

Foundations of Computer Vision

Introduction

EECS 504 Winter 2026

Instructor: Jason Corso
jjcorso@umich.edu



Your turn!

- Form groups of 4.
- You are a new “Computer Vision Startup Design Team”!
- Select one of the three tasks.
- Solve the task using the knowledge you currently have about computer vision, if any.
- What are the challenges? What is easy?
- Discuss for 5 minutes.



Lane Following



Pedestrian Avoidance



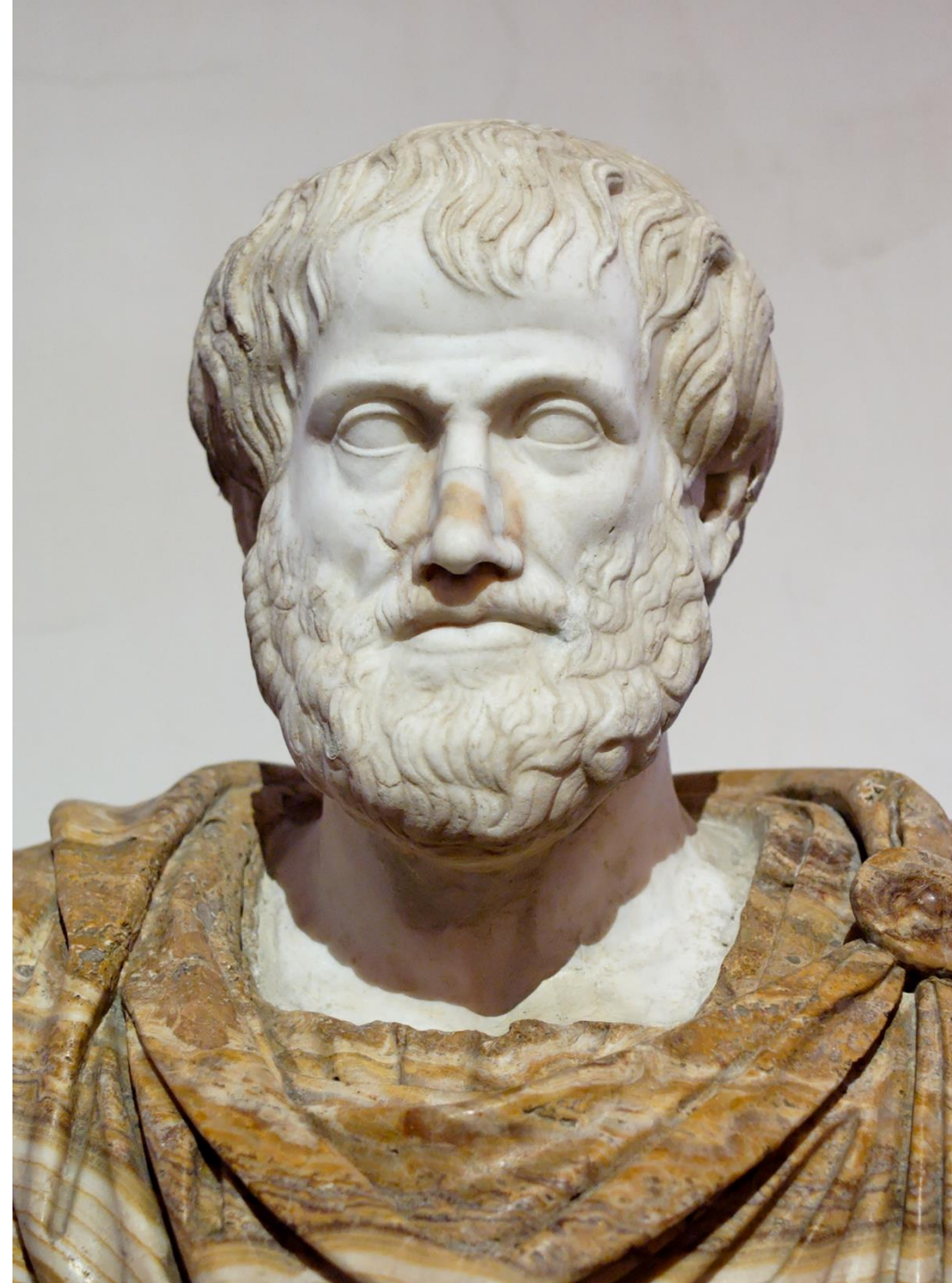
Intersection Detection

What is Computer Vision?

What is ~~Computer~~ Vision?

What is ~~Computer~~ Vision?

- **Aristotle, *De Anima* (On the soul)**
 - "...when once [the sensor, the eye] has been acted upon, it is similar and has the same character as the sensible object."
- **Contrast with modern definition:**
 - "to know what is where by looking."
 - "vision is the *process* of discovering from images what is present in the world, and where it is."
 - Both quotes from Marr 1982.



Vision: A Computational Investigation into the Human Representation and Processing of Visual Information



By David Marr

The MIT Press

DOI: <https://doi.org/10.7551/mitpress/9780262514620.001.0001>

ISBN electronic: 9780262289610

In Special Collection: CogNet

Publication date: 2010

Available again, an influential book that offers a framework for understanding visual perception and considers fundamental questions about the brain and its functions.

David Marr's posthumously published *Vision* (1982) influenced a generation of brain and cognitive scientists, inspiring many to enter the field. In *Vision*, Marr describes a general framework for understanding visual perception and touches on broader questions about how the brain and its functions can be studied and understood. Researchers from a range of brain and cognitive sciences have long valued Marr's creativity, intellectual power, and ability to integrate insights and data from neuroscience, psychology, and computation. This MIT Press edition makes Marr's influential work available to a new generation of students and scientists.

In Marr's framework, the process of vision constructs a set of representations, starting from a description of the input image and culminating with a description of three-dimensional objects in the surrounding environment. A central theme, and one that has had far-reaching influence in both neuroscience and cognitive science, is the notion of different levels of analysis – in Marr's framework, the computational level, the algorithmic level, and the hardware implementation level.

Now, thirty years later, the main problems that occupied Marr remain fundamental open problems in the study of perception. *Vision* provides inspiration for the continuing efforts to integrate knowledge from cognition and computation to understand vision and the brain.

VISION



David Marr

FOREWORD BY
Shimon Ullman

AFTERWORD BY
Tomaso Poggio

What is Computer Vision?

- **Ballard and Brown (1982)**
 - The construction of explicit, meaningful descriptions of physical objects from images.
- **Trucco and Verri (1998)**
 - computing properties of the 3D world from one or more digital images.
- **Stockman and Shapiro (2001)**
 - To make useful decisions about real physical objects and scenes based on sensed images.
- **Forsyth and Ponce (2003)**
 - ...extracting descriptions of the world from pictures or sequences of pictures.
- **Torralba, Isola, and Freeman (2024)**
 - Computer vision studies how to reproduce in a computer the ability to see.
- **Extraction of information from visual content.**

What is Computer Vision?

- **Vision** is the representation and extraction of information from visual content.
- **Computer Vision** is the representation and extraction of information from visual content by a computer system.
 - Or, *vision by a computer* (now that we have adequately defined vision)
- Computer Vision is hence characterized by the study of
 - **Representations** of visual content and *world* that gave rise to it.
 - **Mathematical formalisms and computational algorithms** of the representations, their properties, and the process of information extraction from visual content.
 - **Properties** of the visual content, the production of visual content, and the mathematics and algorithms for processing it. Most important of these is **invariance**.

Some related terms

- **Image Processing:** the study of the properties of operators that produce images from other images
 - we will touch on image filtering and related operators from image processing
- **Machine Vision:** a somewhat outdated term which now tends to refer to industrial vision applications where (usually) a single camera is used to solve a structured inspection task
 - the “reverse CAD” model
- **Pattern Recognition:** typically refers to the recognition of structures in 2D images (usually without reference to any underlying 3D information).
- **Photogrammetry:** the science of measurement through non-contact sensing, e.g. terrain maps from satellite images. Usually is more focused on accuracy issues than interpretation.

What information is in images?



Five classes of information in visual content

1. Early Processes

- extracting basic features, edges, contours, and segmentation.

2. Motion Tracking

- extracting movement, optical flow, tracking, and filtering.

3. Shape

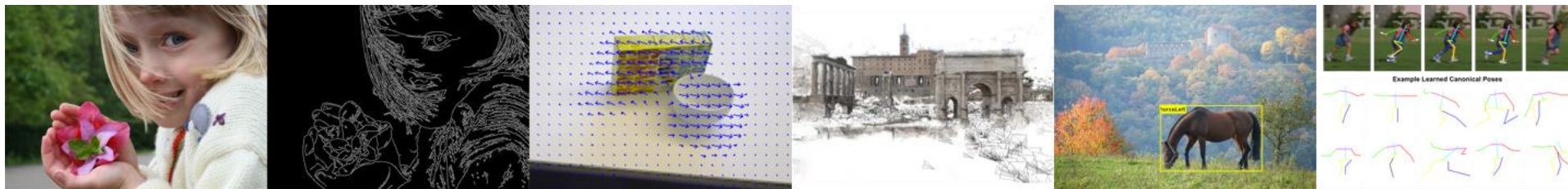
- extracting 2D and 3D structure, epipolar geometry, stereo, SFM, shape from X.

4. Objects

- extracting objects, detection, recognition, and matching.

5. Actions

- extracting actions, space-time localization, detection



Extracting that information is harder than you think



Example of sensors and image formation

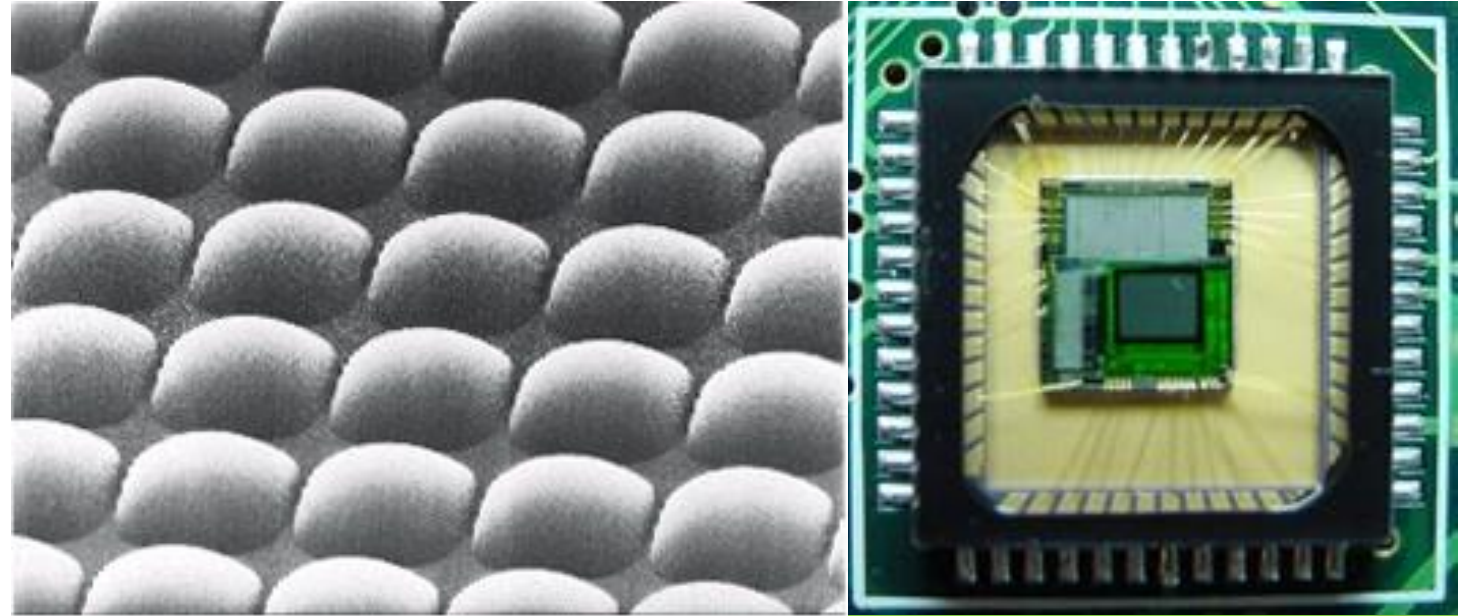


62	70	31	47	100	125	164	166
62	63	40	112	159	140	160	161
50	50	100	143	167	153	150	148
43	73	142	152	165	167	115	114
57	134	170	164	155	114	106	93
111	163	187	144	61	45	50	62
143	180	166	89	51	60	81	176
141	163	105	120	112	99	123	154
167	91	113	135	140	135	135	139

Each pixel is a measure of the brightness (intensity of light) that falls on an area of a sensor (typically a CCD chip)

Example of sensors and image formation

- Basic image sensing process:
 - photons hit a detector
 - the detector becomes charged
 - the charge is read out as brightness
- Sensor types:
 - CCD (charge-coupled device)
 - more common
 - high sensitivity
 - high power
 - cannot be individually addressed
 - blooming
 - CMOS
 - simple to fabricate (cheap)
 - lower sensitivity, lower power
 - can be individually addressed



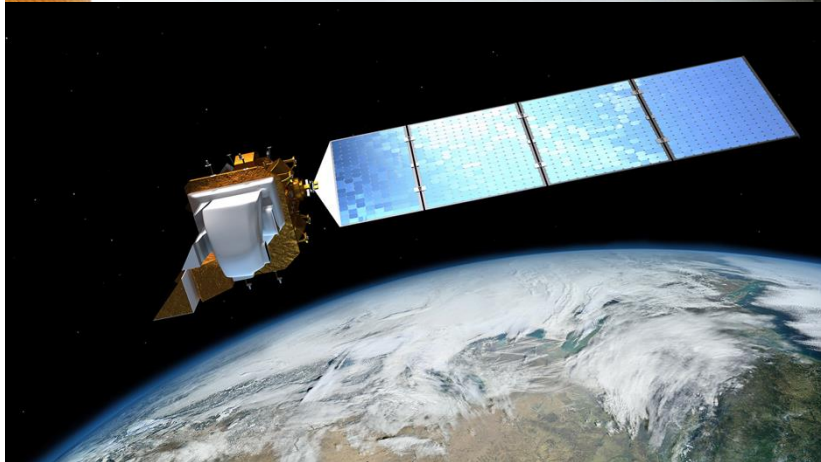
**This core of this
technology has not
changed (to my
knowledge) in decades**

But, visual content can also be produced by

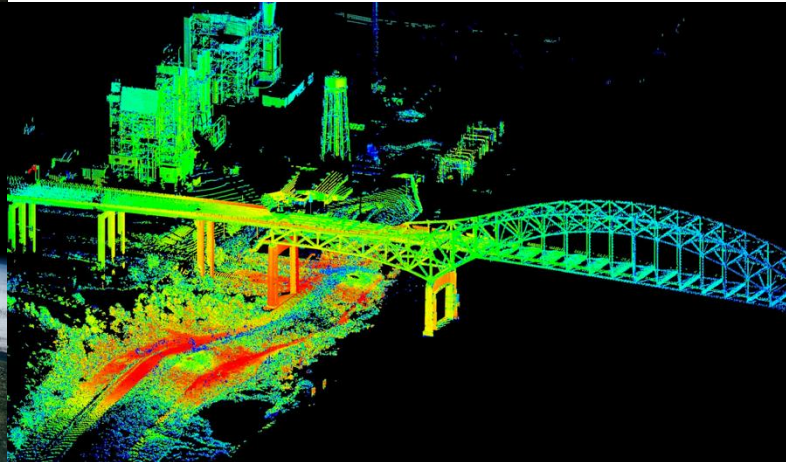
Source: <http://webneel.com/3d-drawings-pencil-art>.



Source: <http://www3.gehealthcare.com/>.



Source: <https://svs.gsfc.nasa.gov/>.



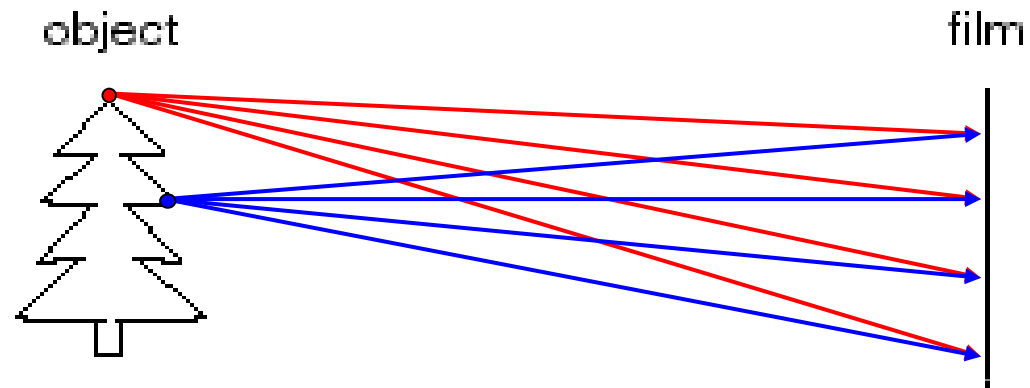
Source: <http://gallery.usgs.gov/>.



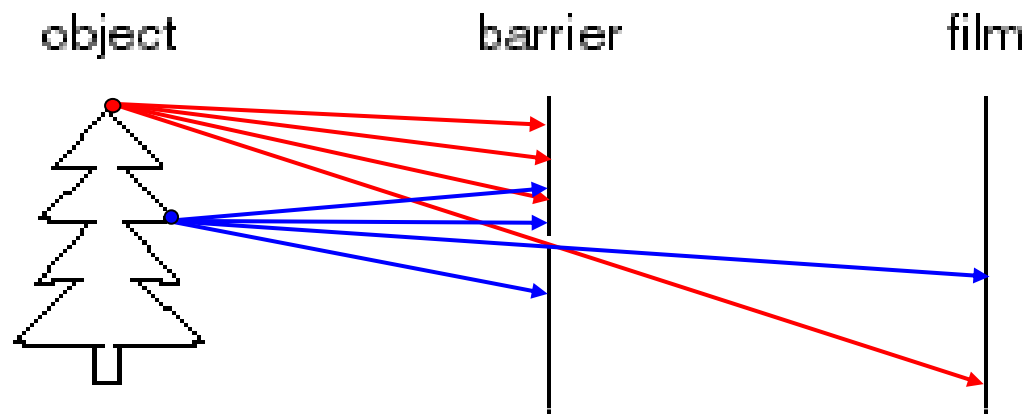
- **And, computer vision is useful in all of these!**

Example of sensors and image formation

- Getting the light to the sensor.



- Add a barrier to block most of the light rays: **aperture**.



The pinhole camera

- Camera obscura
 - First camera; known to Aristotle.
 - Aperture size impacts image.

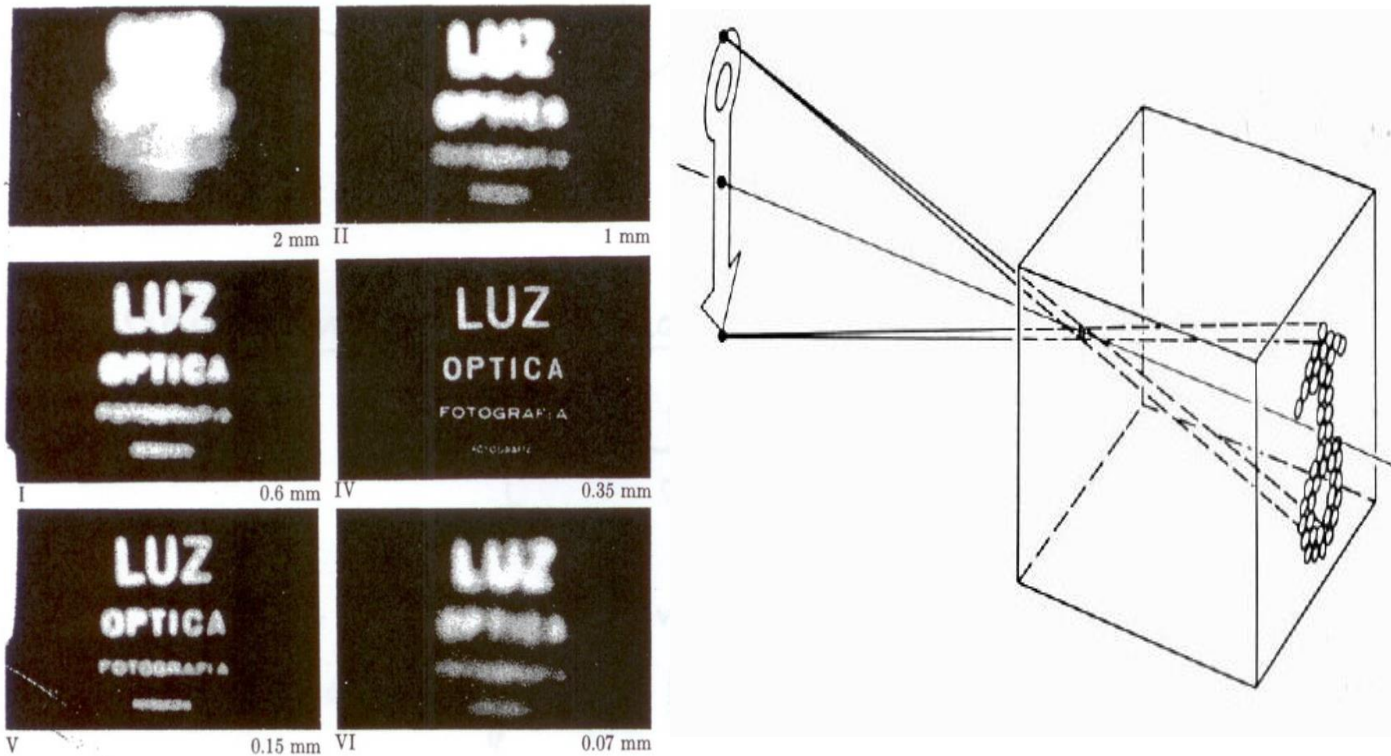
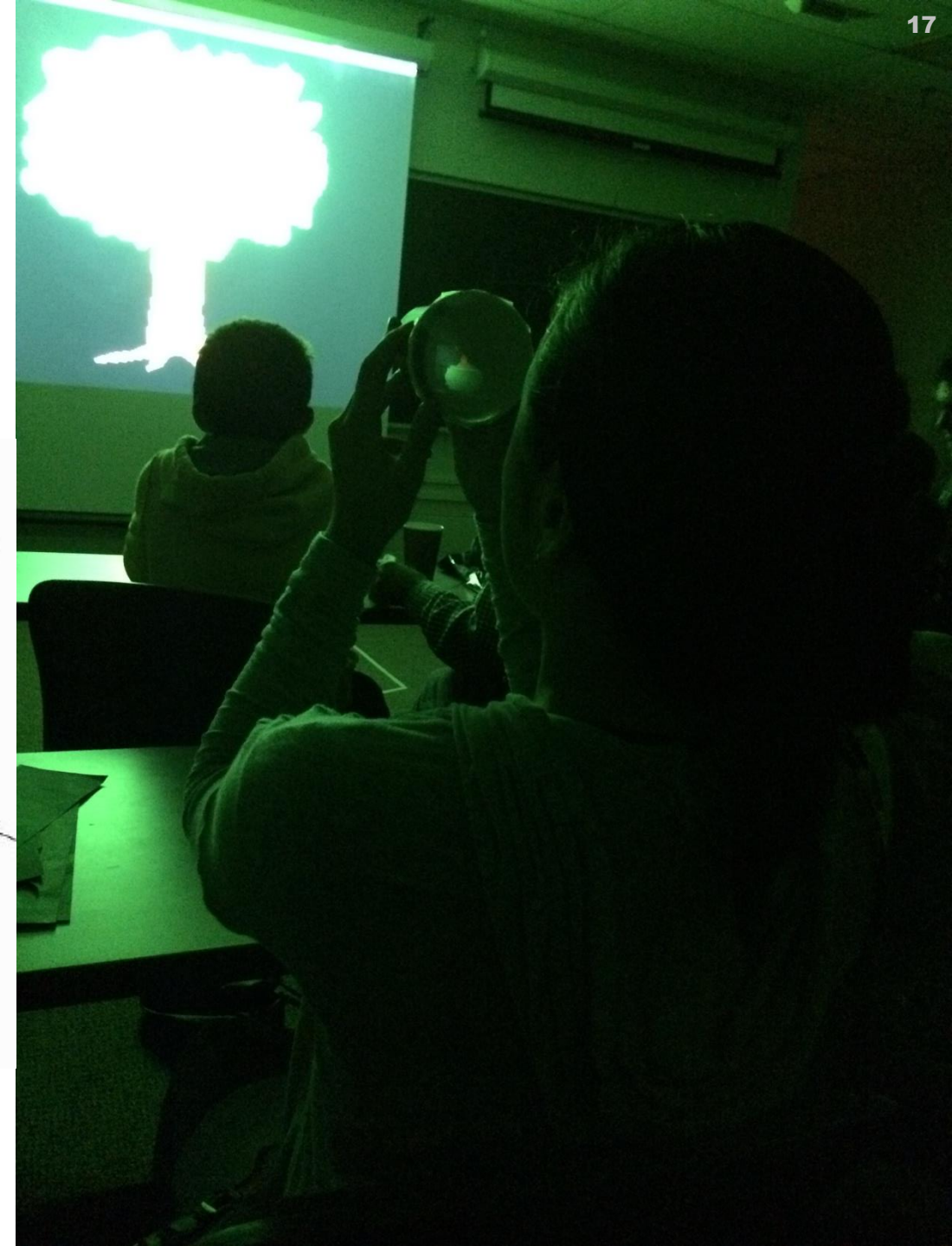


Fig. 5.96 The pinhole camera. Note the variation in image clarity as the hole diameter decreases. [Photos courtesy Dr. N. Joel, UNESCO.]



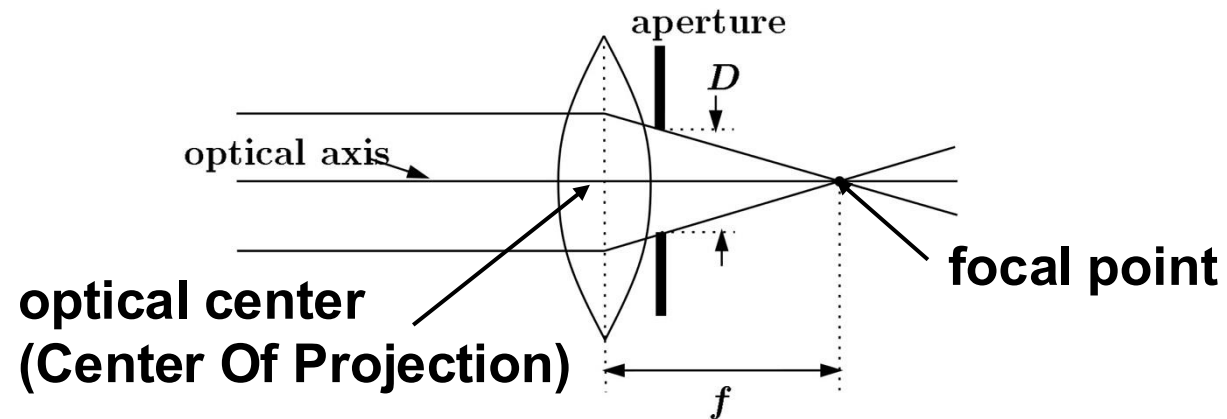
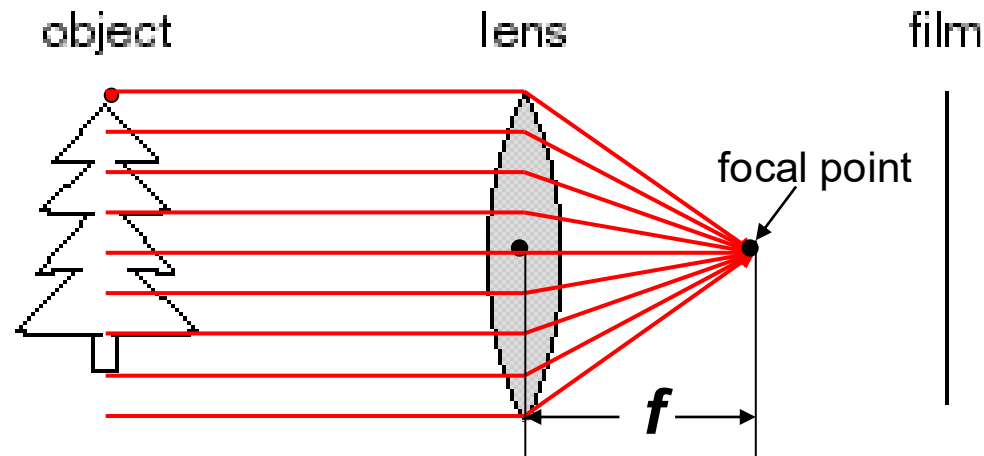
Pinhole Camera

- There are quite a few examples of such videos.
- Maybe you've tried one in real life?
- Maybe you've built one?



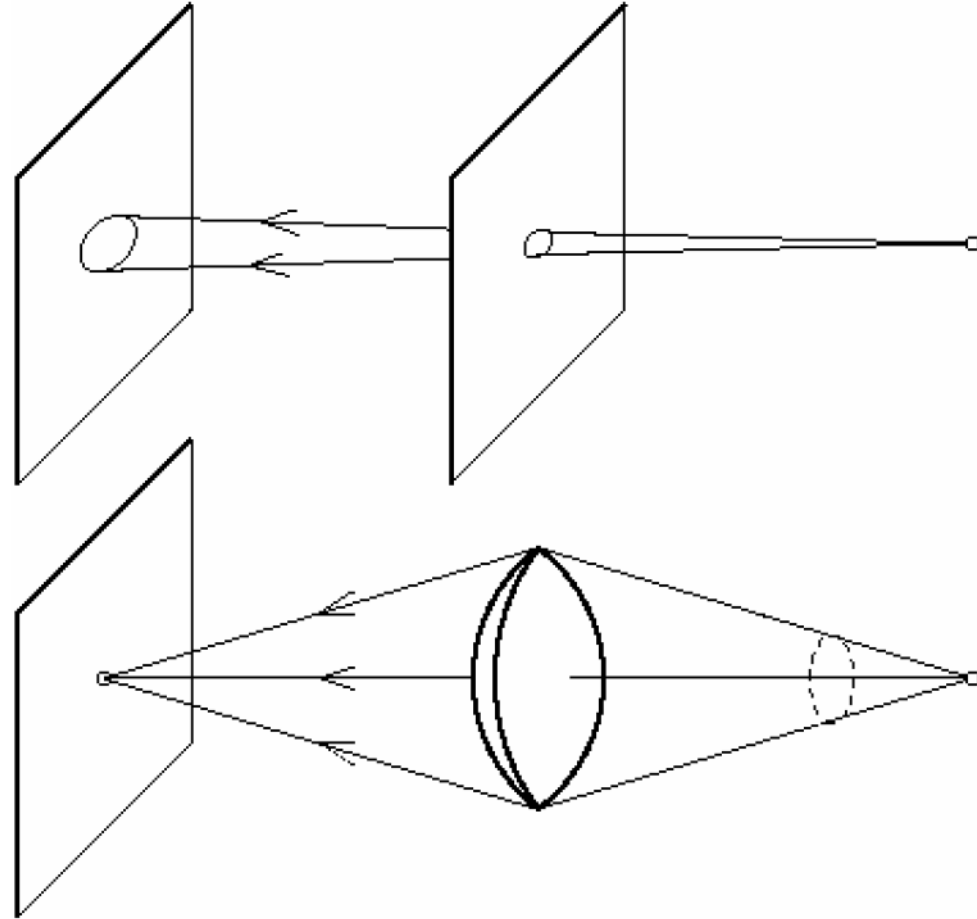
Adding a lens

- A lens focuses the light onto the film/CCD.
- Rays passing through the center are not deviated.
- All parallel rays converge to one point on a plane located at the **focal length f** .

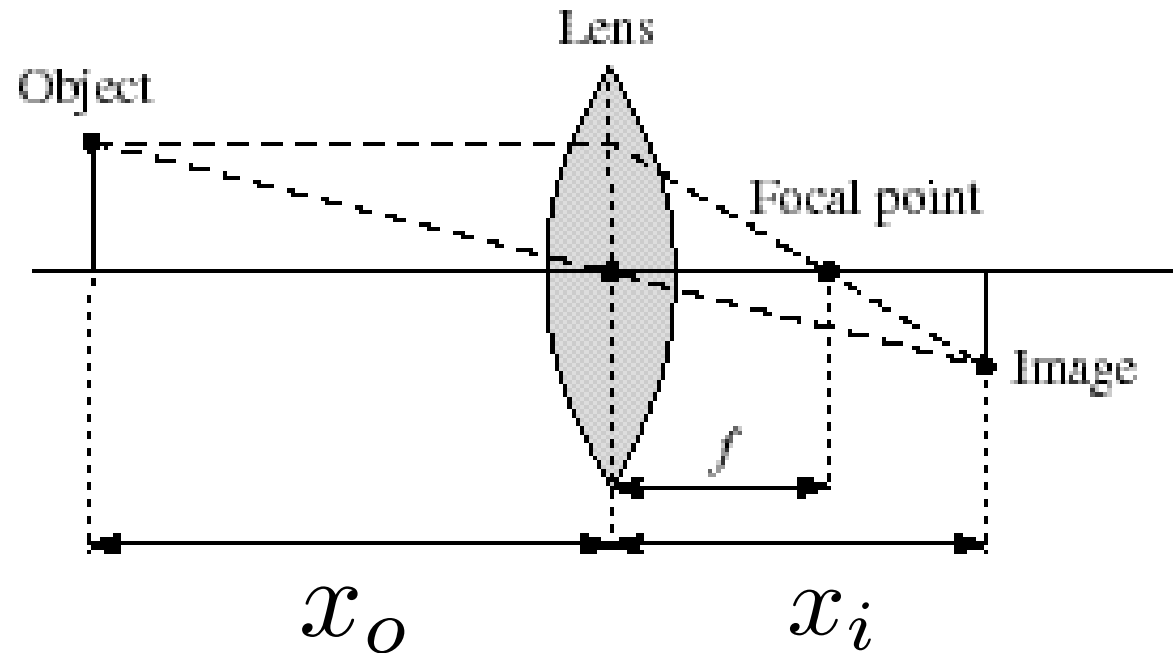


Pinhole vs. lens

- So then, why lenses?

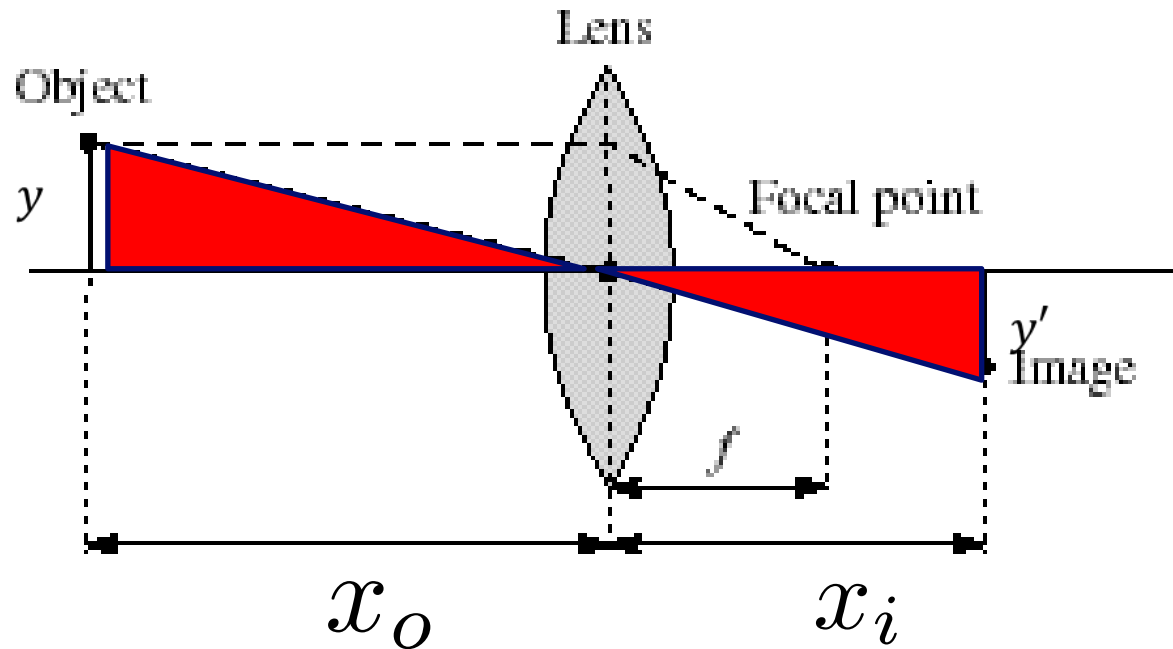


Thin lens equation



- How to relate distance of object from optical center (x_o) to the distance at which it will be in focus (x_i), given focal length f ?

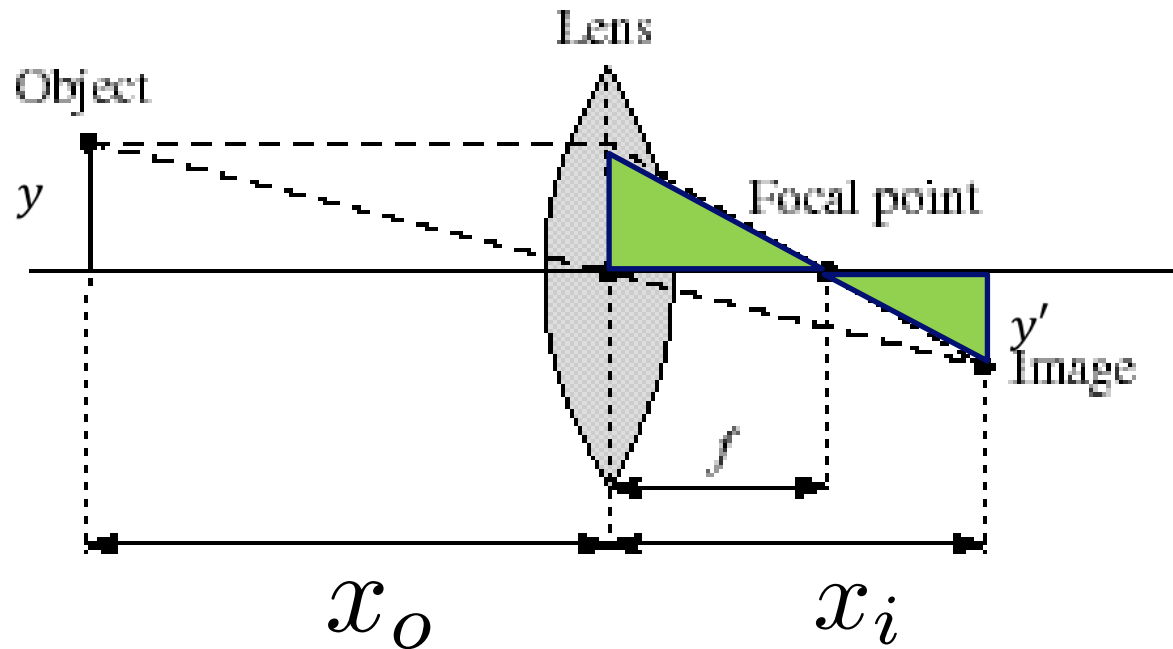
Thin lens equation



$$\frac{y'}{y} = \frac{x_i}{x_o}$$

- How to relate distance of object from optical center (x_o) to the distance at which it will be in focus (x_i), given focal length f ?

Thin lens equation

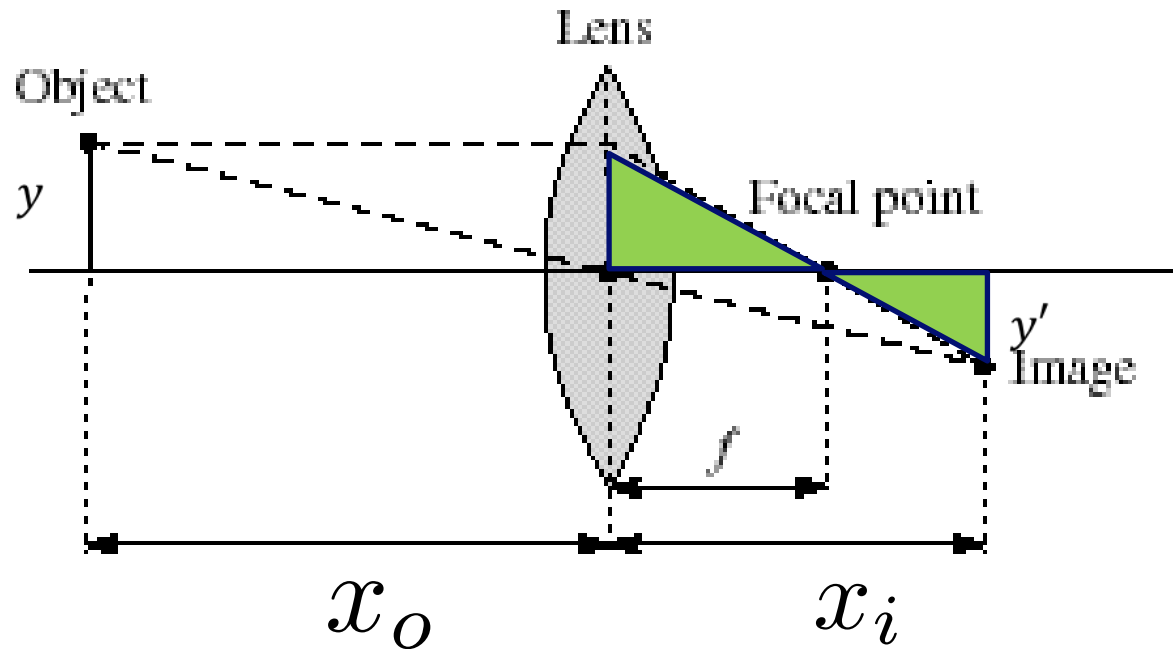


$$\frac{y'}{y} = \frac{x_i}{x_o}$$

$$\frac{y'}{y} = \frac{x_i - f}{f}$$

- How to relate distance of object from optical center (x_o) to the distance at which it will be in focus (x_i), given focal length f ?

Thin lens equation



$$\frac{y'}{y} = \frac{x_i}{x_o}$$

$$\frac{y'}{y} = \frac{x_i - f}{f}$$



$$\frac{1}{f} = \frac{1}{x_o} + \frac{1}{x_i}$$

- Any object point satisfying this equation is in focus

And we have our image...

62	70	31	47	100	125	164	166
62	63	40	112	159	140	160	161
50	50	100	143	167	153	150	148
43	73	142	152	165	167	115	114
57	134	170	164	155	114	106	93
111	163	187	144	61	45	50	62
143	180	166	89	51	60	81	176
141	163	105	120	112	99	123	154
167	91	113	135	140	135	135	139

- Modern digital cameras capture upwards of 25 million pixels per image even at the entry level.



Canon EOS R8 Full-Frame Mirrorless Camera – 24.2MP, 4K60p Video, Dual Pixel AF II, Wi-Fi, Lightweight Body – Body Only (5803C002) + 64GB Memory Card + Canon Shoulder Bag

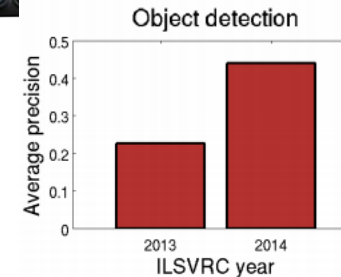
Visit the Canon Store

5.0 ★★★★★ (6) | Search this page

Lowest price in 30 days

-5% \$1,379⁹⁵

Typical price: \$1,459.95 | Price history



- But actual performance still leaves much to be desired for many tasks, such as object detection in the ImageNet large-scale visual recognition challenge.

Object detection



Object detection

Year	Codename	AP (percent)	99.9% Conf Int
2014	GoogLeNet [†]	43.93	42.92 - 45.65
2014	CUHK [†]	40.67	39.68 - 42.30
2014	DeepInsight [†]	40.45	39.49 - 42.06
2014	NUS	37.21	36.29 - 38.80
2014	UvA [†]	35.42	34.63 - 36.92
2014	MSRA	35.11	34.36 - 36.70
2014	Berkeley [†]	34.52	33.67 - 36.12

And human performance is typically but not always the gold standard.

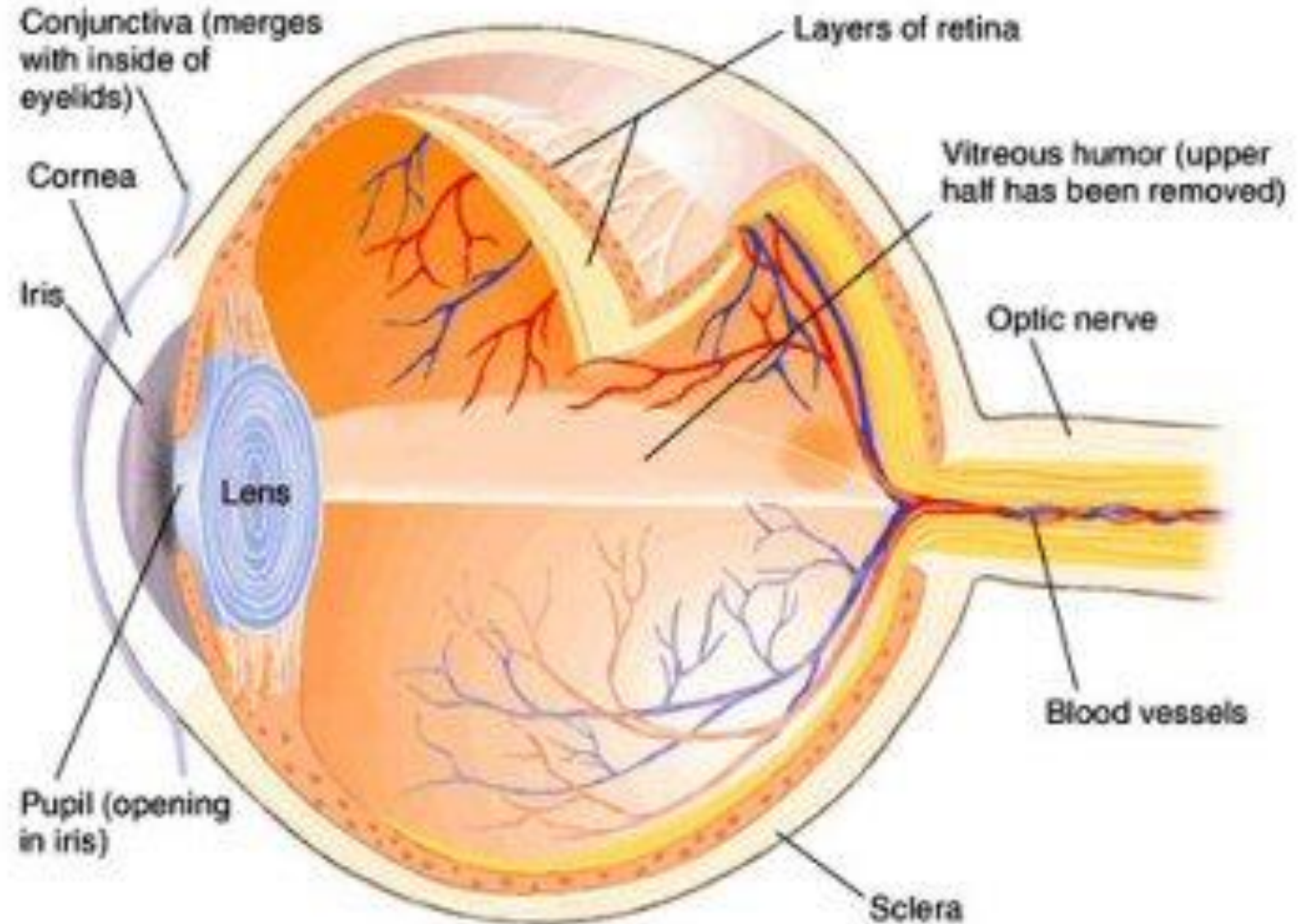
Human vision: a benchmark of sorts

- Is a human eye just that much better than a CCD camera?

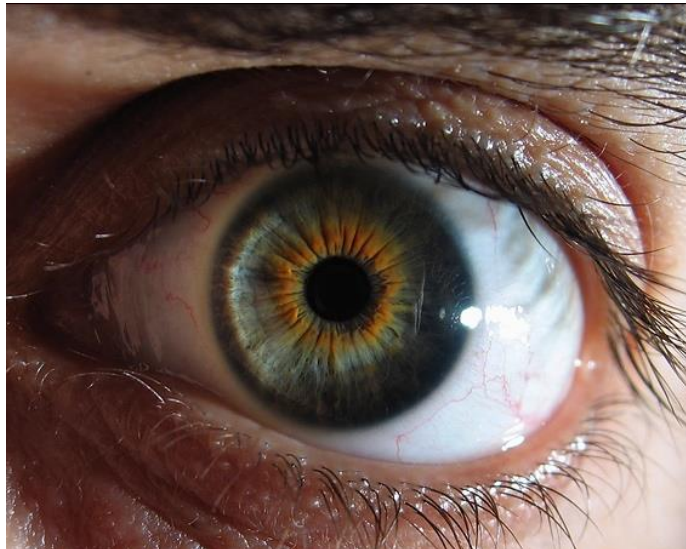
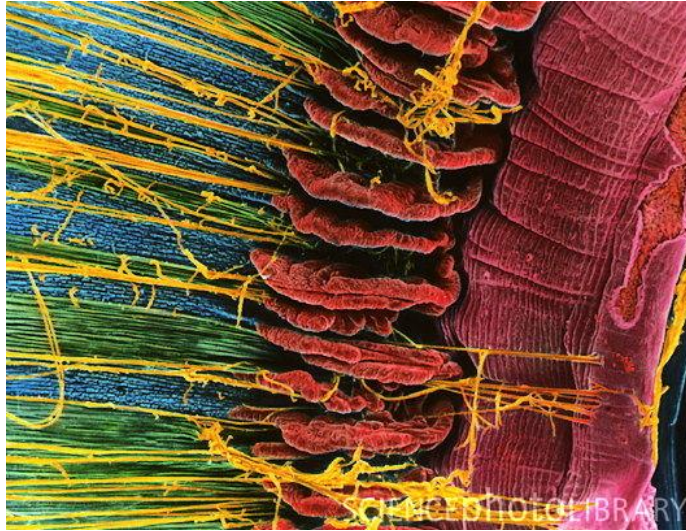


Structure of the eye

The iris is roughly equivalent to the aperture in a camera, the cornea and the lens are both lens-like objects, and the retina is where the image is recorded, similar to a CCD sensor or film.

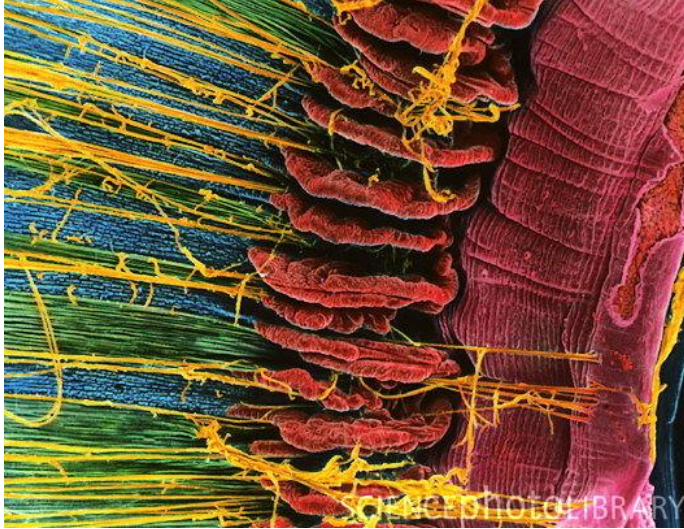


Structure of the eye: Iris



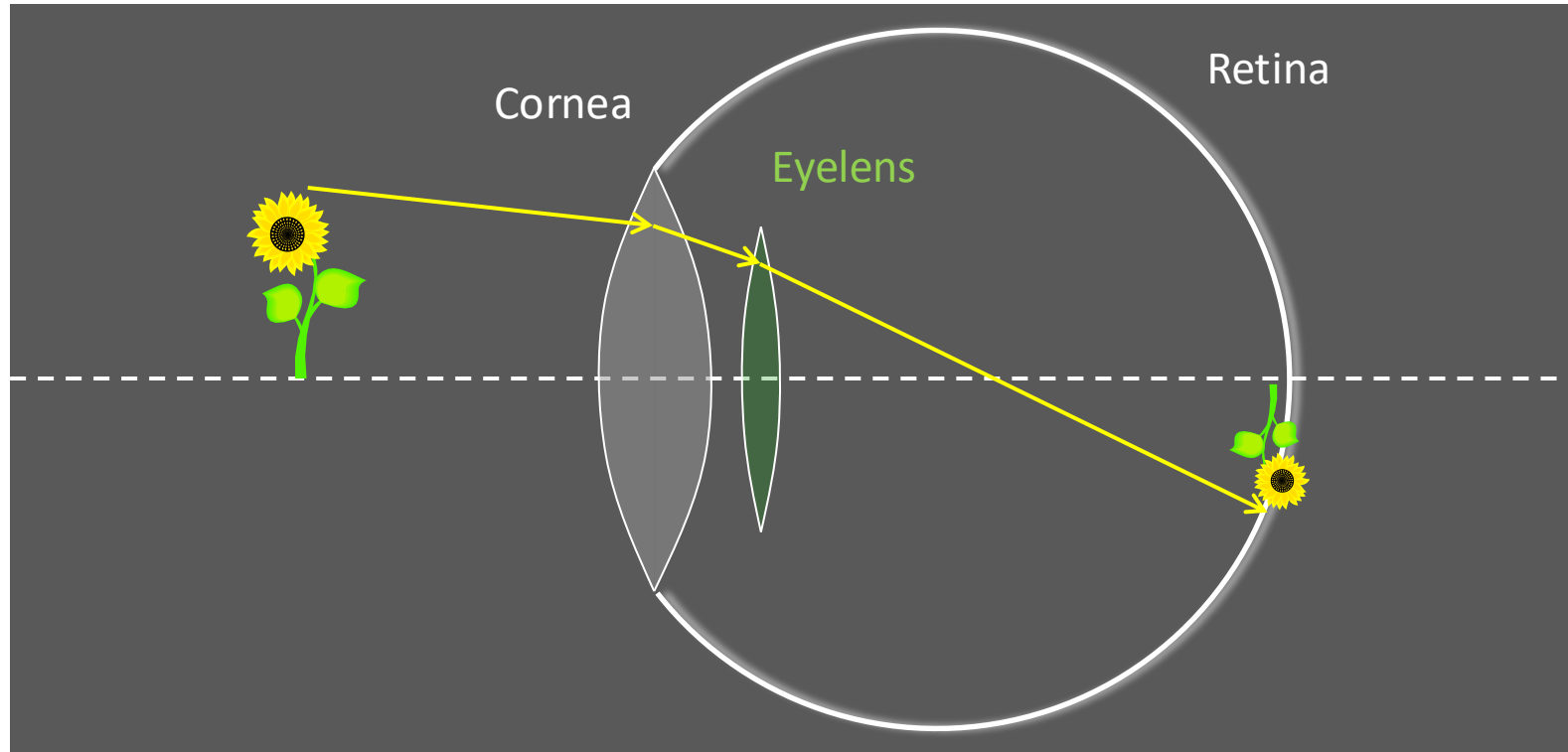
- The iris is similar to the diaphragm/aperture in a camera
- Your iris *widens* in dim light and *narrows* in bright light
- The f-number of your eye varies from **f/2** to **f/8** (large opening) (small opening).
- Compare this to the range of an average camera lens, which may have f-numbers from **f/2.8** to **f/22**.

Structure of the eye: iris



- With a range of only $f/2 - f/8$, your iris can only reduce the light coming into your eye by a factor of 20.
- The range of intensities that your eye can respond to is a factor of 10^{13}
- The main function of the iris is ***not*** to control the intensity of light coming into your eye
- Main functions of iris
 - Reduce aberrations, sharpen image
 - Increase depth of field

Structure of the eye: cornea and lens



- There are two lenses in your eye, the cornea and the eyelens.
- The cornea, the front surface of the eye, does most of the focusing in your eye.
- The eyelens (crystalline lens) provides adjustable fine-tuning of the focus.

How a camera lens focuses

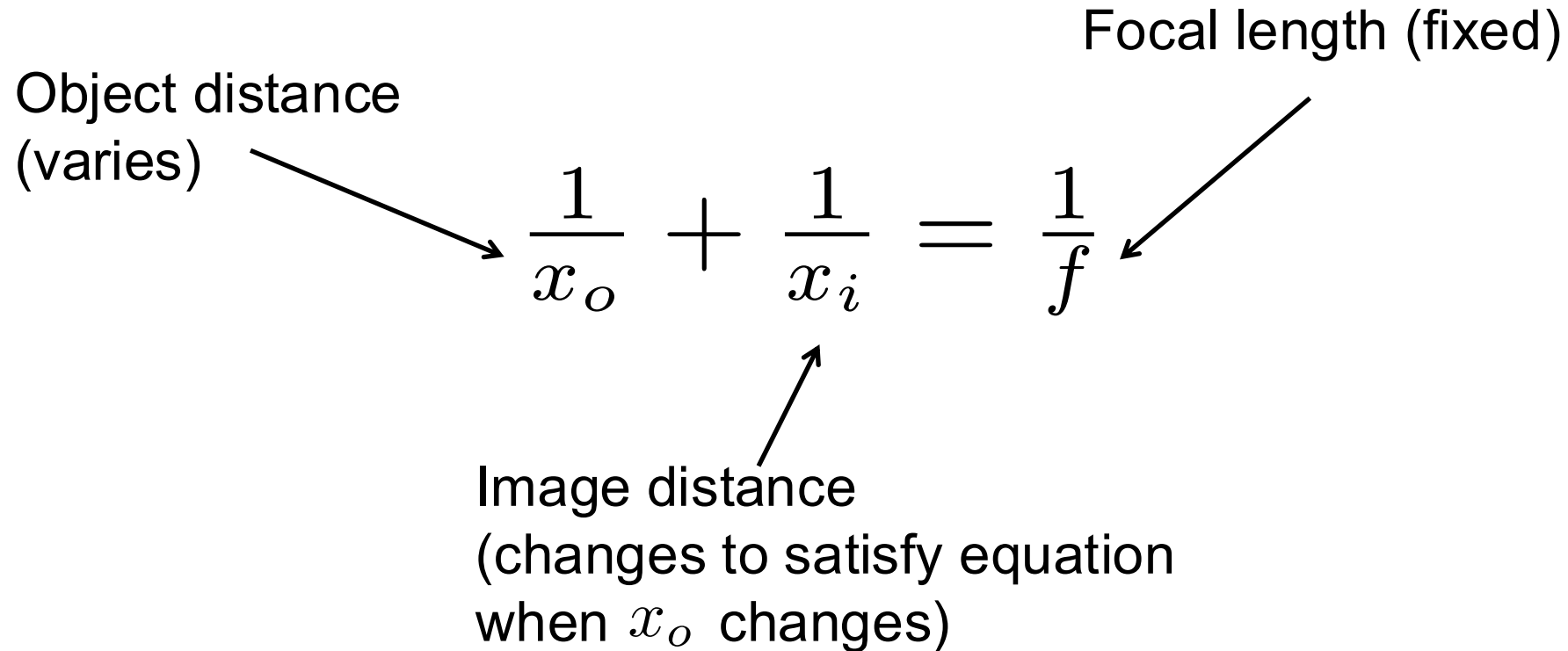
- A typical fixed focal length **camera** is focused by changing the **distance**, x_i , from the lens to the image at the back on the CCD as the distance to the object, x_o , changes.

Object distance
(varies)

Focal length (fixed)

$$\frac{1}{x_o} + \frac{1}{x_i} = \frac{1}{f}$$

Image distance
(changes to satisfy equation
when x_o changes)

The diagram shows the thin lens equation $\frac{1}{x_o} + \frac{1}{x_i} = \frac{1}{f}$. Three arrows point from text labels to the equation: one from 'Object distance (varies)' to $\frac{1}{x_o}$, one from 'Image distance (changes to satisfy equation when x_o changes)' to $\frac{1}{x_i}$, and one from 'Focal length (fixed)' to $\frac{1}{f}$.

How a human eye focuses

- The **eyelens** is a fixed distance, x_i , from the retina at the back of the eyeball where the image is created.

Object distance (varies)

Focal length (changes to satisfy the equation)

$$\frac{1}{x_o} + \frac{1}{x_i} = \frac{1}{f}$$

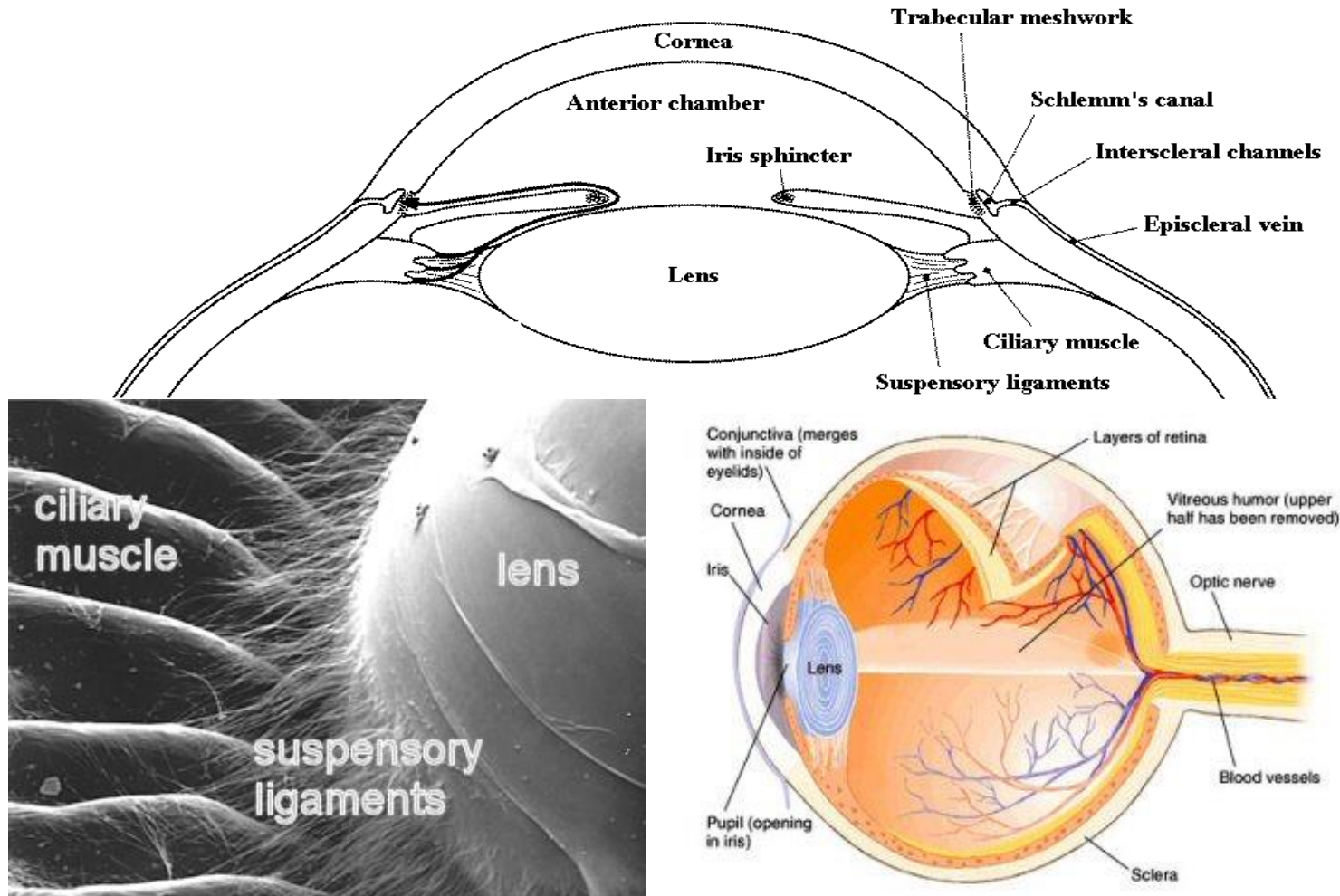
Image distance (fixed)

The diagram shows the thin lens equation $\frac{1}{x_o} + \frac{1}{x_i} = \frac{1}{f}$. Three arrows point from descriptive text to the equation: one from 'Object distance (varies)' to $\frac{1}{x_o}$, one from 'Image distance (fixed)' to $\frac{1}{x_i}$, and one from 'Focal length (changes to satisfy the equation)' to $\frac{1}{f}$.

- Note that many modern cameras do actually automatically adjust the effective focal length as well.

Eyelens: focusing and accommodation

- The eyelens changes its focal length by changing its shape. Ligaments pull on the lens to change the amount of “bulge”.

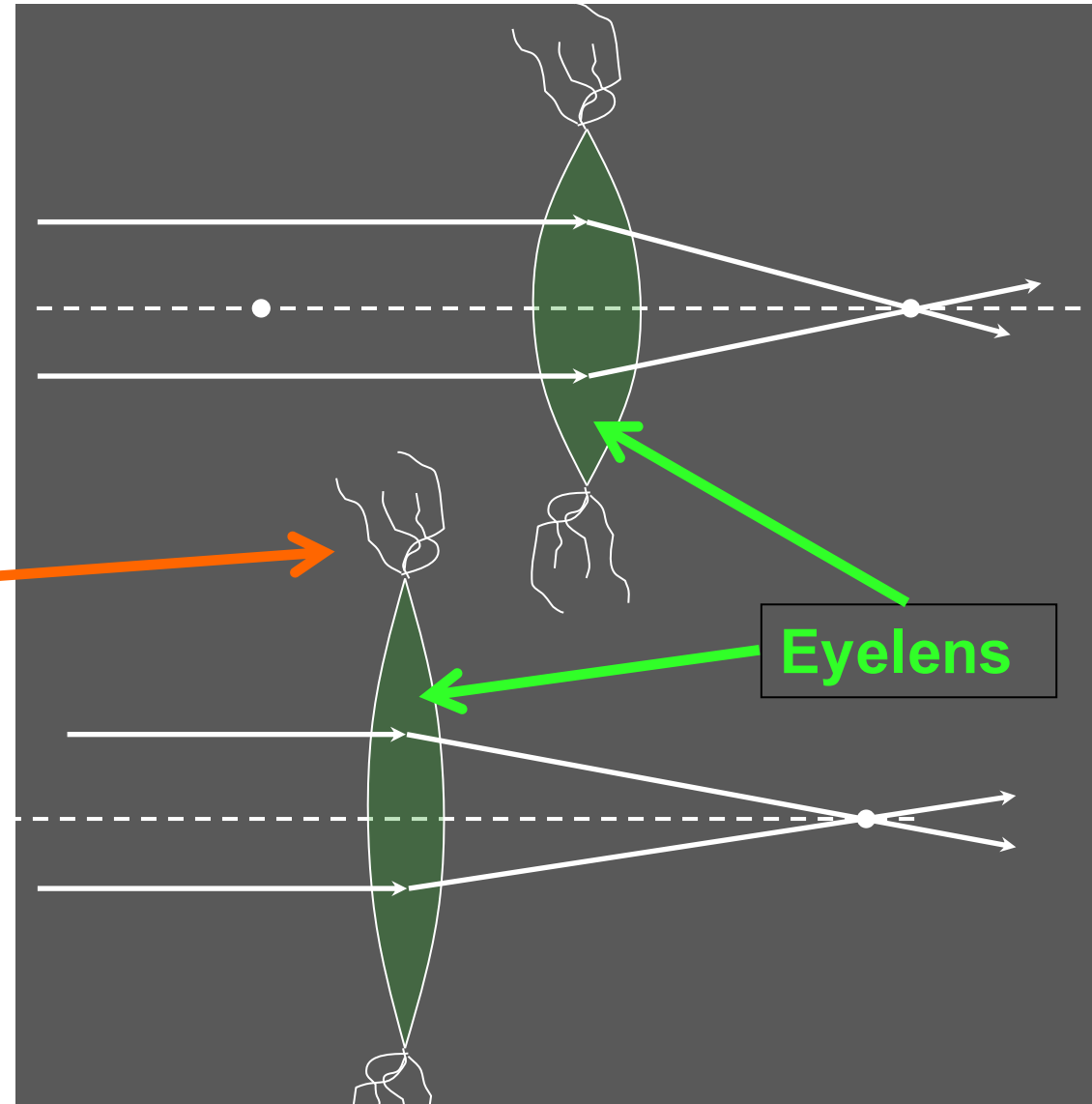


Eyelens: focusing and accommodation

Muscles contract,
ligaments relax, more
bulge, more bending
power, shorter focal
length

Ligaments

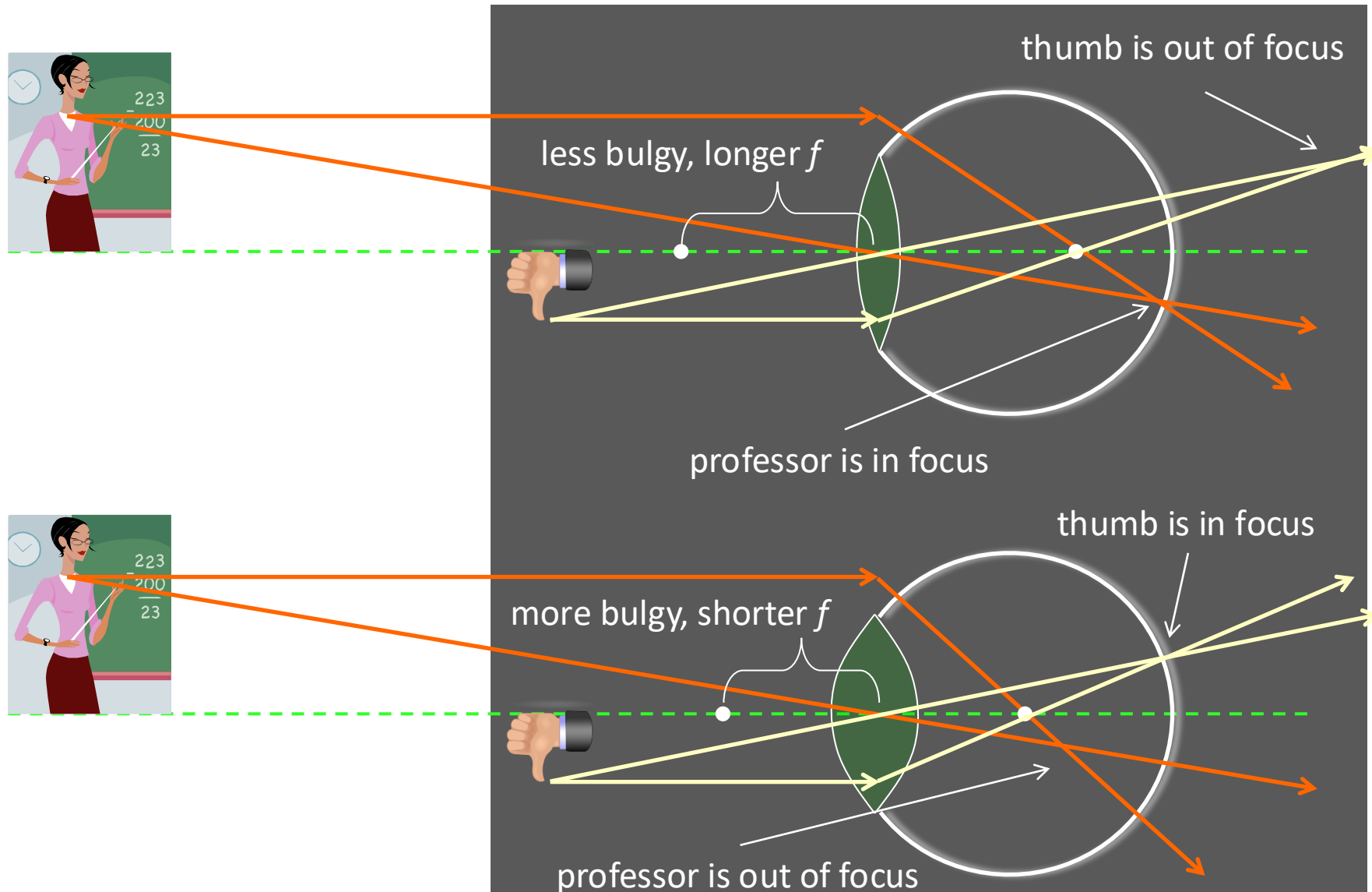
Muscles relax,
ligaments contract,
less bulge, less
bending power, longer
focal length



Eyelens: focusing and accommodation

- Your eyelens has a ***small depth of field***
 - You can't see something close and far with both objects in focus at the same time.
- Hold out your thumb about a foot away from your eye
 - Then, alternately focus on thumb and me (right above your thumb).
- Note that you cannot see ***both*** me ***and*** your thumb sharply (in focus) at the same time
 - You focus on one or the other by changing the bulge of your eyelens.

Eyelens: focusing and accommodation



Concept questions on focusing

You can't see me and your thumb clearly *at the same time*

- a) because your *pupil* is too small
- b) because your iris can't change fast enough
- c) because your eye cannot *accommodate*
- d) because your eye does not have enough depth of field

Concept questions on focusing

When you see someone *out-of-focus*

- a) There is no image anywhere
- b) There is an *in-focus* image on your fovea
- c) There is an *in-focus* image on your retina
- d) There is an image *in-focus* either in front or in back of your retina

Concept questions on focusing

In order to *focus* on close objects

- a) your eyelens needs to bulge
- b) your eyelens needs to flatten
- c) your cornea needs to bulge
- d) your cornea needs to flatten
- e) the distance (x_i) between your eyelens and retina needs to change

Structure of the eye: retina

- The retina is the sensor or film of your eye.
- Its layers do three things
 - Provide blood and nutrients (choroid)
 - Absorb light and convert to an electrical signal (photoreceptors)
 - Transfer the signal to the brain (nerve cells)

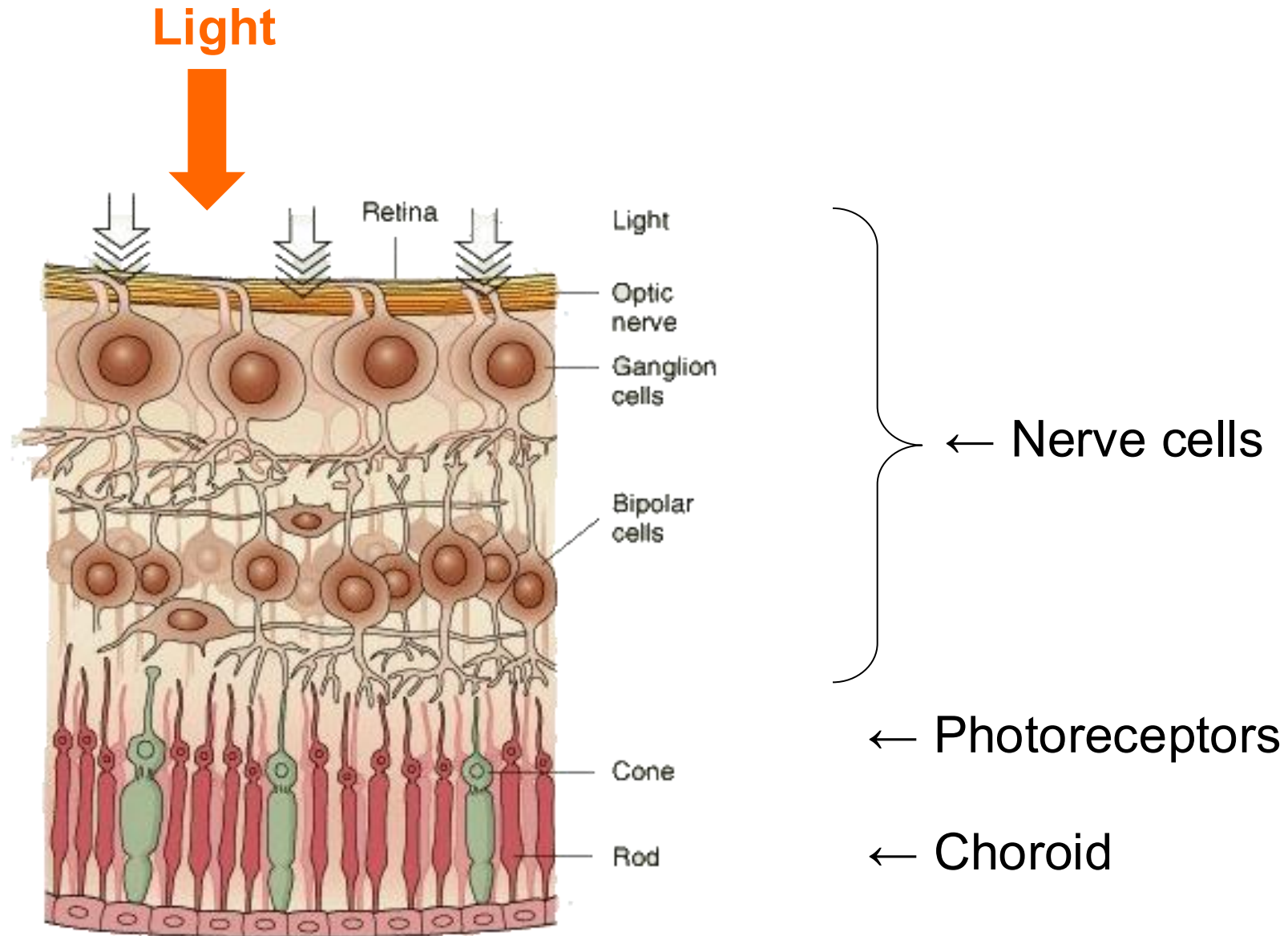
Light
↓

Plexiform layer (nerve cells)

Rods and Cones (photoreceptors)

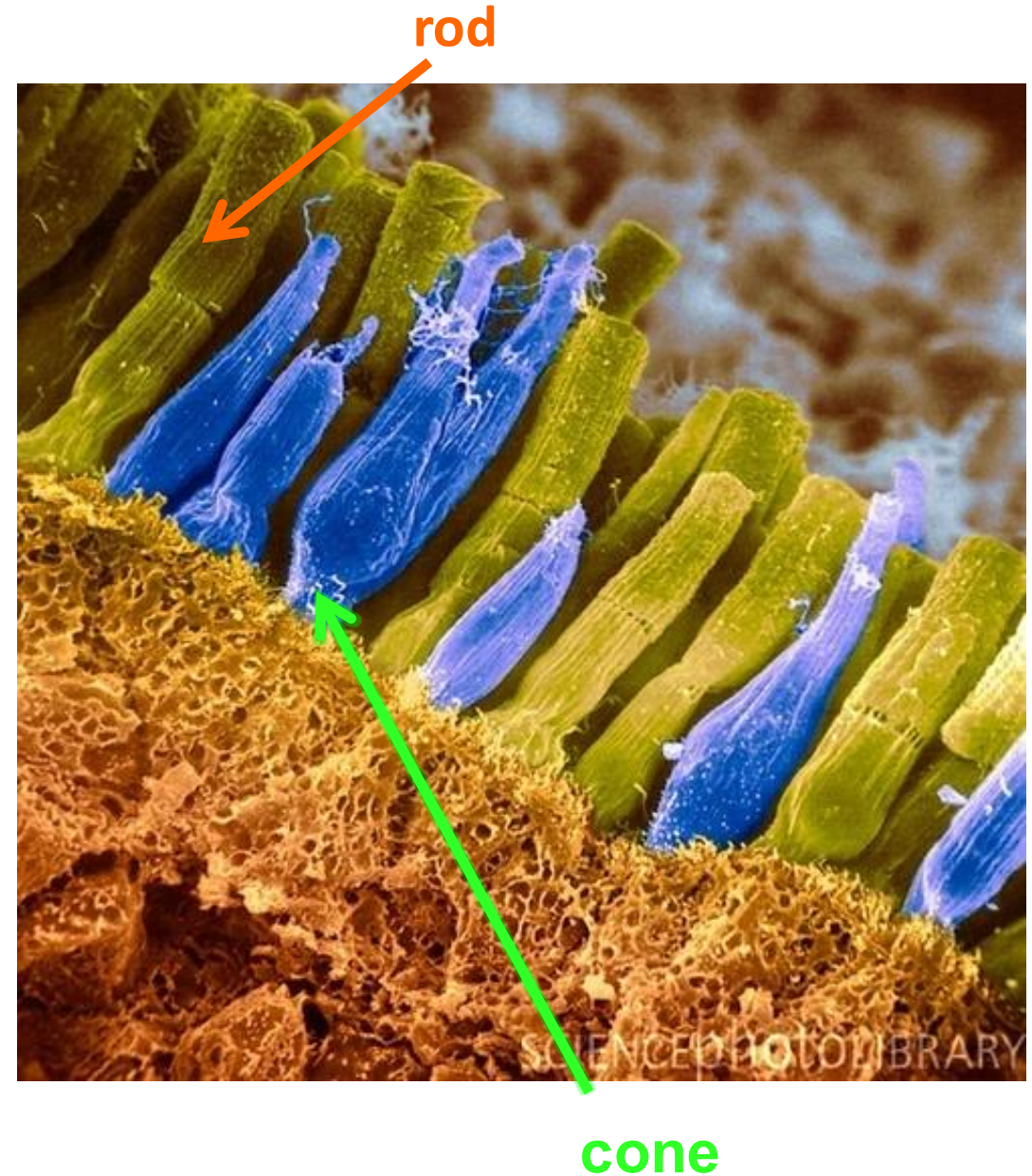
Choroid (blood vessels)

Structure of the retina



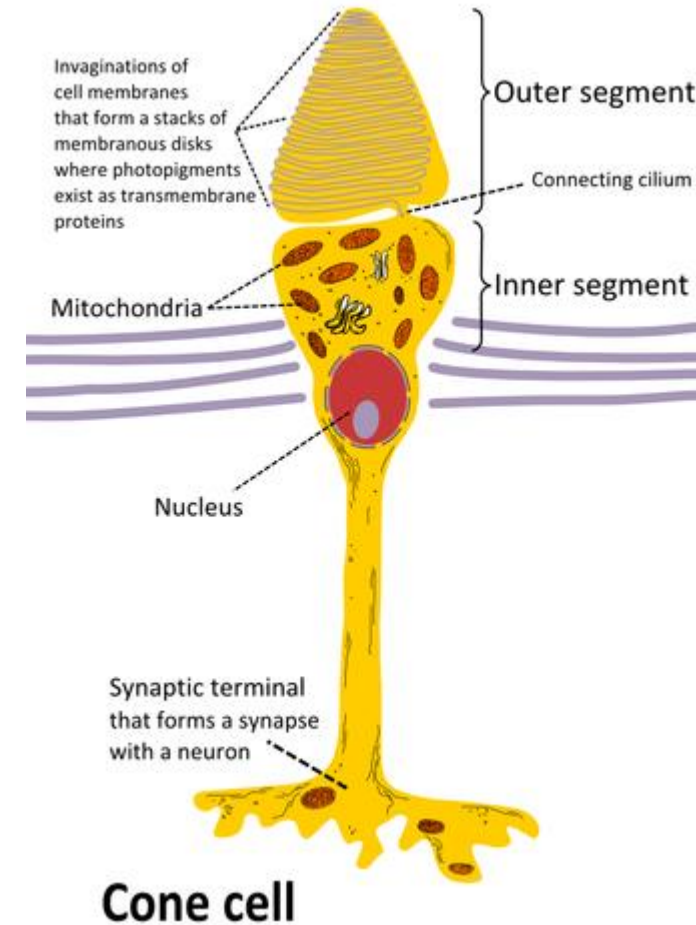
Photoreceptors: rods and cones

- Light is detected and converted to an electrical signal by the photoreceptors in the retina. There are two main kinds of receptors, rods and cones.
- This is a false color image, rods and cones are not actually different colors.



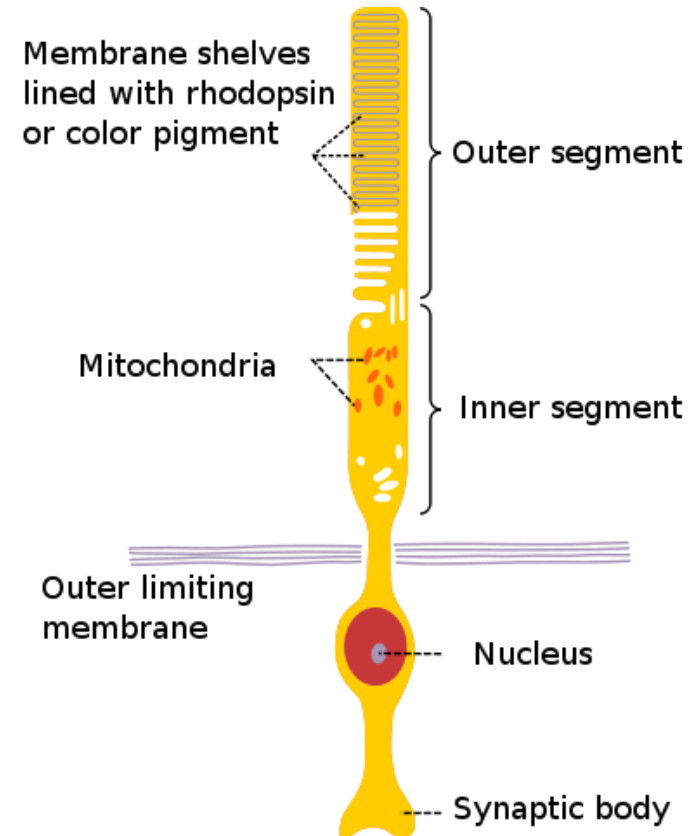
Photoreceptors: cones

- Cones are responsible for our fine detailed and color vision
- Cones are clustered near the center of your retina, called the *fovea*
- There are 5 million cones in the average retina



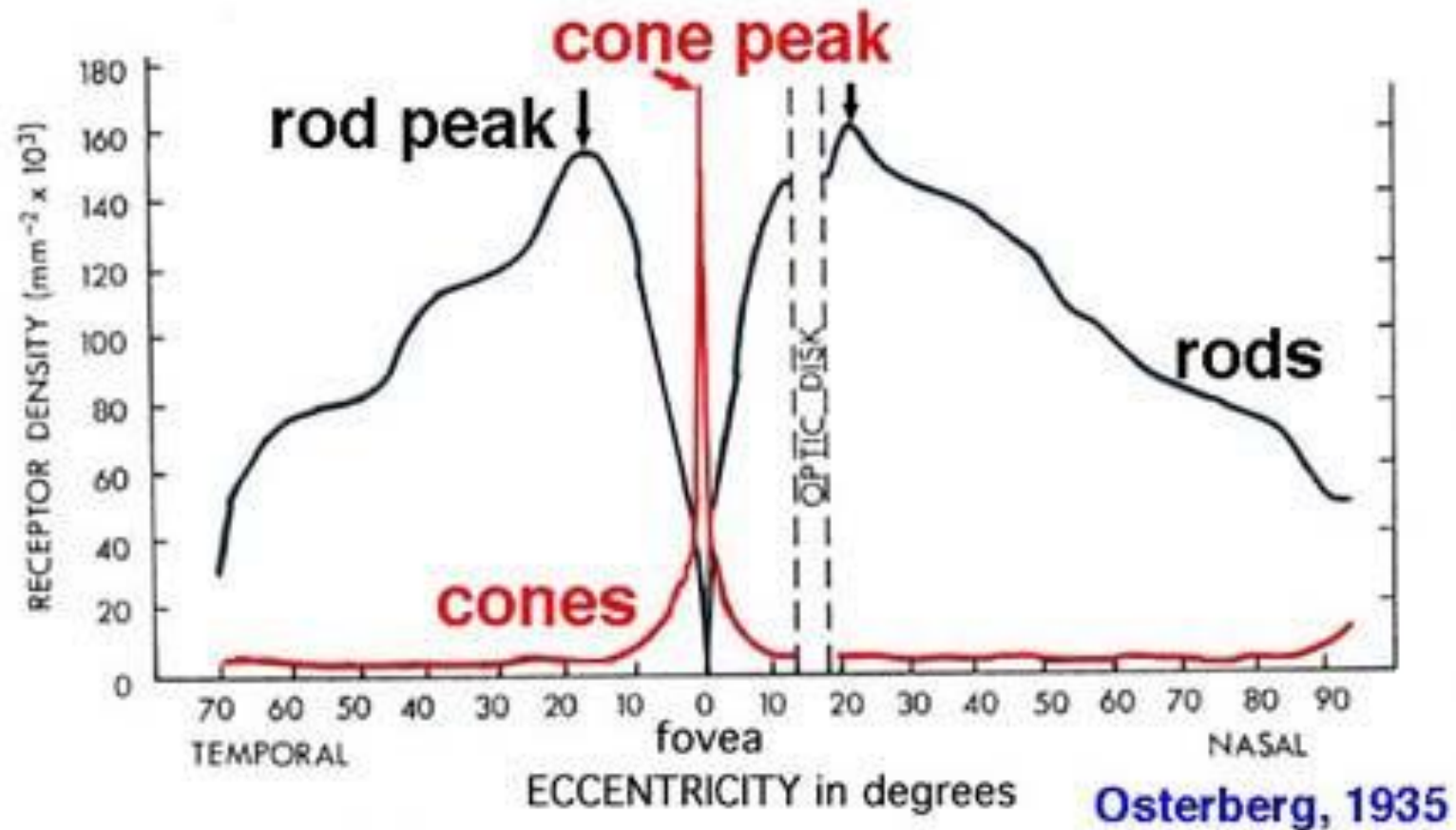
Photoreceptors: rods

- Rods are responsible for low light and peripheral vision
- They are present everywhere in the retina except the fovea
- There are 125 million rods in the average retina



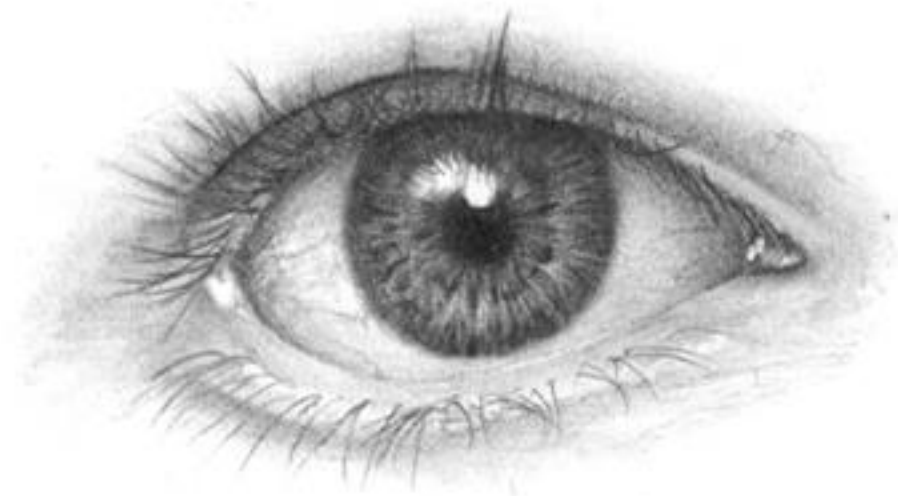
Rods and cones

- Because of their different functions, rods and cones are present in varying densities in the retina. The blind spot is due to the connection of the optic nerve



Human vision: our benchmark

- Is a human eye just that much better than a CCD camera?

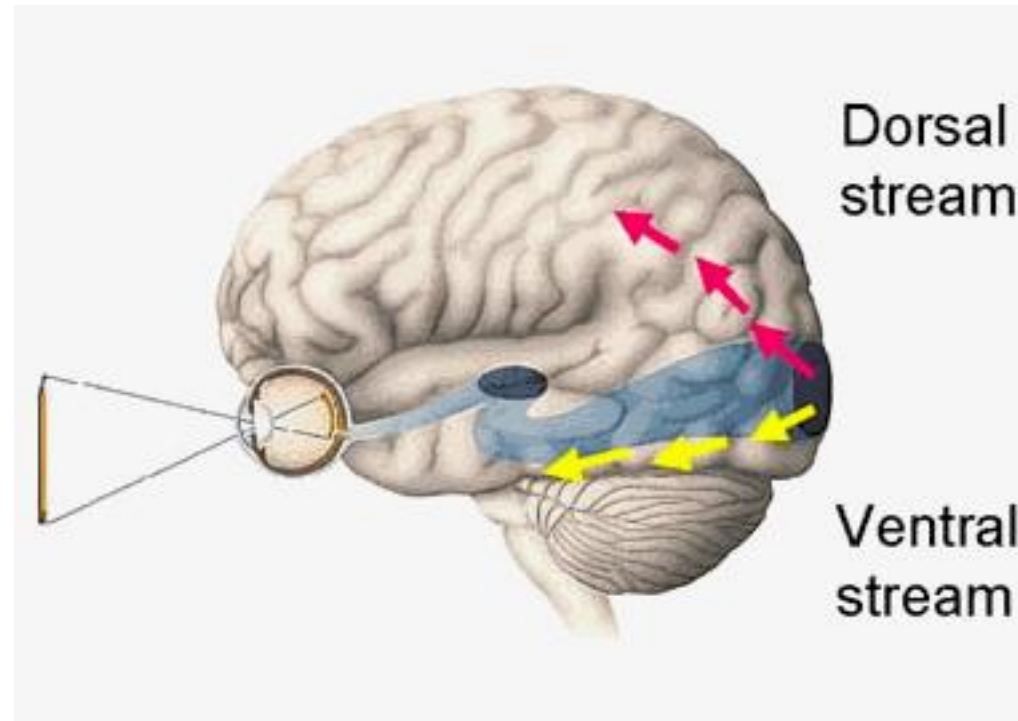


Definitely NOT!!!

Fewer pixels, less depth-of-field, etc.

Human vision: our benchmark

- Is a human eye just that much better than a CCD camera?

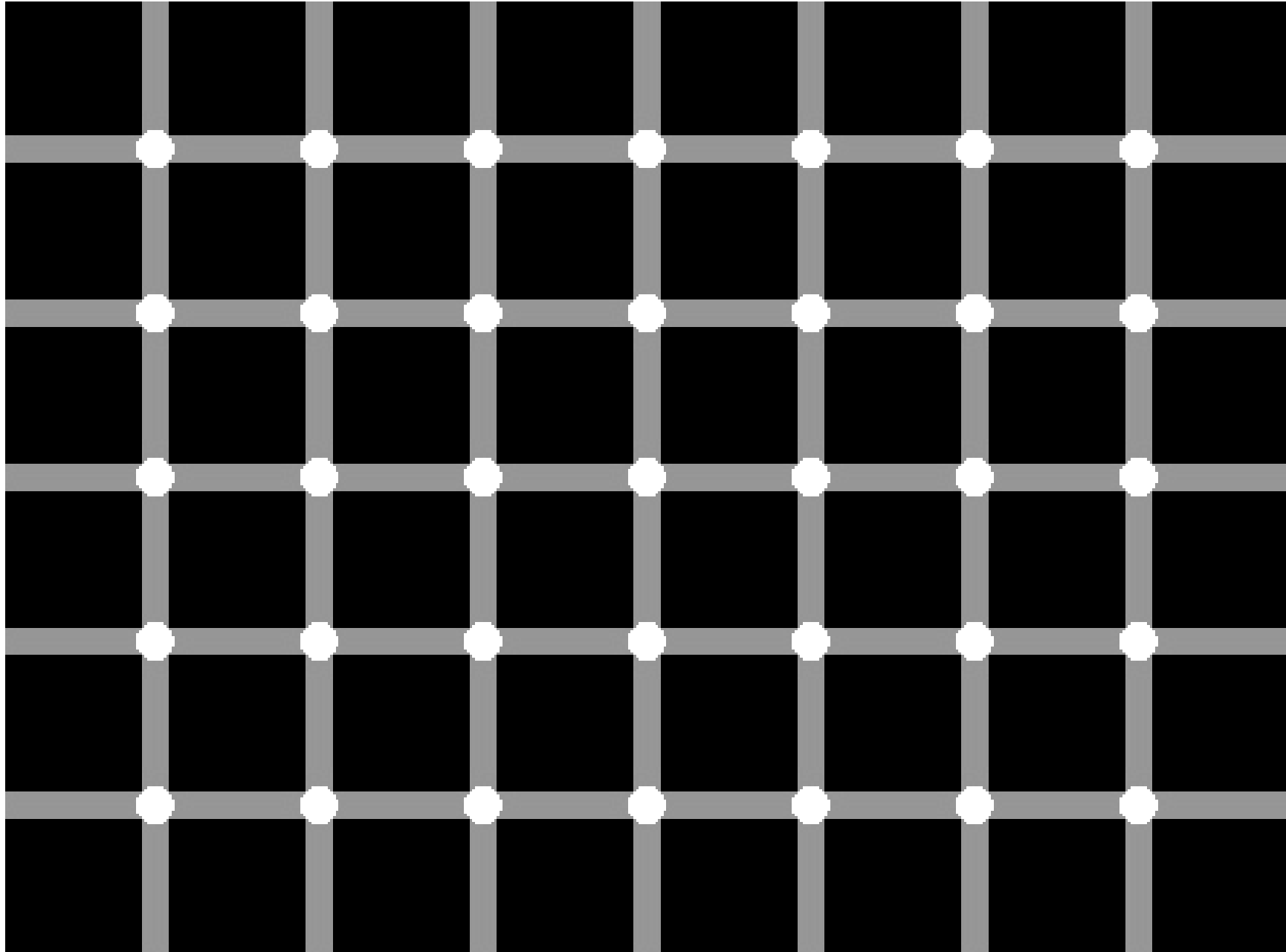


- Then how do humans **see** so well?

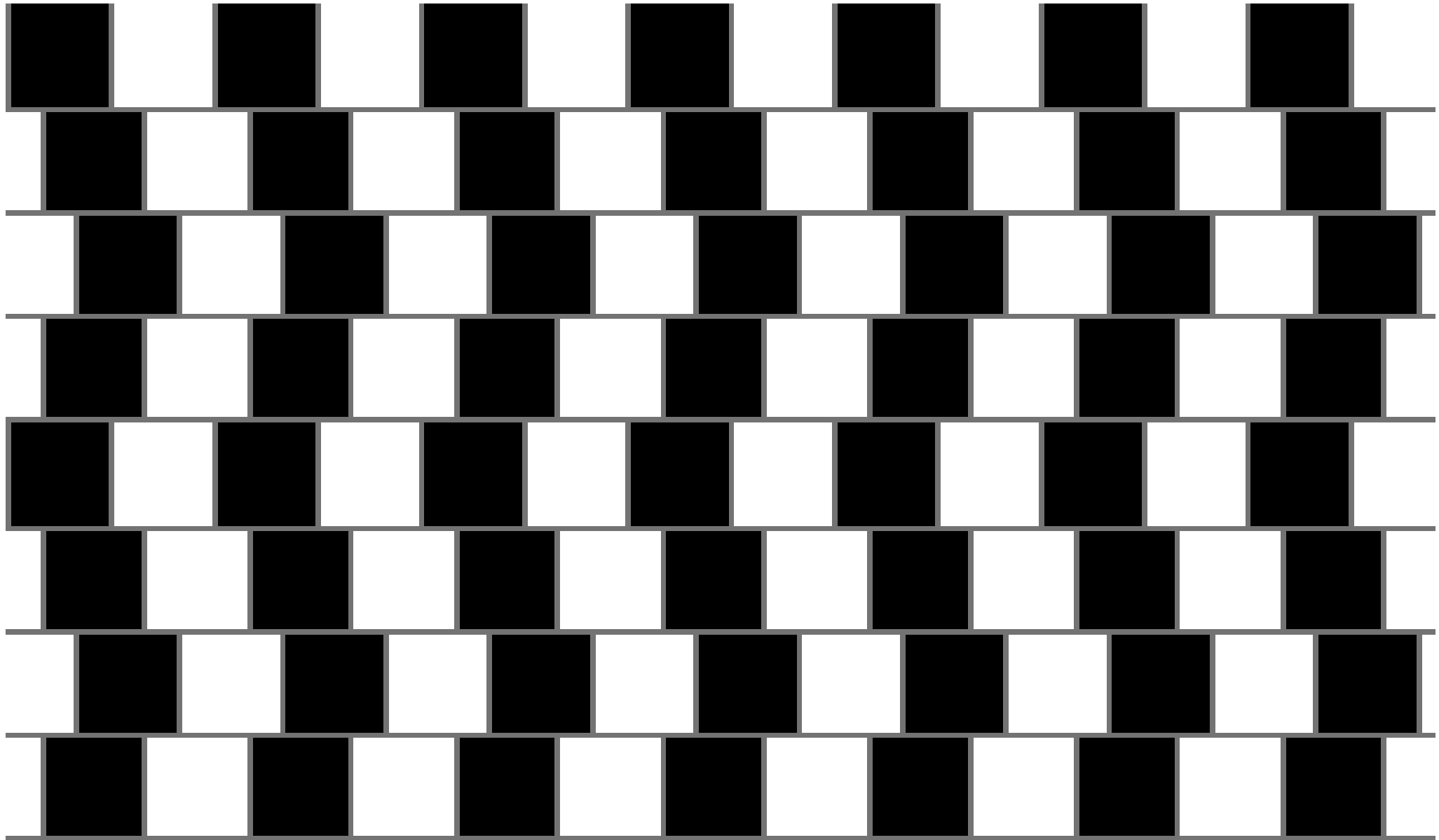
Human vision: our benchmark



Illusions: what do they tell us?



Illusions: what do they tell us?



Why is computer vision hard?



Why is computer vision hard?



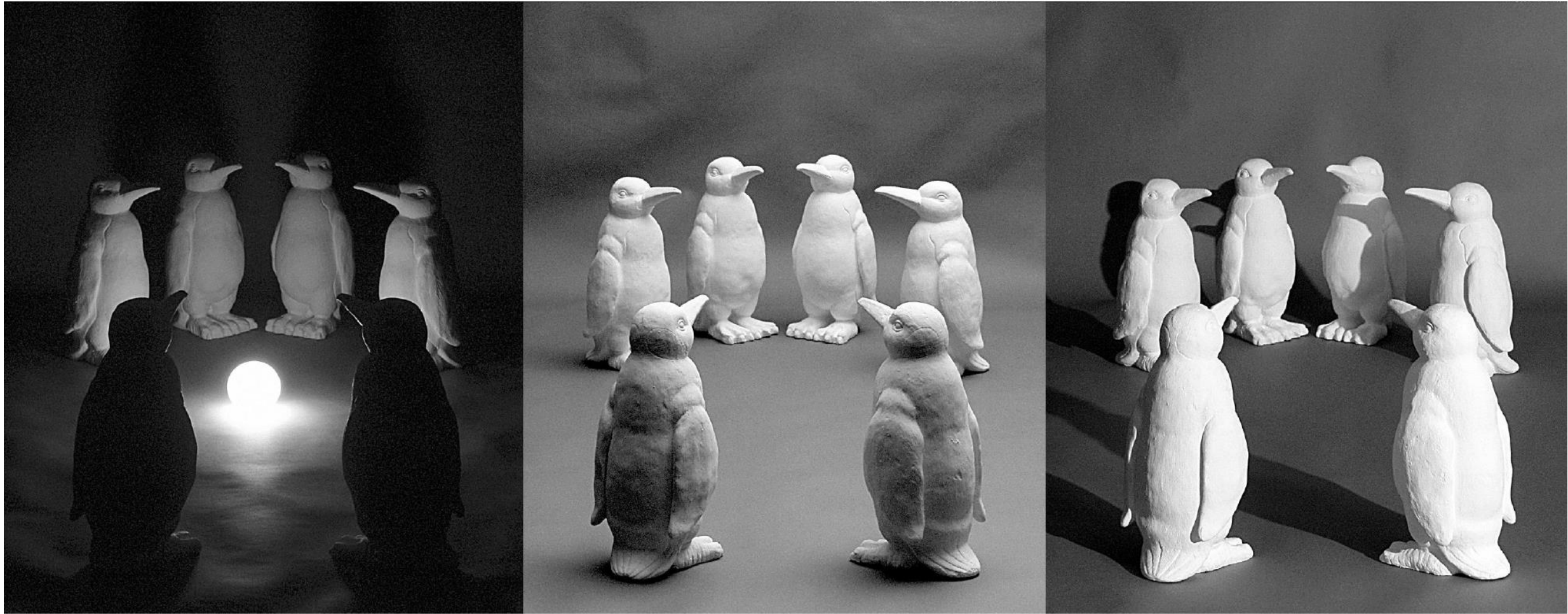
Why is computer vision hard?

- “Context” counts for as much as appearance.
- Huge amounts of prior knowledge (learned and innate).
- AI complete (hard to solve a cleanly defined “simple” problem without invoking unrealistic assumptions).
- Lack of a clear metric for success (indeed, we are often completely wrong as you’ve seen).
- The diversity of the natural world.

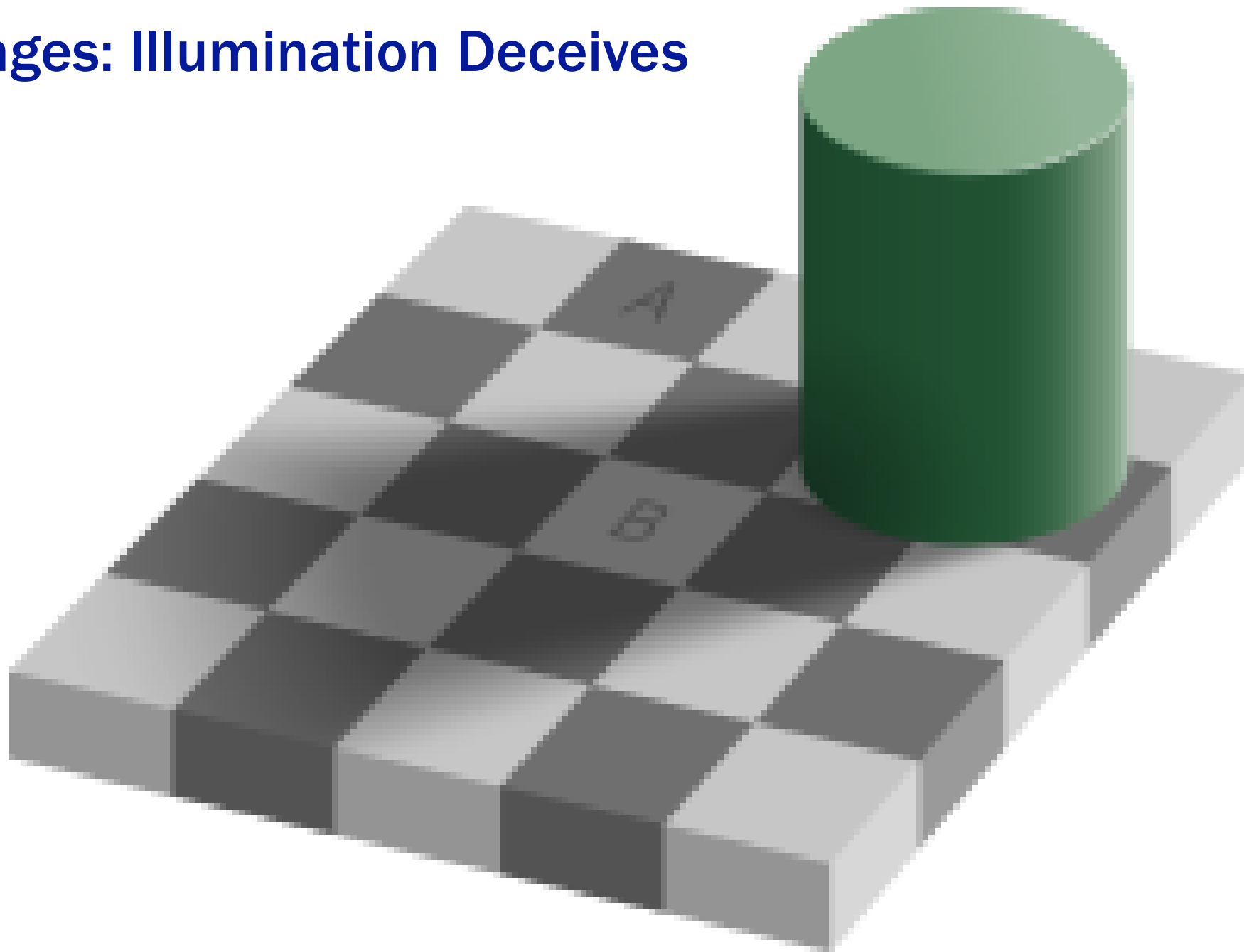
Challenges: viewpoint variation



Challenges: illumination variation



Challenges: Illumination Deceives





Challenges: scale variation

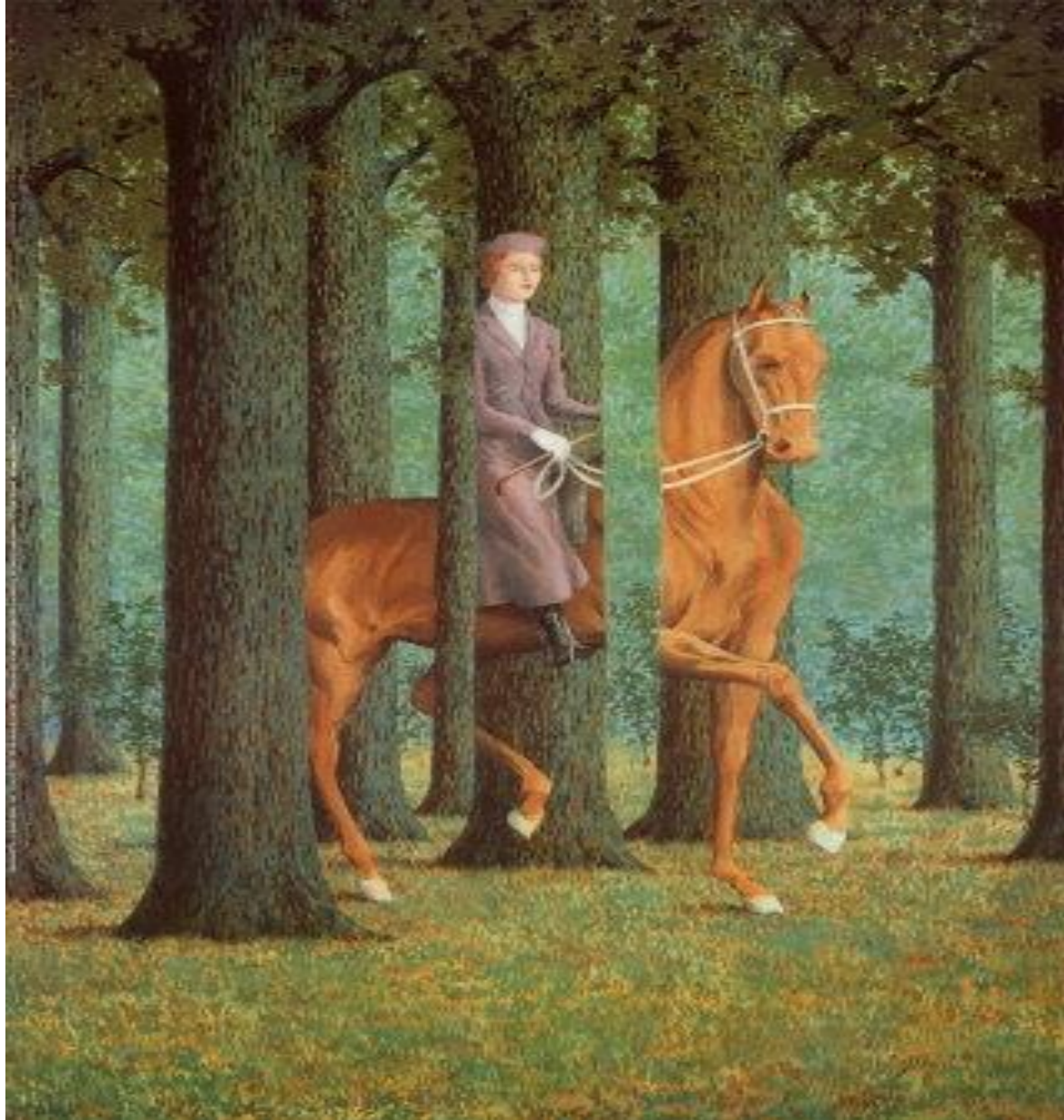


slide

Challenges: deformation and structure variation



Challenges: occlusion



Challenges: background clutter or context variation



Challenges: intra-class variation



Challenges: class variation



~10,000 to 30,000

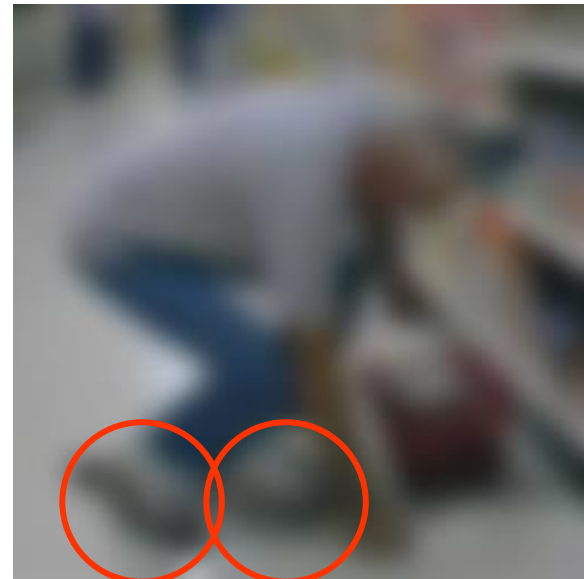
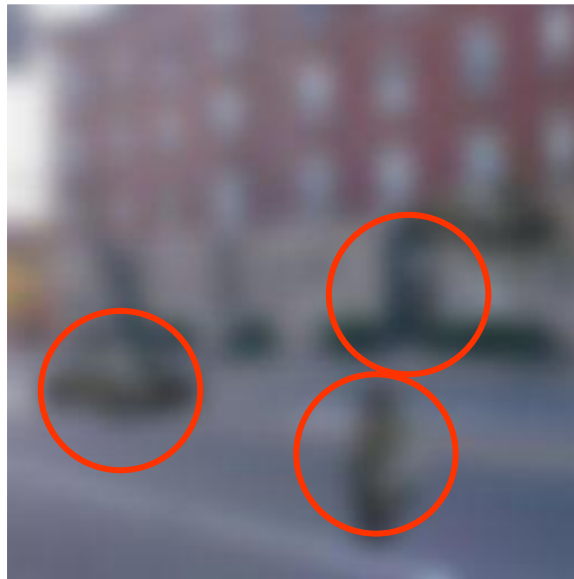
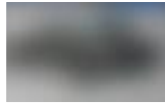
Challenges: ambiguity



Challenges: ambiguity

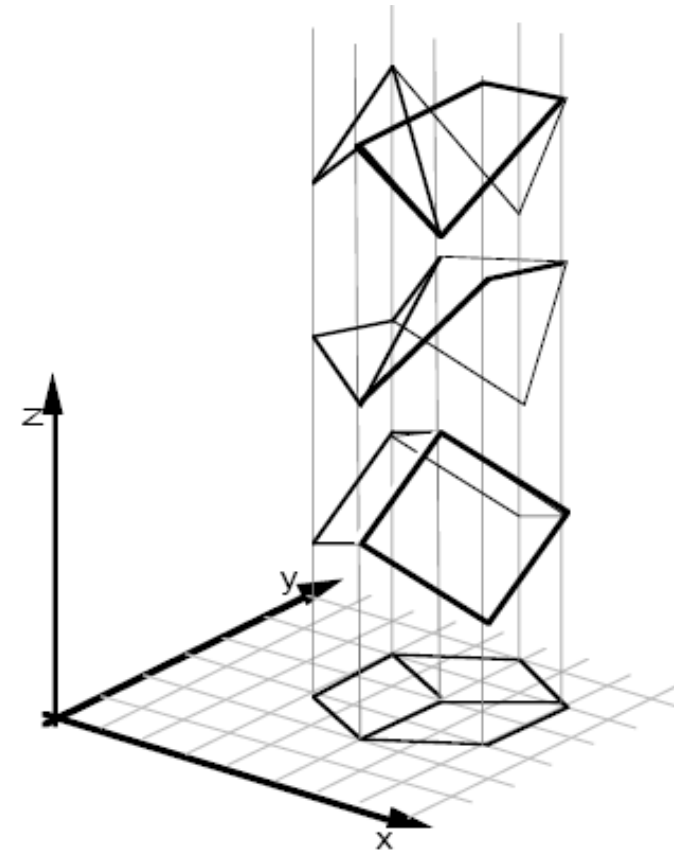


Challenges: ambiguity



Challenges: ambiguity

- Many different 3D scenes could have given rise to a particular 2D picture



[Sinha and Adelson 1993]

Scenes are unique



But not all scenes are so original



But not all scenes are so original



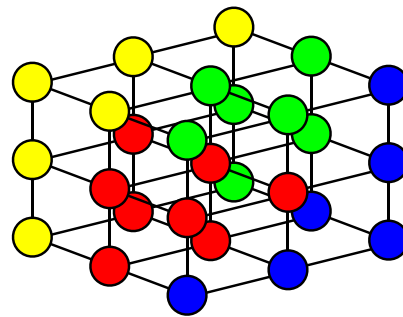
Computer Vision: A Search for Invariance

- One major theme of the semester (and computer vision) is the search for invariance.
 - **Invariant/Invariance**: a function, quantity, or property that remains unchanged when a specified transformation is applied.
 - Note usage, e.g.,
 - Algorithm X is scale invariant.
 - Algorithm X demonstrates scale invariance.
- Examples we will work through
 - Shift, scale, rotation, affine, photometric, projective, deformation, structure, class, category, rate

Admittedly, this notion of vision as a search for invariance has decreased in "popularity" in recent years due to the great successes of transformer-based big vision.

Major Themes

- Computer Vision as a search for visual invariants.
- Computer Vision as optimization.
 - Nearly all (well-done) examples of computer vision methods/systems can be posed a minimization of some objective.
- Representation of visual content is an important issue.
 - Images as functions $I: \mathbb{R}^2 \rightarrow \mathbb{R}$
 - Continuous vs. discrete in domain and range
 - Images as points (or coefficients)
 - Standard bases: Cartesian basis, Fourier basis
 - Learned bases.
 - Images as graphs
 - Pixels induce nodes.
 - Connectivity?



Core Abstracted Problems in Computer Vision

Module 2

- **Reduction**

$$X^* = \arg \min_X f(I)$$

- Transforming visual content into an alternative representation, often with the loss of data, but with the retention of important characteristics (think invariance).
- E.g., segmentation, feature extraction, local edge detection

- **Matching**

$$\hat{X} = \min_{\{X'\}} d(X, X')$$

- Finding local correspondences between comparable representations of visual content, with the potential of incorporating global regularizing constraints.
- E.g., stereo correspondence, image retrieval

- **Fitting**

$$\theta^* = \arg \min_{\theta} g(X, X'; \theta)$$

- Estimating the parameters of a model from a set of representations of visual content and potentially correspondences.
- E.g., image registration, 3D reconstruction

Computer Vision Problems and Applications

- Image denoising
- Image search / retrieval
- Reconstruction and Structure from Motion
- Recognition of objects, actions, scenes
- Tracking of objects and articulation
- Segmentation
- Detection
- Optical flow

**Will cover
in a good
bit of
detail this
Friday.**



Course Information

- Syllabus **CANVAS site will be published tomorrow, with syllabus.**
 - Contains all necessary information for course structure and grading. (below ratios may vary, still deciding)
 - 30% Homework (3 or 4; likely drop 1 if 4)
 - 10% PACES (Weekly paper reading output)
 - 15% Project
 - 15% Quizzes (~10; two dropped)
 - 20% Exam
 - 10% Participation (Read/View/Annotate) (~10)
- Canvas
 - Canvas is the primary source for course information.
- Perusall (through Canvas)
 - Used for reading/annotating.
 - Used for watching and annotating video content
 - Used for ALL Discussion

Debating
Project vs
Second Exam

AI Makes
Things Harder
for Courses

Debating
Discussion in
Perusall or
Piazza

Maybe
1.5 GSIs

**Note we have 1 GSI
who will host office
hours in addition to
the instructor.**

Course Information

- Content Types
 - **Prerecorded Video Lectures**
(materials for quiz and exam sourced from here and in-class)
 - **Lecture Reading** – textbook and notes materials supporting the recorded prerecorded and live lecture
 - **Live Lecture** – parts of prerecorded video lecture iterated, questions answered, etc.
 - **Quiz** – Held in class, participation and learning focused, not assessment focused. Covers material from prerecorded video lecture and reading BEFORE discussed in class!
 - **Paper Reading** – A primary, maybe even seminal, paper. One per week, read and discussed on Perusall, group discussion in class.

WHAT?!

Yes, this course is partially flipped!

Course Information

- Typical week
 - Friday before week:
 - Readings and prerecorded lectures posted to Canvas and Perusall.
 - PACES paper of the week posted.
 - Saturday—Wednesday: watching and reading and discussing
 - Wednesday 2PM time: Comments captured from Perusall for participation credit (ie due)
 - Wednesday 4:30PM time: Live class begins
 - 4:30-5:30 Live Lecture
 - 5:30-6:00 Outside Lecturer (Industry, Research, etc.)
 - 6:00-6:30 PACES Active Work
 - Friday 10:30AM time:
 - 10:30-10:50 In-class quiz and discussion
 - Selected Perusall and Readings questions presentation
 - Time-permitting: modern vision demo and AMA

Roughly 1.5 hours
prerecorded lecture
weekly

LLMs and Generative AI

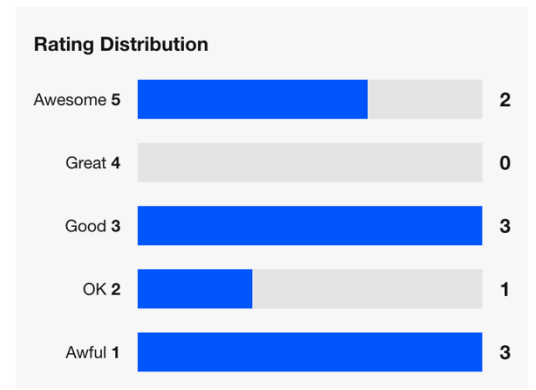
- From Syllabus:

Use the tools at your disposal. We are here to seek the truth. To understand. To communicate. But, you are expected to do your own work. Plagiarism is never acceptable. Copying the output of an automated system is the same as copying the output of a classmate.

Admittedly, this is not something easily translated into practice beyond (don't really use AIs, which is not acceptable). So, I will refine this policy...

Rate Your Professor

I've taught no fewer than 5,000 students in the last two decades. My teaching style differs from many Engineering faculty. I choose loose vectors and self-directed learning whereas many choose exacted blueprinted learning.



Some like this, some do not

QUALITY **1.0** **EECS442** Jan 16th, 2019

For Credit: **Yes** Attendance: **Mandatory** Would Take Again: **No** Textbook: **No**

Never replied questions on Piazza and he has no idea what homework and the exam is.

DIFFICULTY **5.0**

SKIP CLASS? YOU WON'T PASS. LOTS OF HOMEWORK TEST HEAVY

Helpful 0 0

QUALITY **3.0** **EECS598** Nov 14th, 2015

For Credit: **Yes** Attendance: **Not Mandatory** Textbook: **No**

Prof. Corso is a very clever and hard-working. I received emails from him on 5am for many times. He is also very passionate to what he taught. Because EECS598 is a new course created by Corso(2015Fall is second semester since this course was opened) and Computer Vision is a relatively new topic, this course still need to be improved in many places.

DIFFICULTY **4.0**

Helpful 1 1

9 of those students have given some colorful commentary in public.

Rate Your Professor

I've taught no fewer than 5,000 students in the last two decades. My teaching style differs from many Engineering faculty. I choose loose vectors and self-directed learning whereas many choose exacted blueprinted learning.

QUALITY





5.0

EECS442 Jan 9th, 2018

For Credit: **Yes** Attendance: **Mandatory** Would Take Again: **Yes** Grade: **A-** Textbook: **No**

Prof. Corso is a very inspiring teacher. I have never learnt computer version before. It's hard to catch up at the beginning but Prof makes it interesting and easier to learn. Homework is a little heavy but it really worth it. I would love to take his course again.

LOTS OF HOMEWORK **INSPIRATIONAL**

Helpful  2  11  

QUALITY

3.0

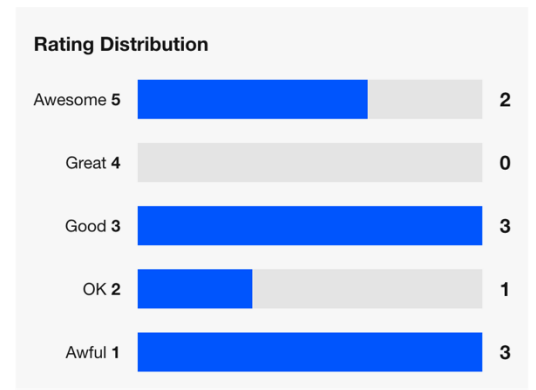
EECS442 Nov 25th, 2018

QUALITY

1.0

EECS442 Nov 14th, 2018

For Credit: **Yes** Attendance: **Mandatory** Would Take Again: **No** Textbook: **No**



Some like this, some do not

Taking this (or any course) is a bidirectional agreement in which we embark on a learning journey together. I am the guide on that journey, not the oracle.

Coming up in Week 2

Images as Functions

- Prerecorded video lecture and readings posted Friday
- Mathematical preliminaries posted (optional)
 - Linear least squares
 - Lagrange Multipliers
 - Matrices and Linear Algebra
 - Eigendecomposition and SVD

- Homework 0 (2pts) – quick survey about the course;
 - Assigned Friday and due Monday

- Homework 1 will be released next week.