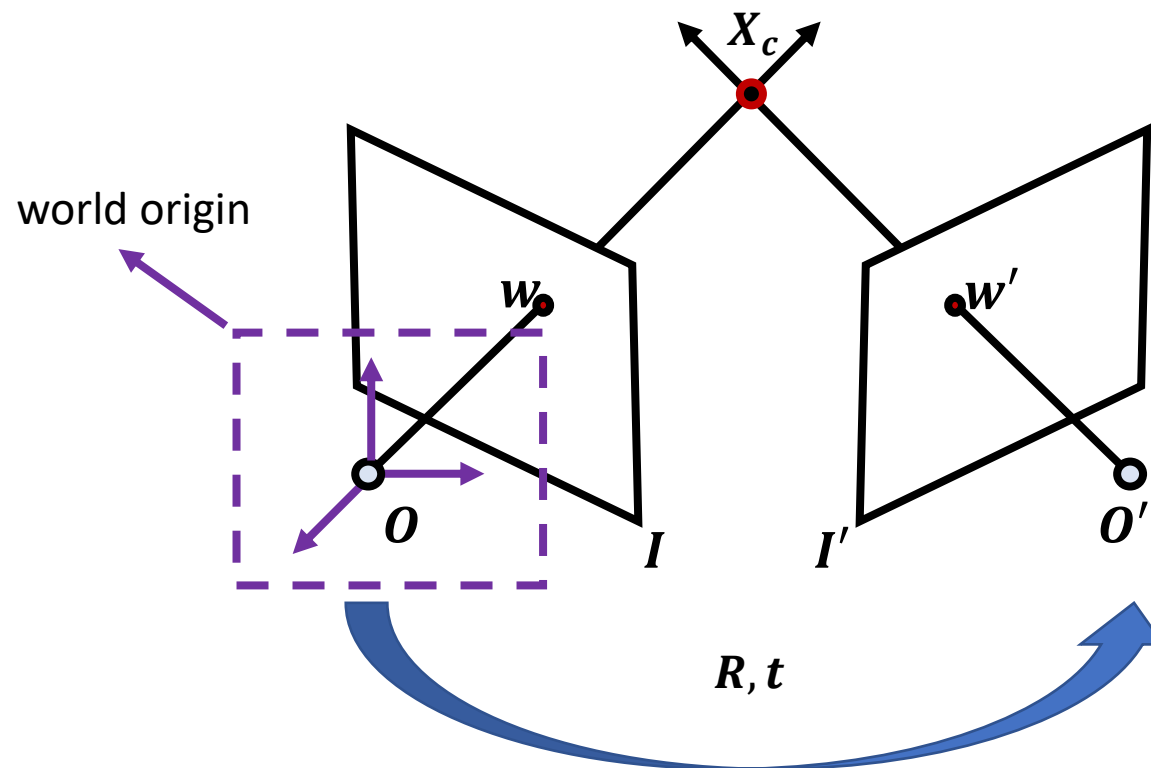
- Firstly, let's assume that the 1<sup>st</sup> camera is located at the **world origin  $O$**
- This means that every point  $X$  is expressed in the camera coordinates of the 1<sup>st</sup> camera, since its optical center coincides with the world origin  $\Rightarrow X \rightarrow X_c$





# Essential Matrix

- If this is the case, the 2<sup>nd</sup> camera is at a location  $\mathbf{O}'$  in world coordinates, which can be expressed by a **translation  $\mathbf{t}$**  and **rotation  $\mathbf{R}$** , w.r.t. 1<sup>st</sup> camera, i.e., the world origin





# Essential Matrix

- Every  $3D$  point  $\mathbf{X}$  can be expressed as  $\mathbf{X}_c$ , which is the camera-centered coordinate system of the 1<sup>st</sup> camera, and as  $\mathbf{X}'_c$ , which is the camera-centered coordinate system of the 2<sup>nd</sup> camera.
- Since  $\mathbf{X} = \mathbf{X}_c$ , we can relate  $\mathbf{X} \leftrightarrow \mathbf{X}'_c$  or  $\mathbf{X}_c \leftrightarrow \mathbf{X}'_c$  using a **Euclidean transformation** composed by the **translation  $\mathbf{t}$**  and **rotation  $\mathbf{R}$**  of the 2<sup>nd</sup> camera, w.r.t. the 1<sup>st</sup> camera

# Essential Matrix

Here is how to find an expression for the epipolar line:

$$\begin{aligned}
 \widetilde{X}'_c &= P_e \widetilde{X}_c \\
 \Leftrightarrow X'_c &= RX_c + t \\
 \Leftrightarrow t \times X'_c &= t \times RX_c + \cancel{t \times t}^0 \\
 \Leftrightarrow X'_c \cdot (t \times X'_c) &= X'_c \cdot (t \times RX_c) \\
 \Leftrightarrow \mathbf{0} &= X'_c \cdot (t \times RX_c)
 \end{aligned}$$

- ← This is in homogeneous coordinates
- ← Express it in cartesian coordinates
- ← Apply cross product with  $t$  to both sides
- ← Apply dot product with  $X'_c$  to both sides

# Essential Matrix

This can be rewritten in matrix form:

$$\begin{aligned} \mathbf{X}'_c \cdot (\mathbf{t} \times \mathbf{R}\mathbf{X}_c) &= \mathbf{0} \\ \Leftrightarrow \mathbf{X}'_c{}^T \mathbf{E} \mathbf{X}_c &= \mathbf{0} \end{aligned}$$

Where  $\mathbf{E} = \mathbf{T}_\times \mathbf{R}$  is the **essential matrix**, and

$$\mathbf{T}_\times = \begin{bmatrix} 0 & -t_z & t_y \\ t_z & 0 & -t_x \\ -t_y & t_x & 0 \end{bmatrix}$$

is a matrix representing the cross product with  $\mathbf{t}$  such that  $\mathbf{t} \times \mathbf{v} = \mathbf{T}_\times \mathbf{v}$

# Fundamental Matrix

Recall that for  $\mathbf{X}_c \leftrightarrow \mathbf{w}$ :  $\tilde{\mathbf{w}} = \mathbf{K}\mathbf{X}_c \Leftrightarrow \mathbf{X}_c = \mathbf{K}^{-1}\tilde{\mathbf{w}}$

and similarly, for  $\mathbf{X}'_c \leftrightarrow \mathbf{w}'$ :  $\tilde{\mathbf{w}}' = \mathbf{K}'\mathbf{X}'_c \Leftrightarrow \mathbf{X}'_c = \mathbf{K}'^{-1}\tilde{\mathbf{w}}'$

Combining the two equations yields the equation of **the two epipolar lines in pixel coordinates**:

$$\begin{aligned} \mathbf{X}'^T \mathbf{E} \mathbf{X}_c &= 0 \\ \Rightarrow (\mathbf{K}'^{-1}\tilde{\mathbf{w}}')^T \mathbf{E} (\mathbf{K}^{-1}\tilde{\mathbf{w}}) &= 0 \\ \Rightarrow \tilde{\mathbf{w}}'^T (\mathbf{K}'^{-T} \mathbf{E} \mathbf{K}^{-1}) \tilde{\mathbf{w}} &= 0 \\ \Rightarrow \tilde{\mathbf{w}}'^T \mathbf{F} \tilde{\mathbf{w}} &= 0 \end{aligned}$$

where  $\mathbf{F} = \mathbf{K}'^{-T} \mathbf{E} \mathbf{K}^{-1}$  is the fundamental matrix.

# Fundamental Matrix

This is the equation of the epipolar line in either camera:

$$\tilde{\mathbf{w}}'^T \mathbf{F} \tilde{\mathbf{w}} = 0$$

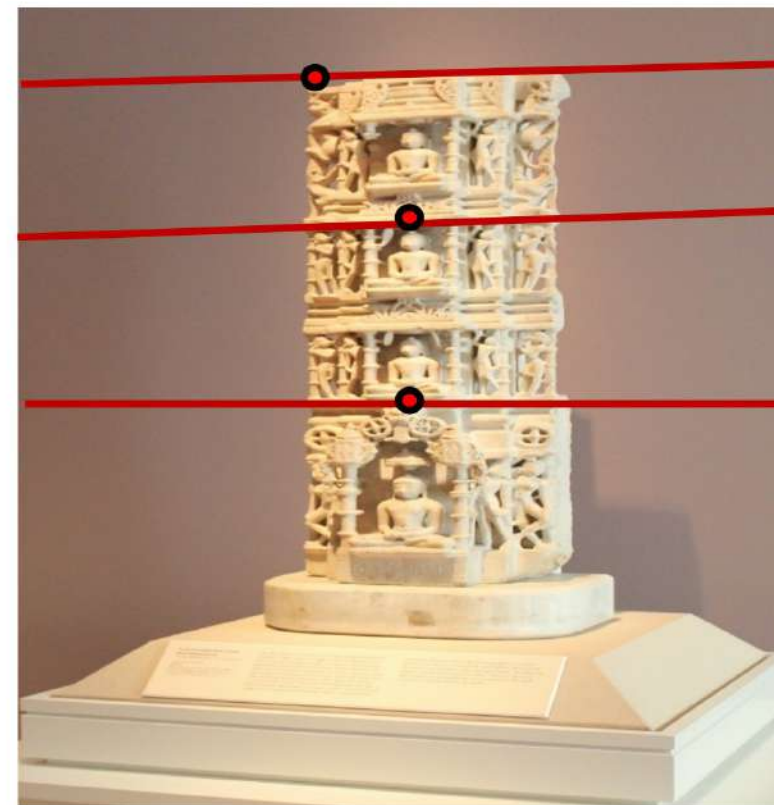
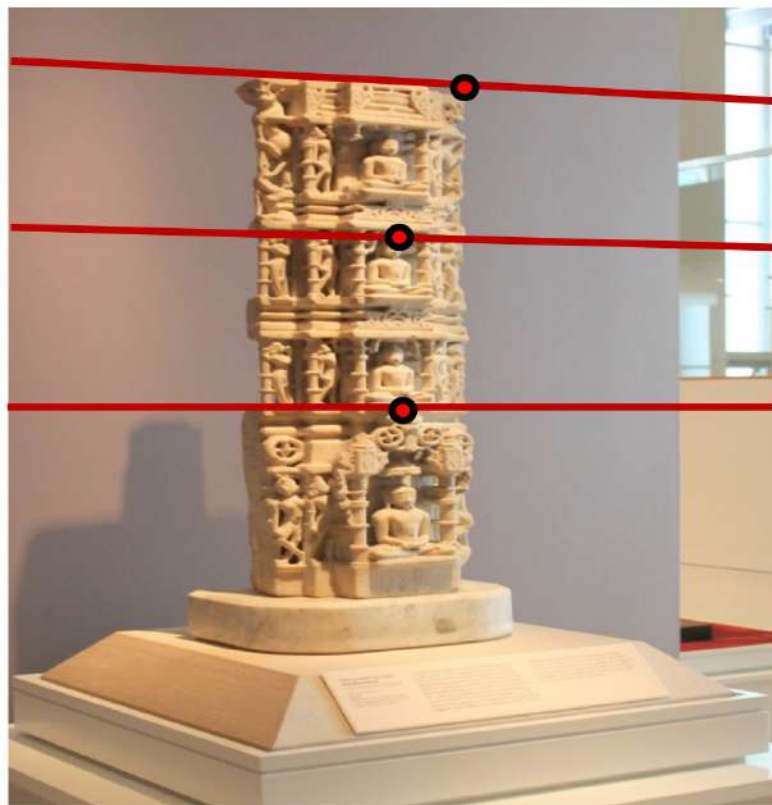
Assuming we know the fundamental matrix, for every point  $\mathbf{w}$  in image 1, this expression gives us the line in image 2 on which the corresponding  $\mathbf{w}'$  must lie, and *vice versa*.

- $l' = \mathbf{F}\tilde{\mathbf{w}}$  is the epipolar line in the 2<sup>nd</sup> image, associated with  $\mathbf{w}$
- $l = \mathbf{F}^T \tilde{\mathbf{w}}'$  is the epipolar line in the 1<sup>st</sup> image, associated with  $\mathbf{w}'$
- $l_i: ax + by + c = 0$



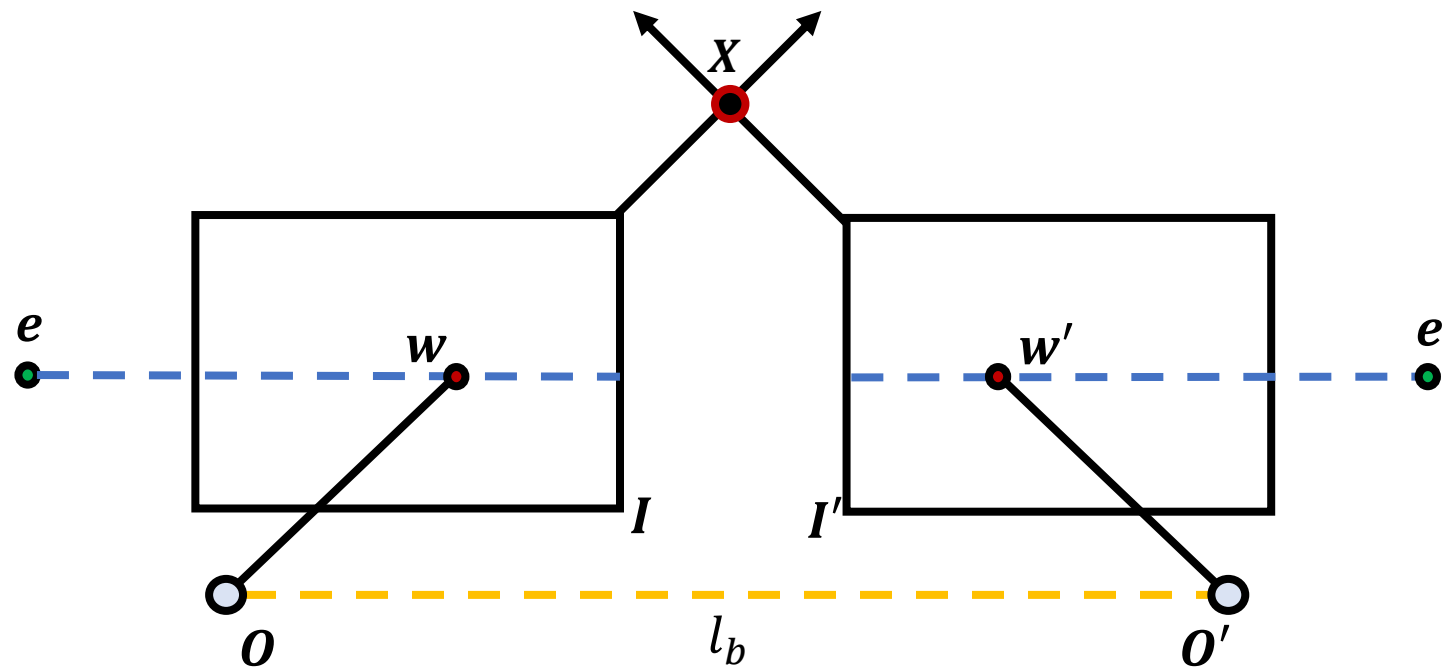
# Examples

Here are a few examples of the epipolar constraint



# Examples

Epipolar constraint examples: Parallel image planes



- **Baseline** intersects the image planes at infinity
- **Epipoles** are at infinity
- **Epipolar lines** are parallel to the  $u$ -axis of each image plane

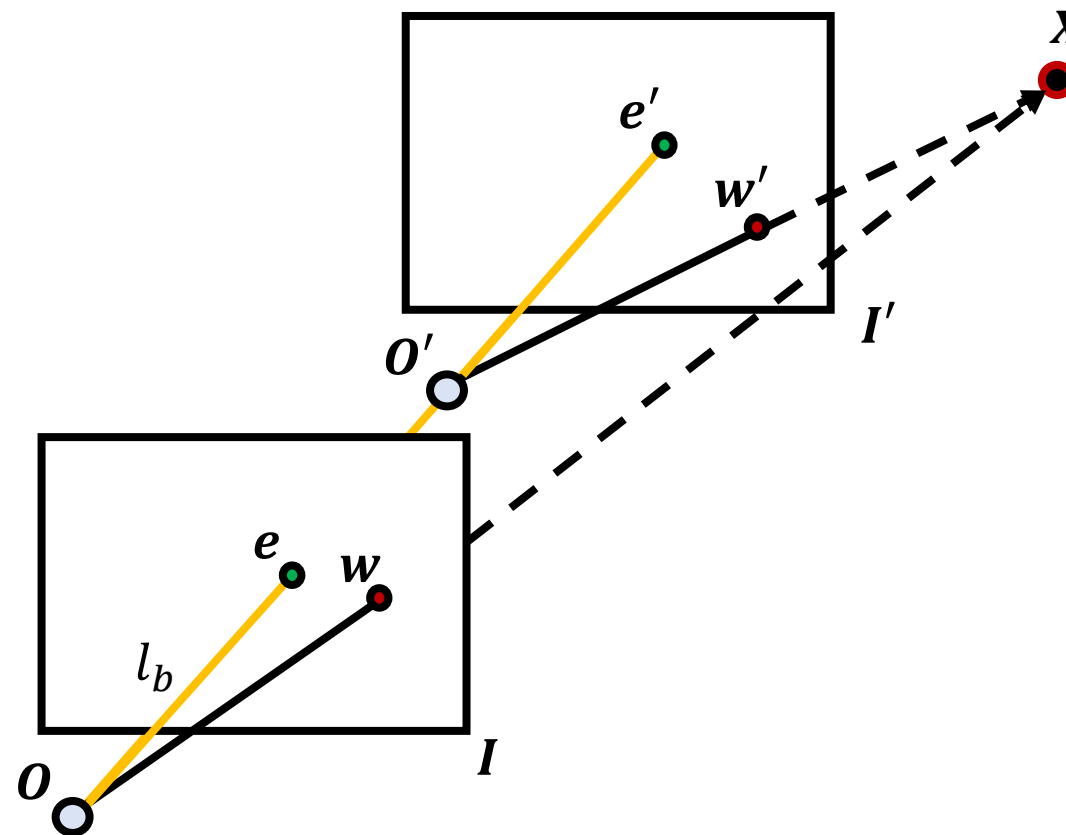
# Examples

Epipolar constraint examples: **Parallel image planes**



# Examples

## Epipolar constraint examples: Forward translation

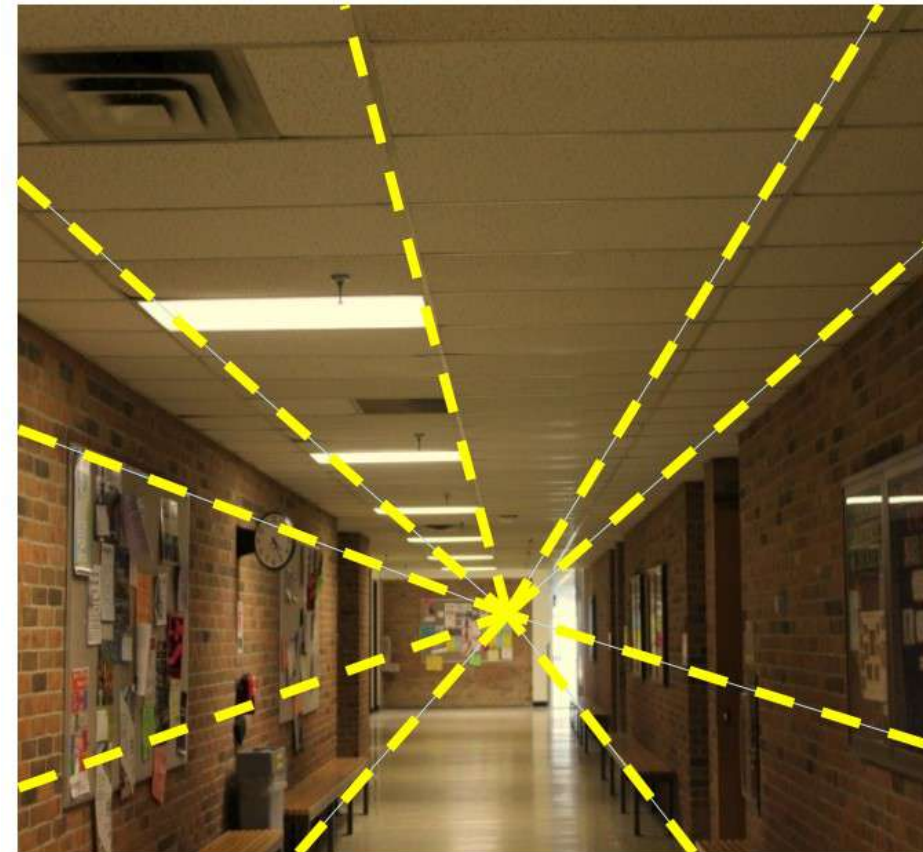
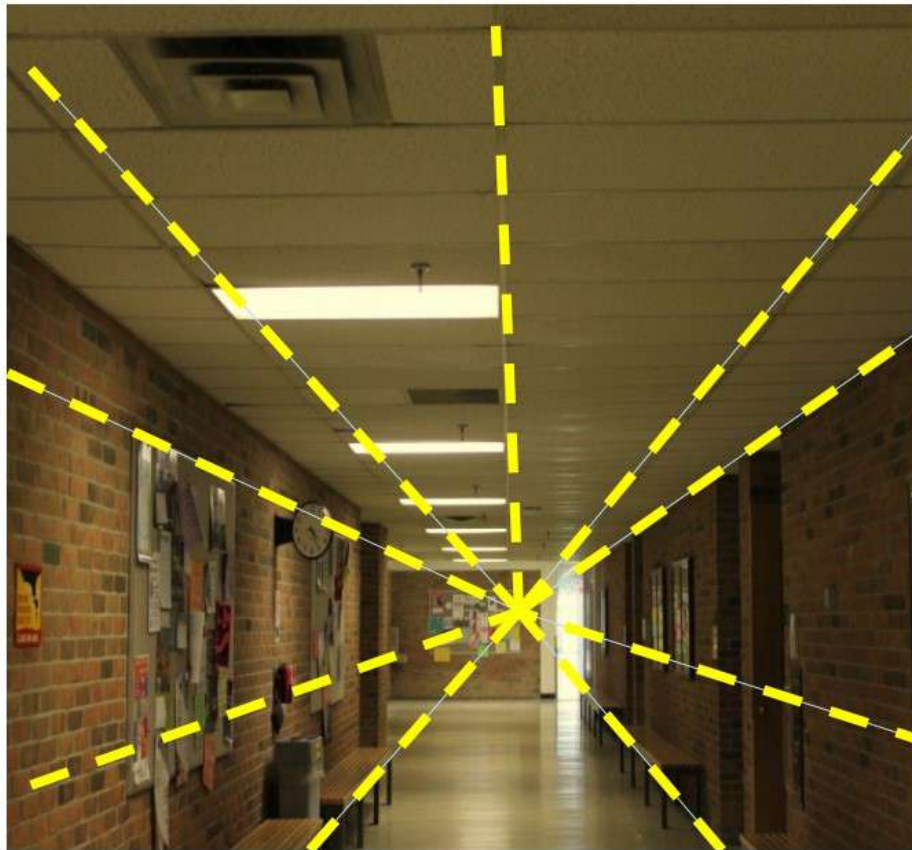


- The **epipoles** have the same position in both images



# Examples

## Epipolar constraint examples: Forward translation





# Fundamental Matrix

In order to apply the epipolar constraint we need to know the fundamental matrix:

$$\begin{aligned}\mathbf{F} &= \mathbf{K}'^{-T} \mathbf{E} \mathbf{K}^{-1} \\ &= \mathbf{K}'^{-T} \mathbf{T}_{\times} \mathbf{R} \mathbf{K}^{-1}\end{aligned}$$

All the parameters of  $\mathbf{F}$  come from the calibration of the two cameras.

Remember that the perspective camera model  $\mathbf{P}_{ps} = \mathbf{K}[\mathbf{R}|\mathbf{T}]$  contains this information.

If, however, these are not available (e.g., because we used a projective camera model to calibrate our cameras),  $\mathbf{F}$  must be estimated using known image correspondences.

# Estimating the fundamental matrix

To estimate the fundamental matrix, we follow a similar approach as in camera calibration

$$\tilde{\mathbf{w}}'^T \mathbf{F} \tilde{\mathbf{w}} = 0 \Rightarrow [u' \ v' \ 1] \begin{bmatrix} f_{11} & f_{12} & f_{13} \\ f_{21} & f_{22} & f_{23} \\ f_{31} & f_{32} & f_{33} \end{bmatrix} \begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = 0$$

Each point correspondence  $\mathbf{w} \leftrightarrow \mathbf{w}'$  generates a single equation for estimating the parameters inside  $\mathbf{F}$ .

# Estimating the fundamental matrix

Here are  $N$  such correspondences:

$$\begin{bmatrix} u_1 u'_1 & v_1 u'_1 & u'_1 & u_1 v'_1 & v_1 v'_1 & v'_1 & u_1 & v_1 \\ & & \vdots & & & & & \\ u_N u'_N & v_N u'_N & u'_N & u_N v'_N & v_N v'_N & v'_N & u_N & v_N \end{bmatrix} \begin{bmatrix} f_{11} \\ f_{12} \\ f_{13} \\ f_{21} \\ f_{22} \\ f_{23} \\ f_{31} \\ f_{32} \end{bmatrix} = \begin{bmatrix} -1 \\ \vdots \\ -1 \end{bmatrix}$$

The minimal number for  $N$  is 8, because  $\mathbf{F}$  has 8 degrees of freedom (we can set  $f_{33} = 1$ )

# Estimating the fundamental matrix

- In practise, 8 clearly indicated correspondences are never available
- What is available is a set of correspondences proposed by SIFT, for a large set of features, in which there are **always** errors

# Estimating the fundamental matrix

- We can use RANSAC to solve this problem
  1. Obtain 8 random SIFT correspondences. Use them to estimate  $\mathbf{F}$ .
  2. For every feature  $\mathbf{w}_i$  in image 1 calculate the epipolar line in image 2. Check if the corresponding feature  $\mathbf{w}'_j$  in image 2 proposed by SIFT falls “close” to the epipolar line. Count the number  $S$  of such “inliers”.
  3. If  $S \geq T$  where  $T$  is a threshold, then there is consensus with the random sample taken in the first step. Calculate  $\mathbf{F}$  for all inliers and terminate here.
  4. If  $S < T$  then no consensus is reached. Repeat from step 1.
  5. If after  $N$  iterations no consensus is reached, select the model that gave the highest  $S$ , calculate  $\mathbf{F}$  using all inliers in  $S$  and terminate



**MAI4CAREU**

Master programmes in Artificial  
Intelligence 4 Careers in Europe



# Thank you.

