

the quest for robotic vision

Peter Corke



Queensland University
of Technology



in the beginning



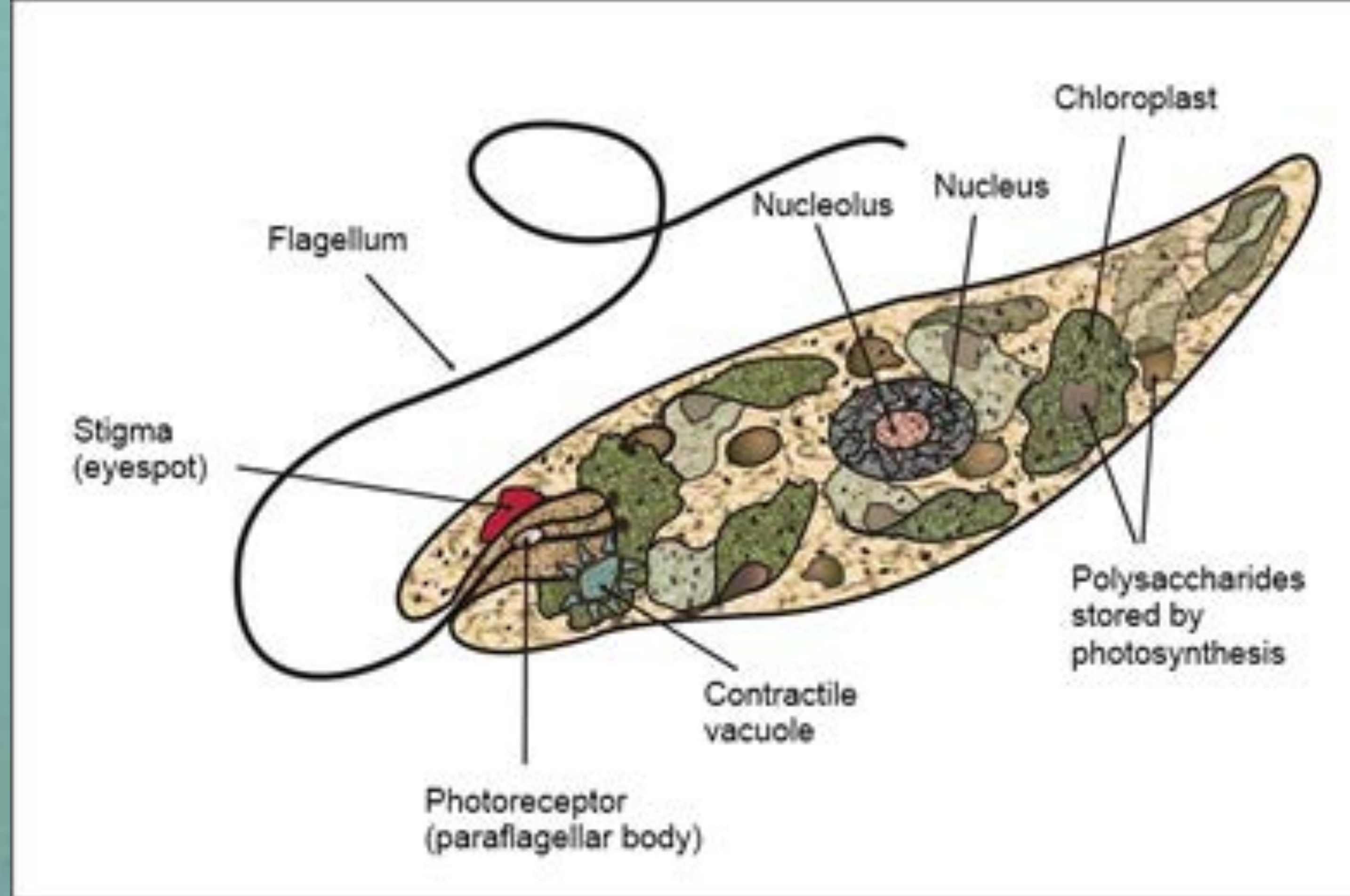
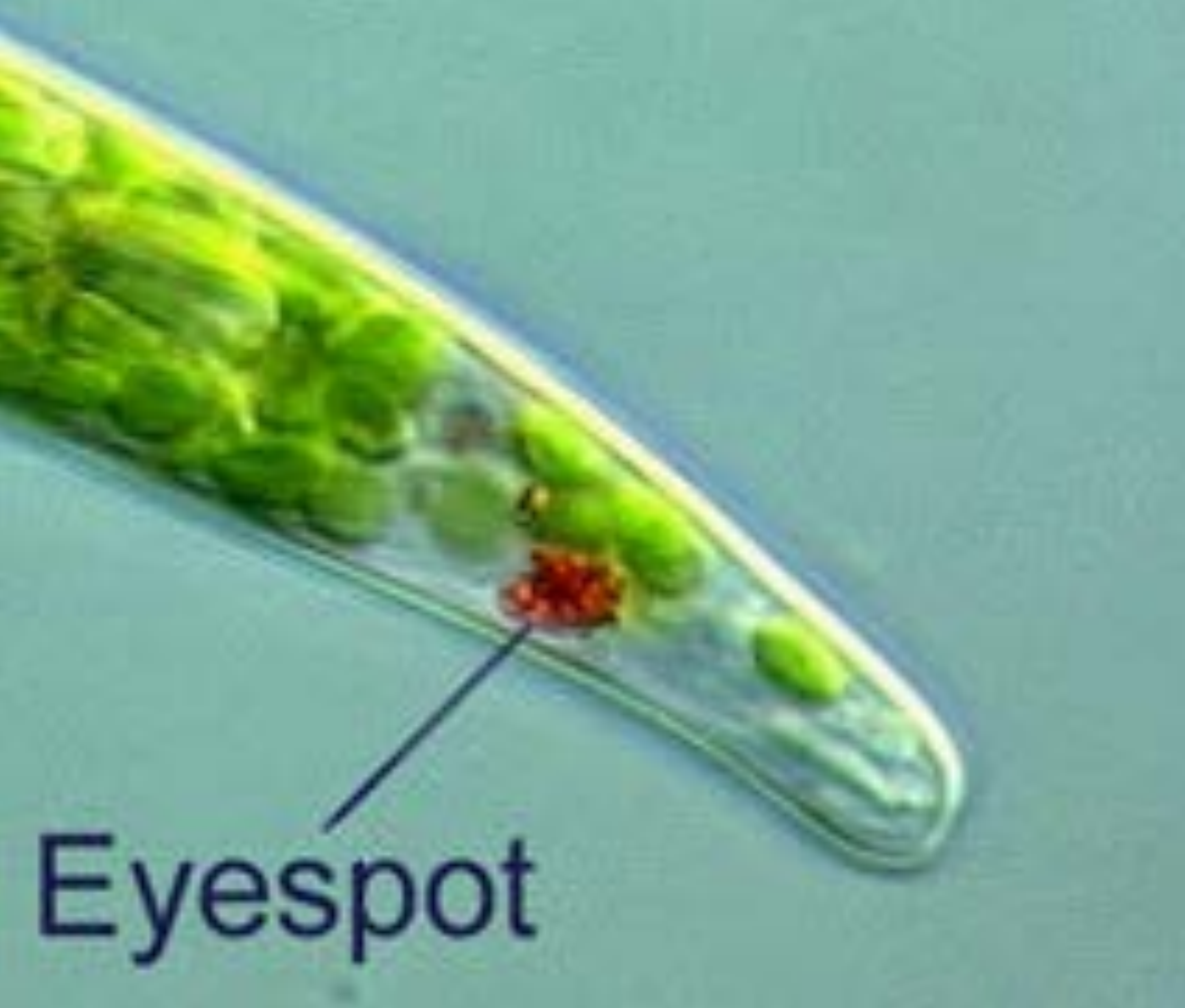
t = -13.8B years



t = -4B years



$t = -3.5$ years

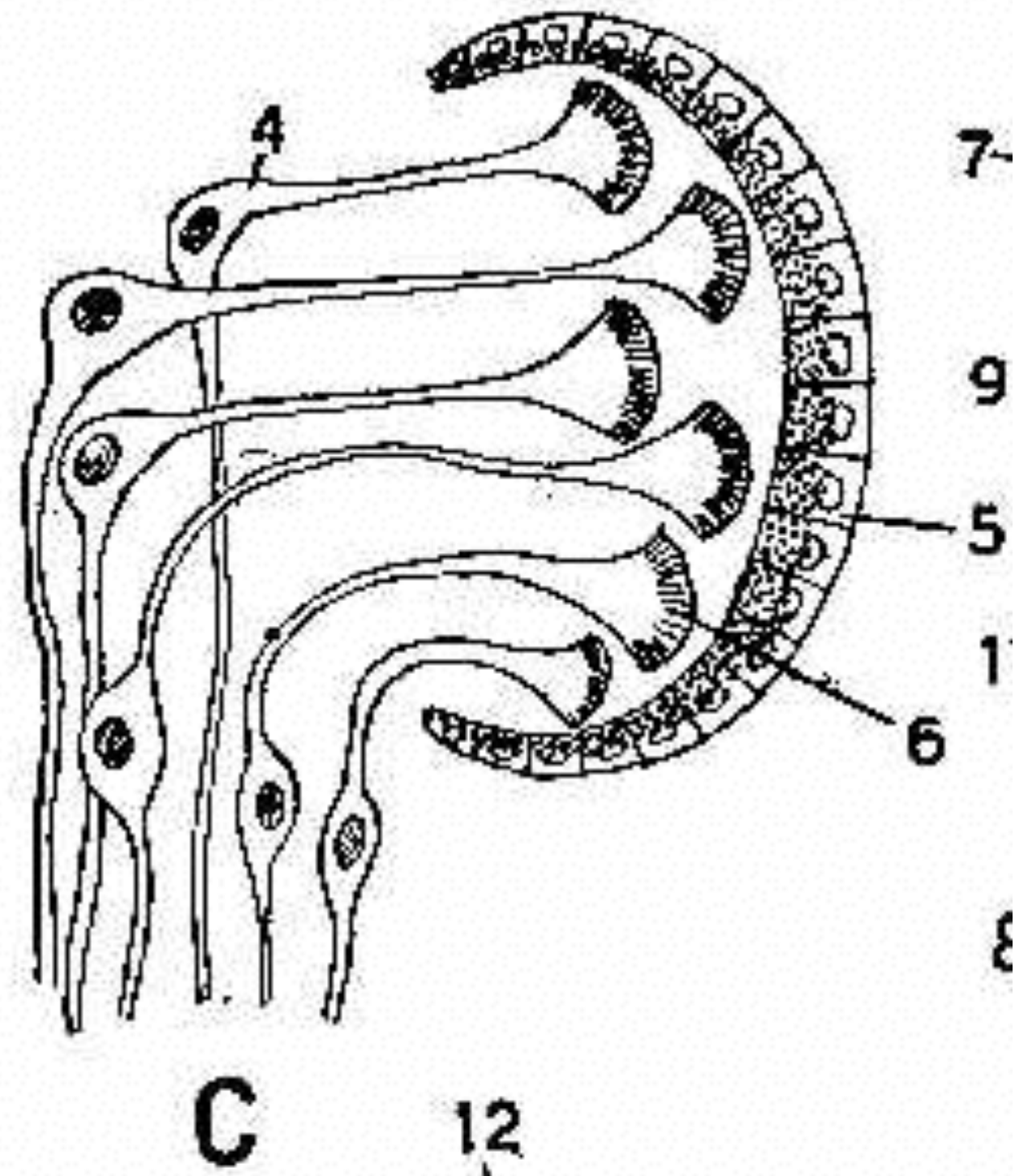


evolution invents a
light sensor

$t = -600M$ years



Planarium flat worm




t = -550M years

**several eyespots
make an eye**



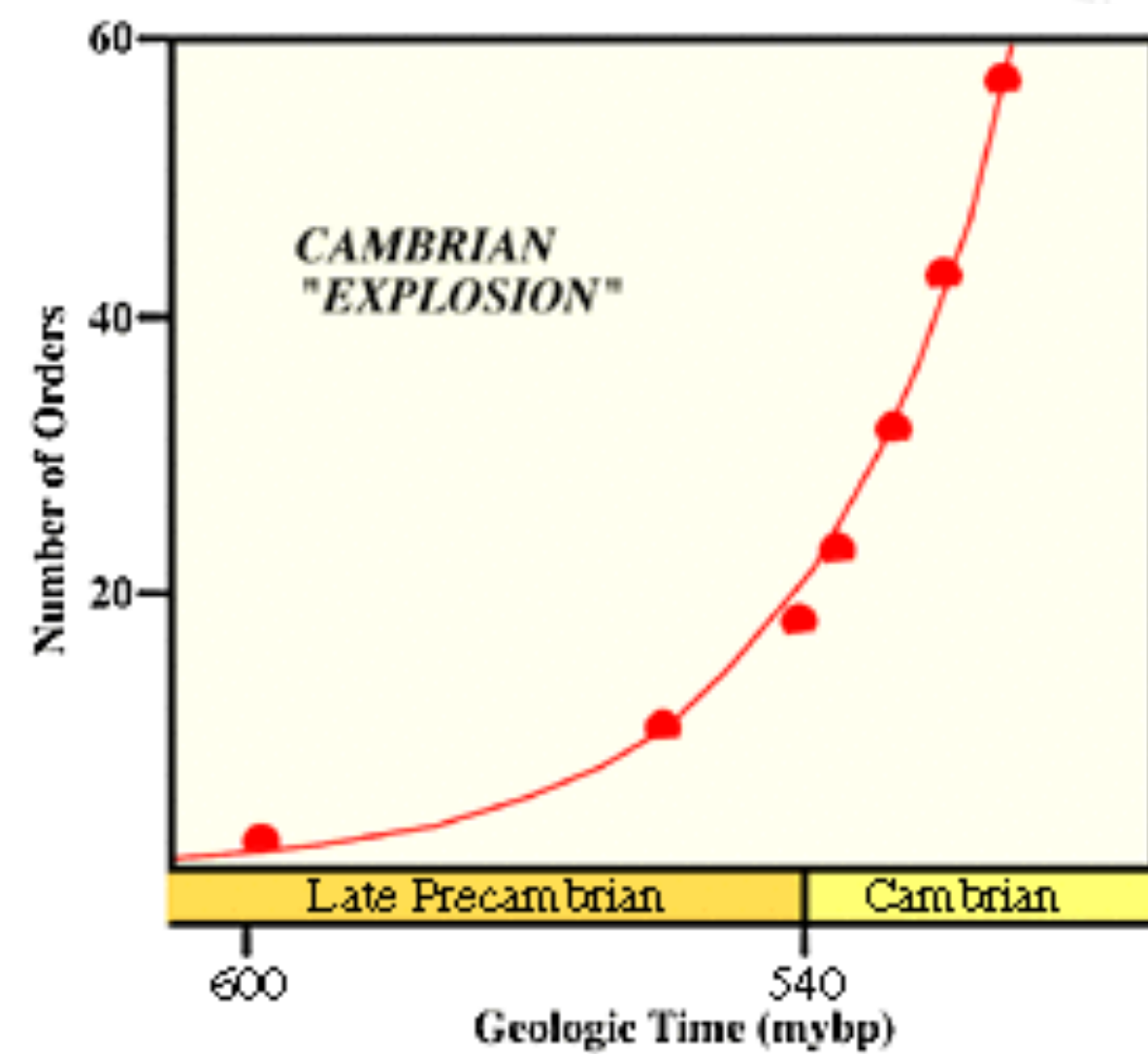
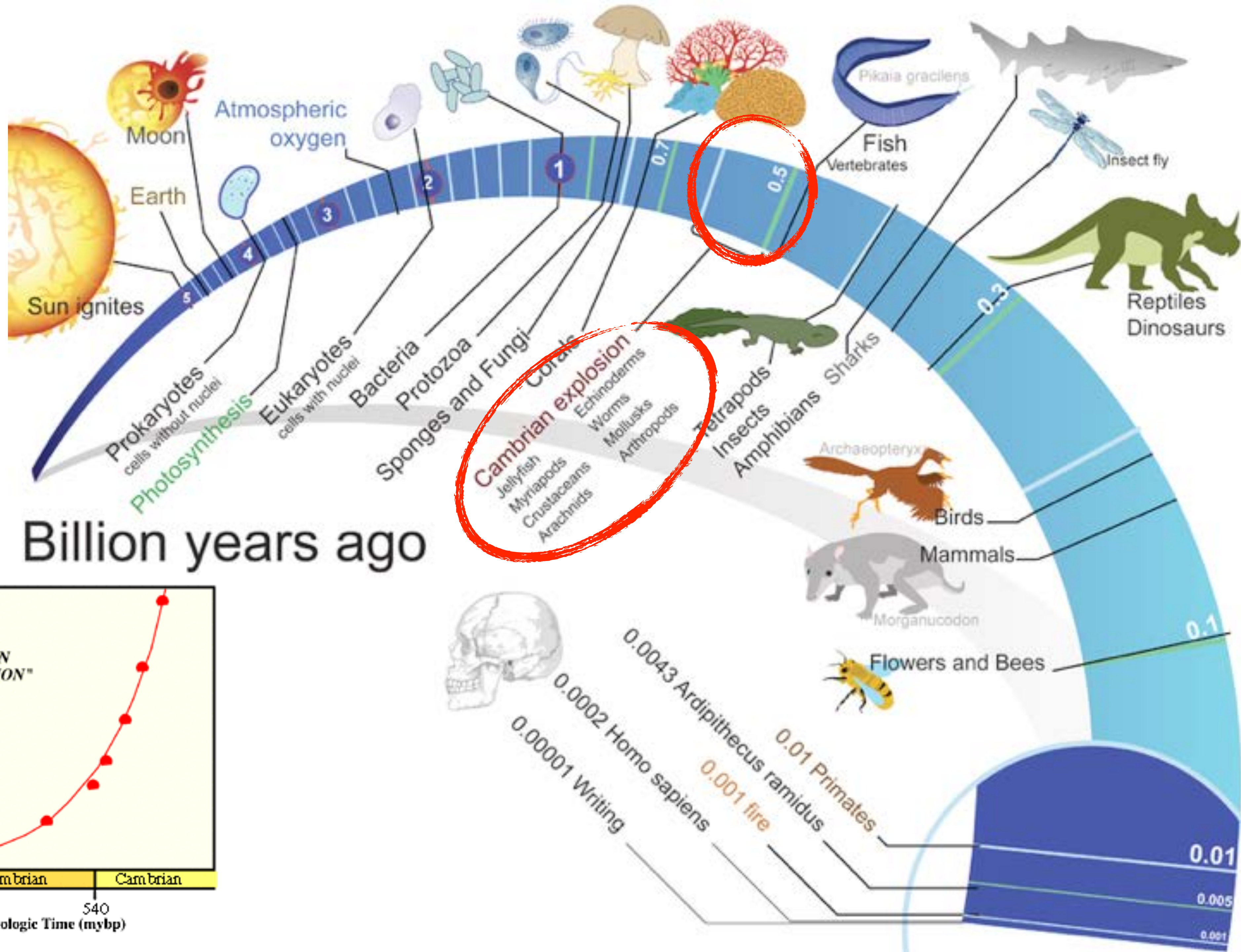
t = -521M years

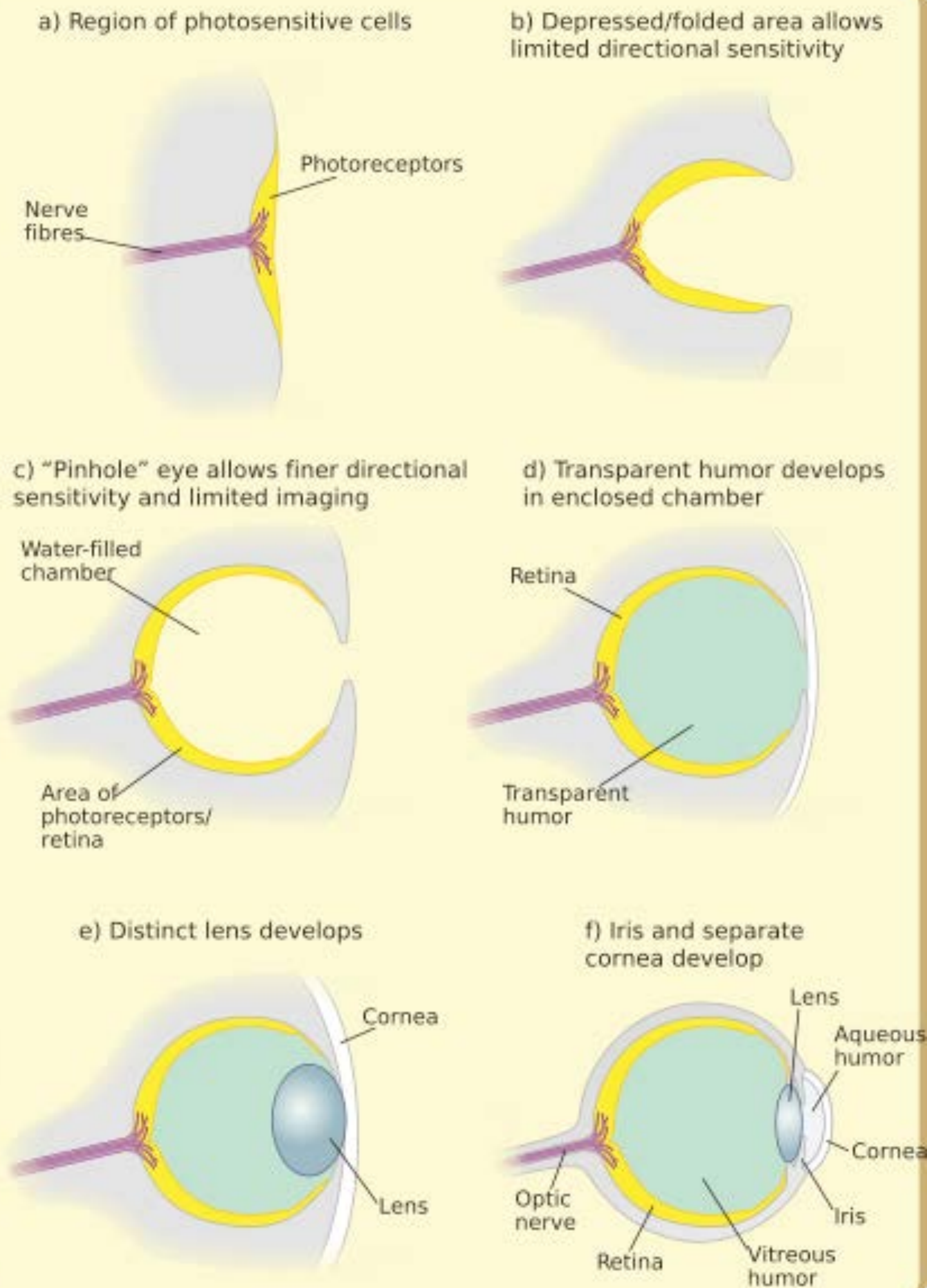
**trilobite: successful marine
animal for 270M years**

A photograph of two lampreys in an aquarium. One lamprey is in the foreground, swimming towards the right, showing its head with a single eye and a row of seven barbels. The second lamprey is in the background, swimming towards the left. The tank floor is covered with small, colorful gravel. The background is a dark, solid color.

lamprey eye is very
similar to our own eye

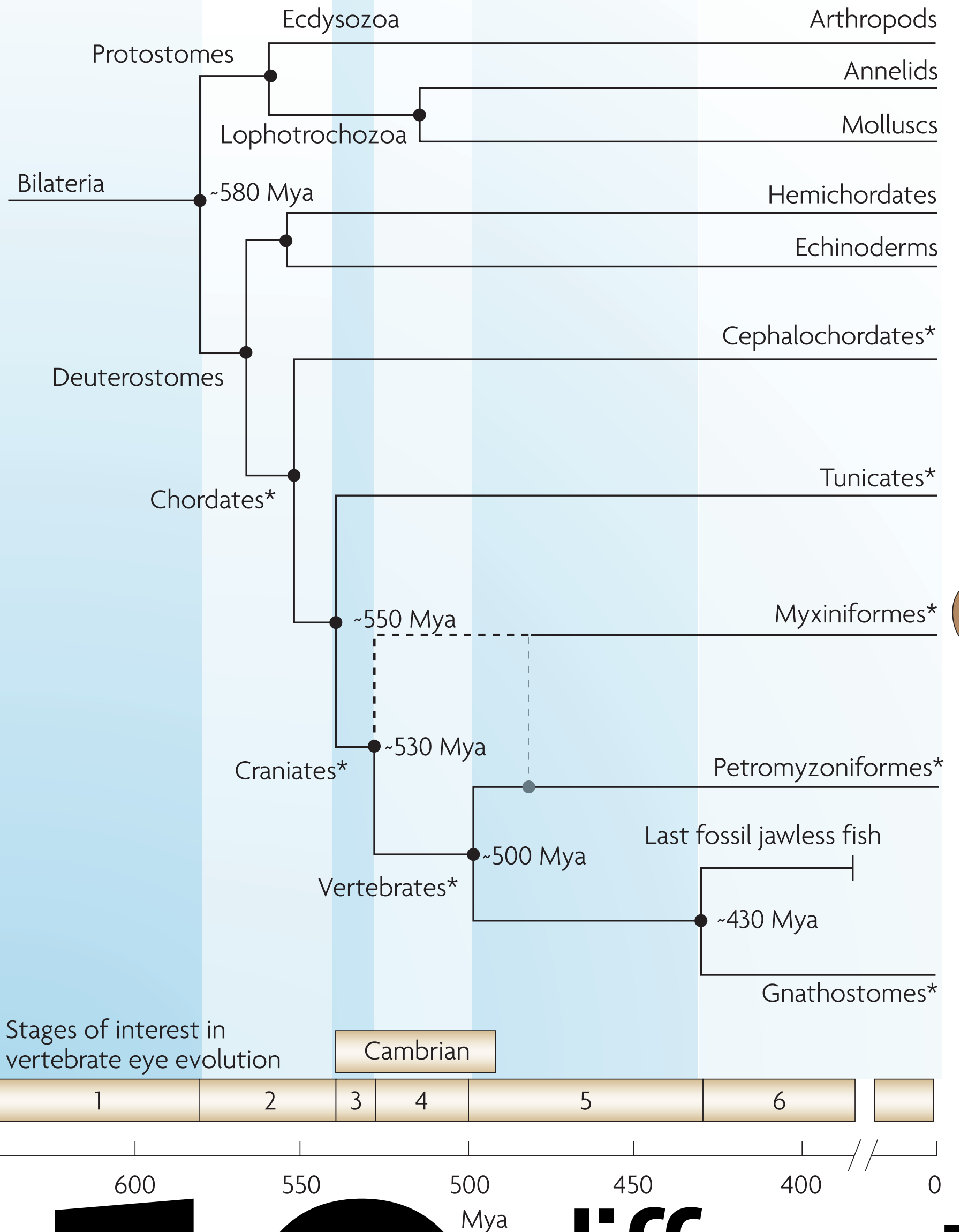
t = -500M years





To suppose that the eye, with all its inimitable contrivances... could have been formed by natural selection, seems, I freely confess, absurd in the highest possible degree... Yet reason tells me, that if numerous gradations from a perfect and complex eye to one very imperfect and simple, each grade being useful to its possessor...

- Charles Darwin (1809–1882)



10 different eye designs

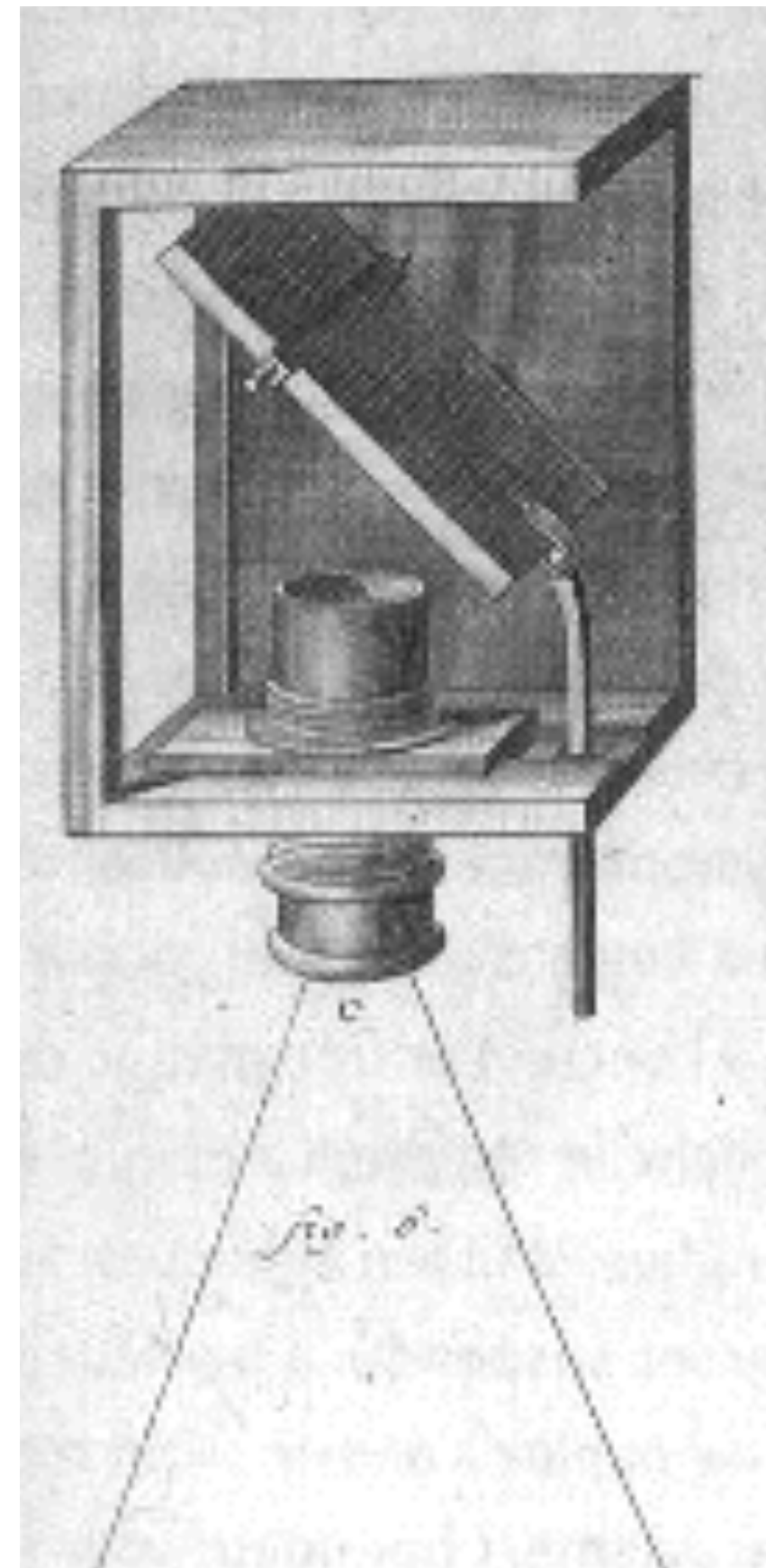
96% of animal species have eyes

retinal projection



$$(X, Y, Z) \mapsto \left(\frac{fX}{Z}, \frac{fY}{Z} \right)$$

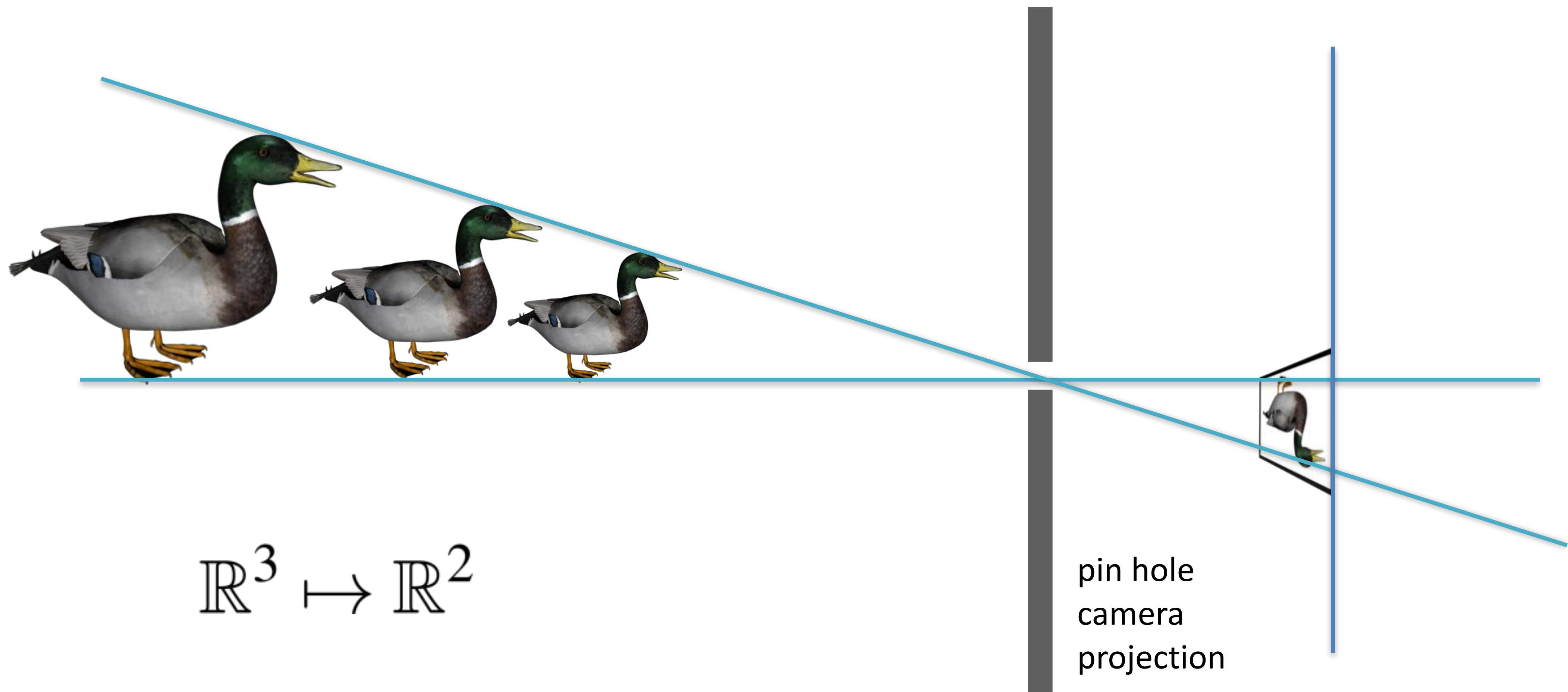
$$\mathbb{R}^3 \mapsto \mathbb{R}^2$$



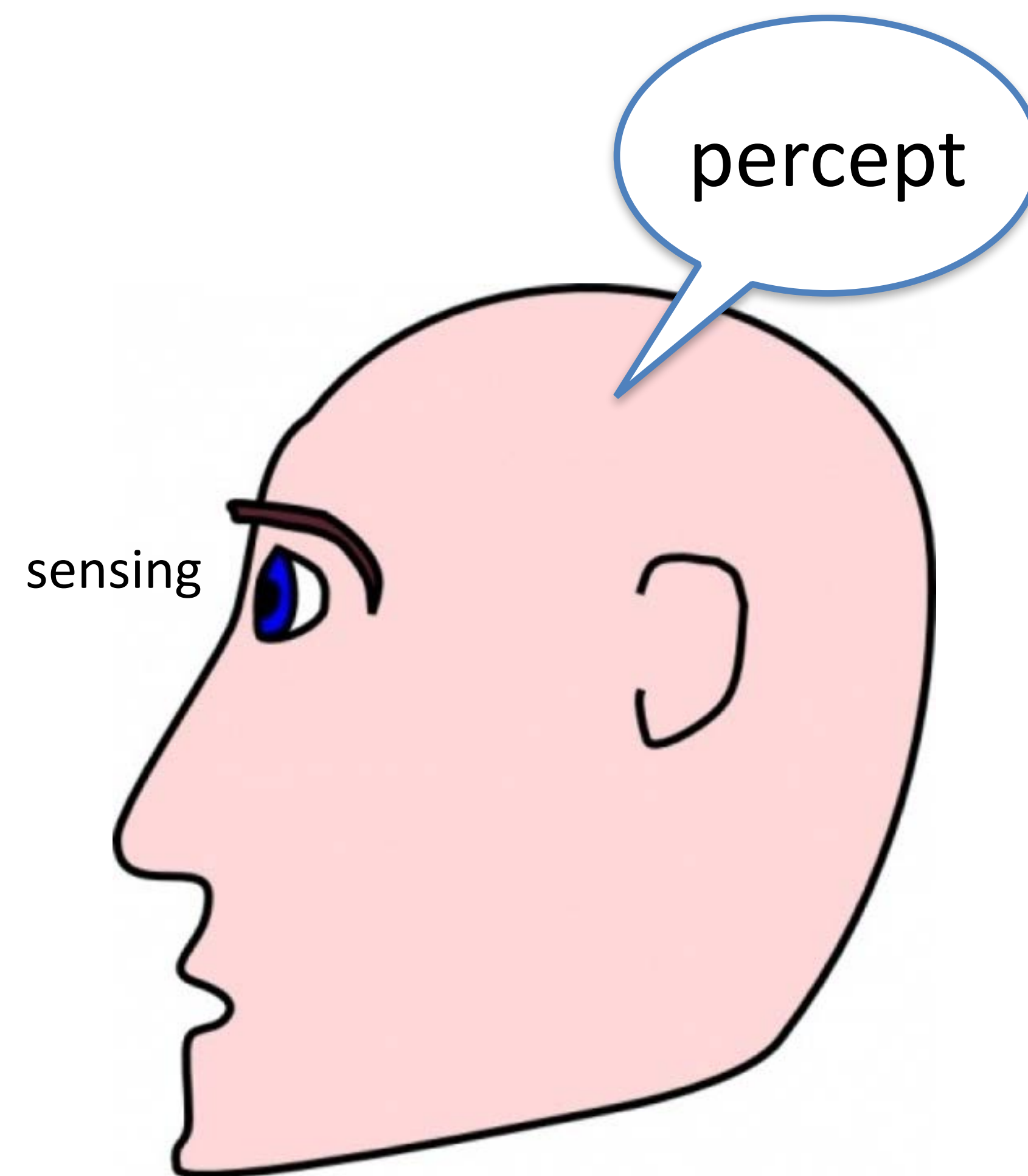
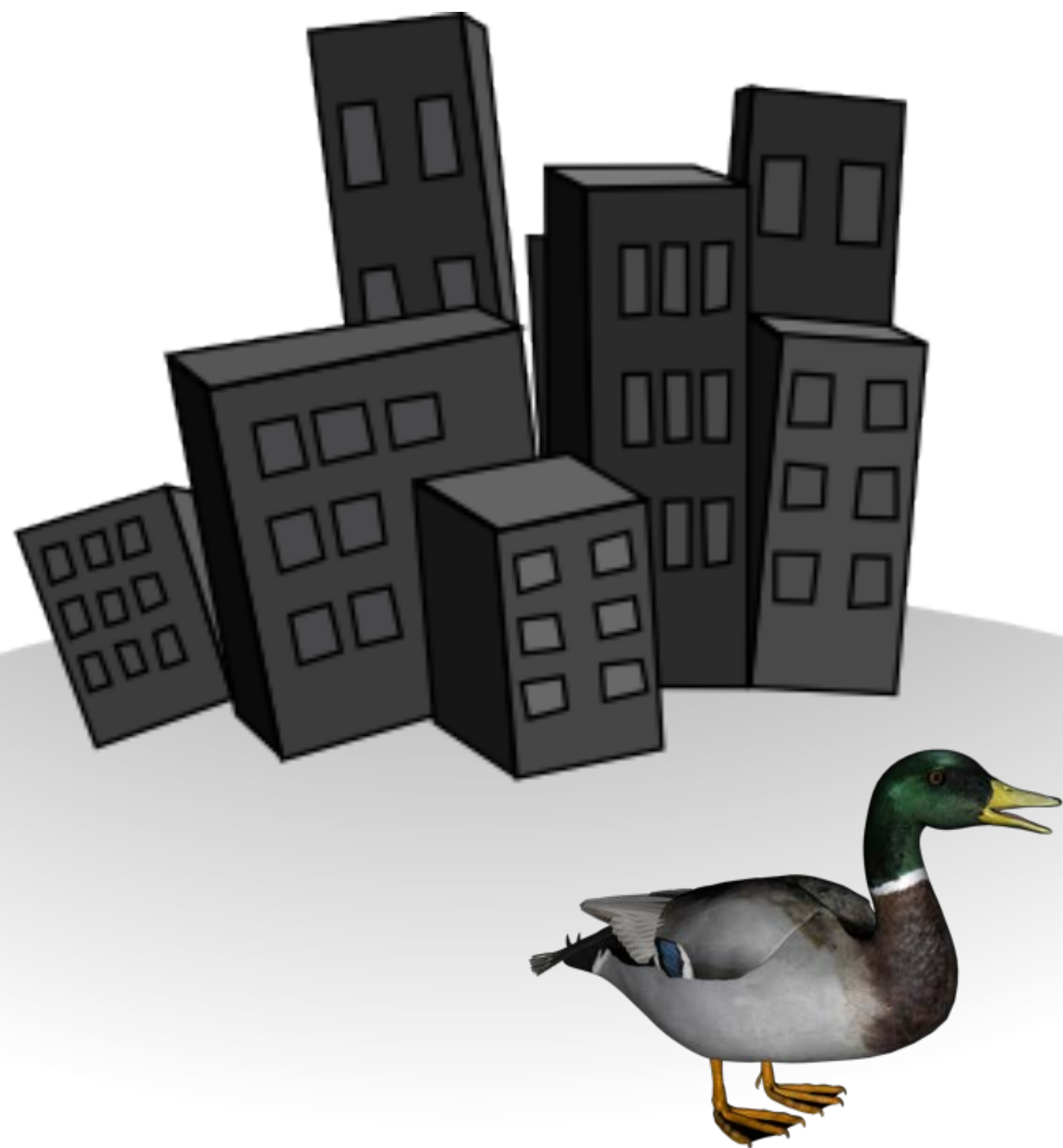




there is no unique inverse



theories of vision





Plato
428-348 BCE



Aristotle
384-322 BCE



Euclid
350-250 BCE?

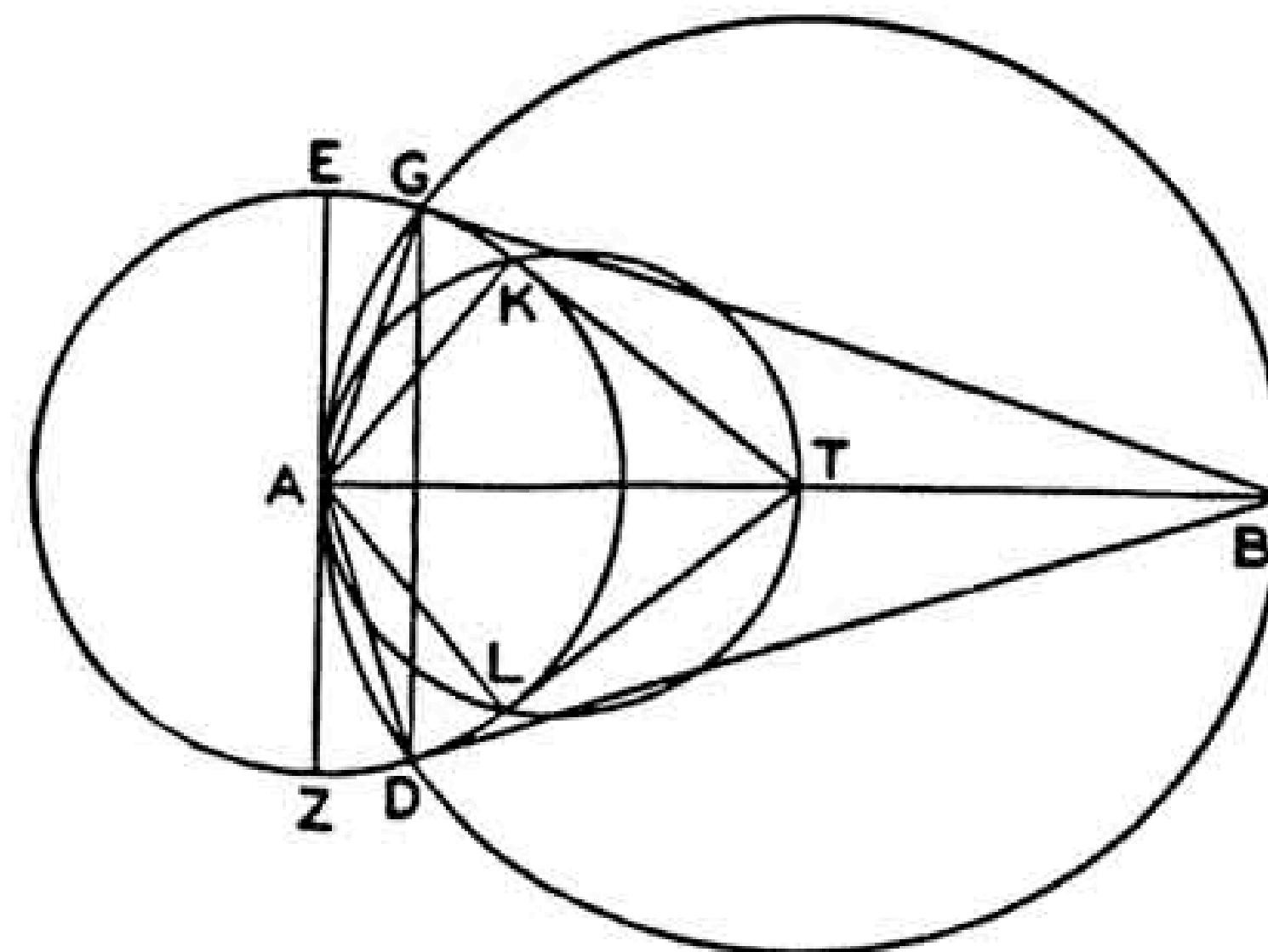
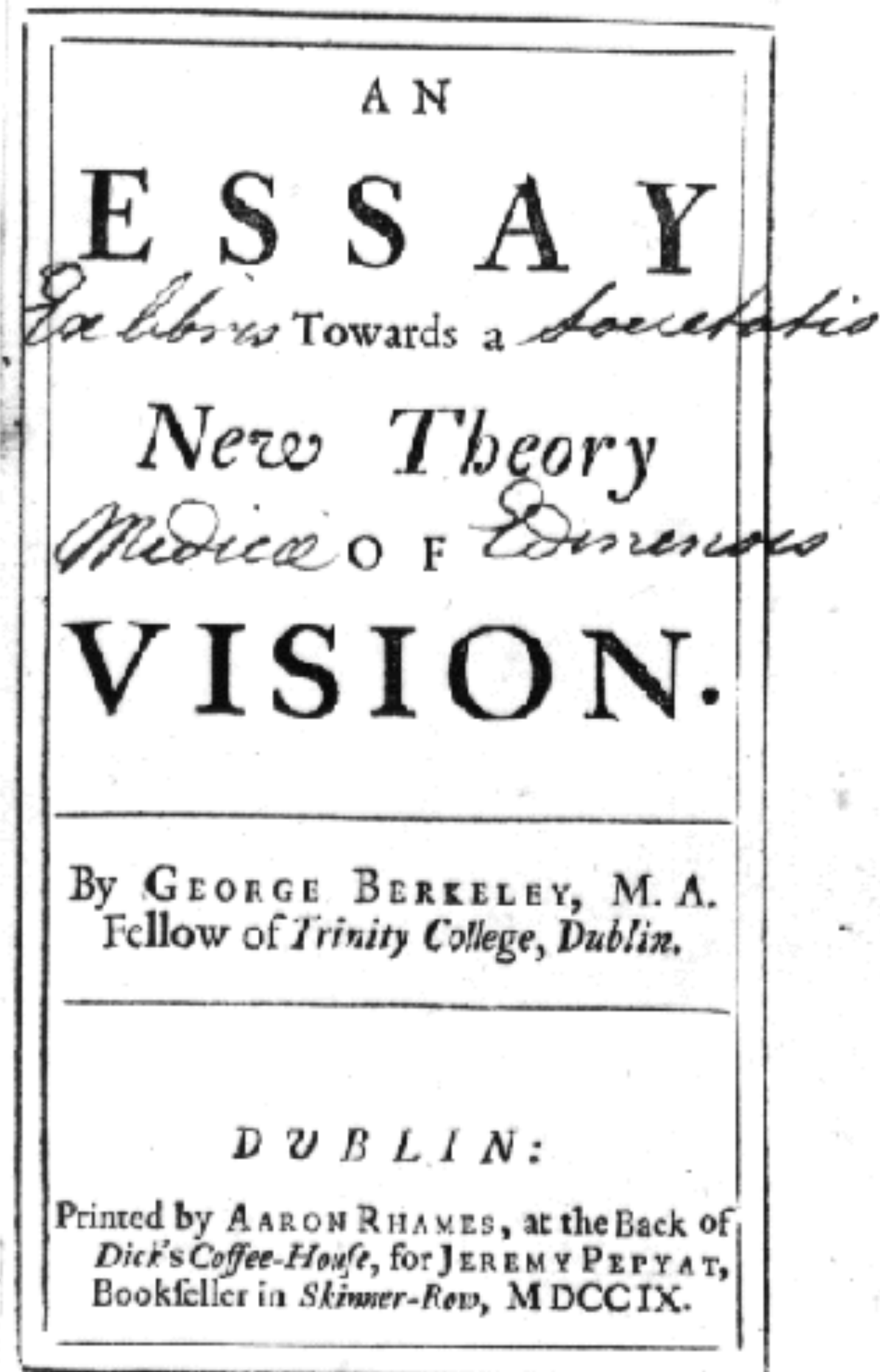


FIG. 24.

When the eye approaches the sphere, the part seen will be less, but will seem to be more. (Fig. 24.)

Let there be a sphere, of which the center is A , and let the eye be B , from which let the straight line AB be drawn. And around AB let the circle GBD be inscribed, and from the point A let the straight line EZ be drawn, perpendicular to the straight line AB in either direction, and let the plane be produced along EZ and AB . So it will make a circular section. Let it be $GEZD$, and let GA , AD , DB , BG , and GD be drawn. So, according to the theorem given before, the angles at the points G and D are right angles. Thus, BG and BD , whatever rays there are, touch the sphere. And the part of the sphere, GD , is seen by the eye, B . Now let the eye be moved nearer to the sphere, and let it be at T , from which let the straight line TA be drawn, and let the circle ALK be inscribed, and let the straight lines TK , KA ,



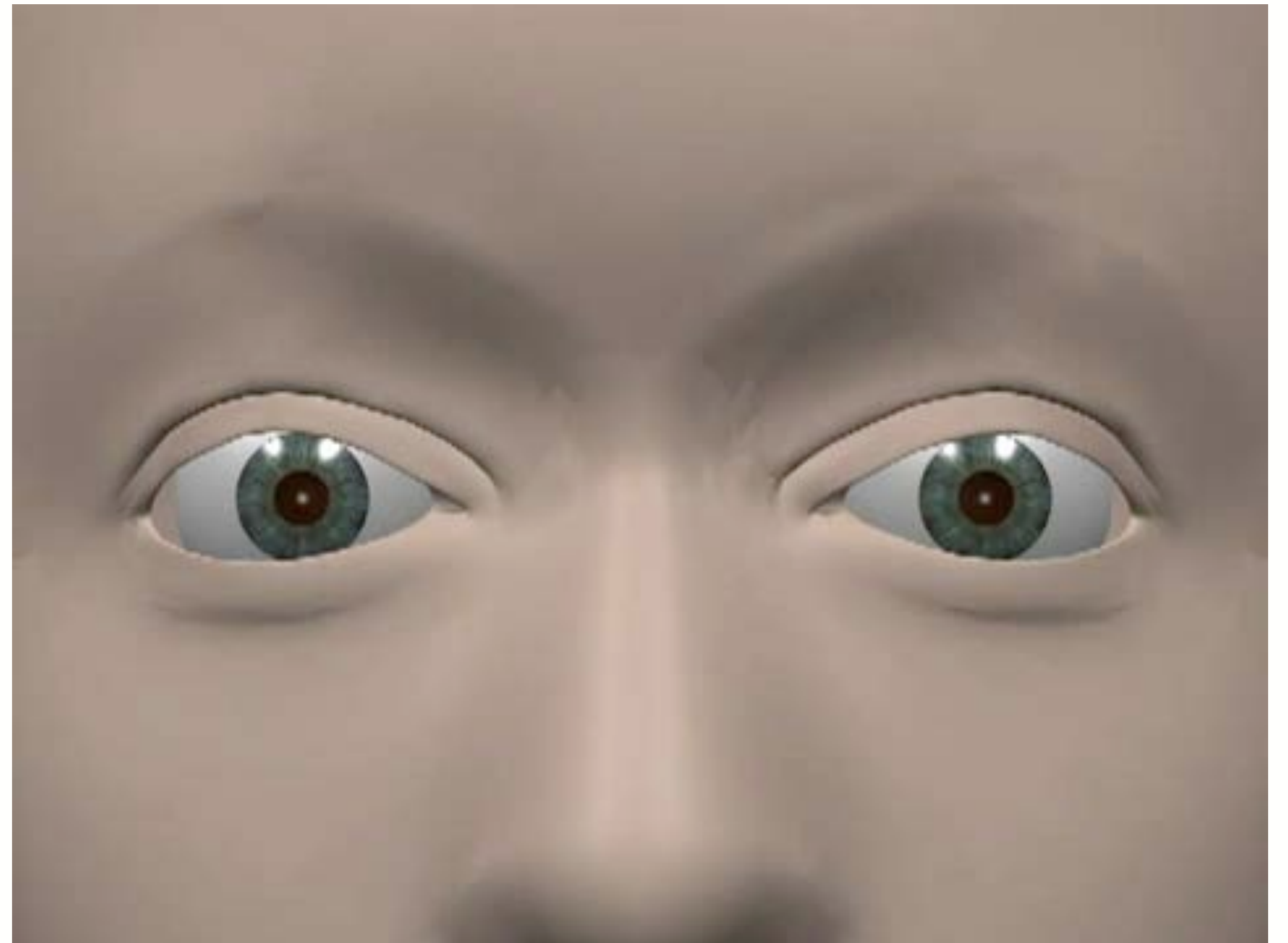
- Distance is determined indirectly by visual cues.
- We learn to separate size from distance.



Binocular stereo



Accommodation



Convergence



Occlusion

distance: bush < cat < pillar



Height in visual field



Relative size



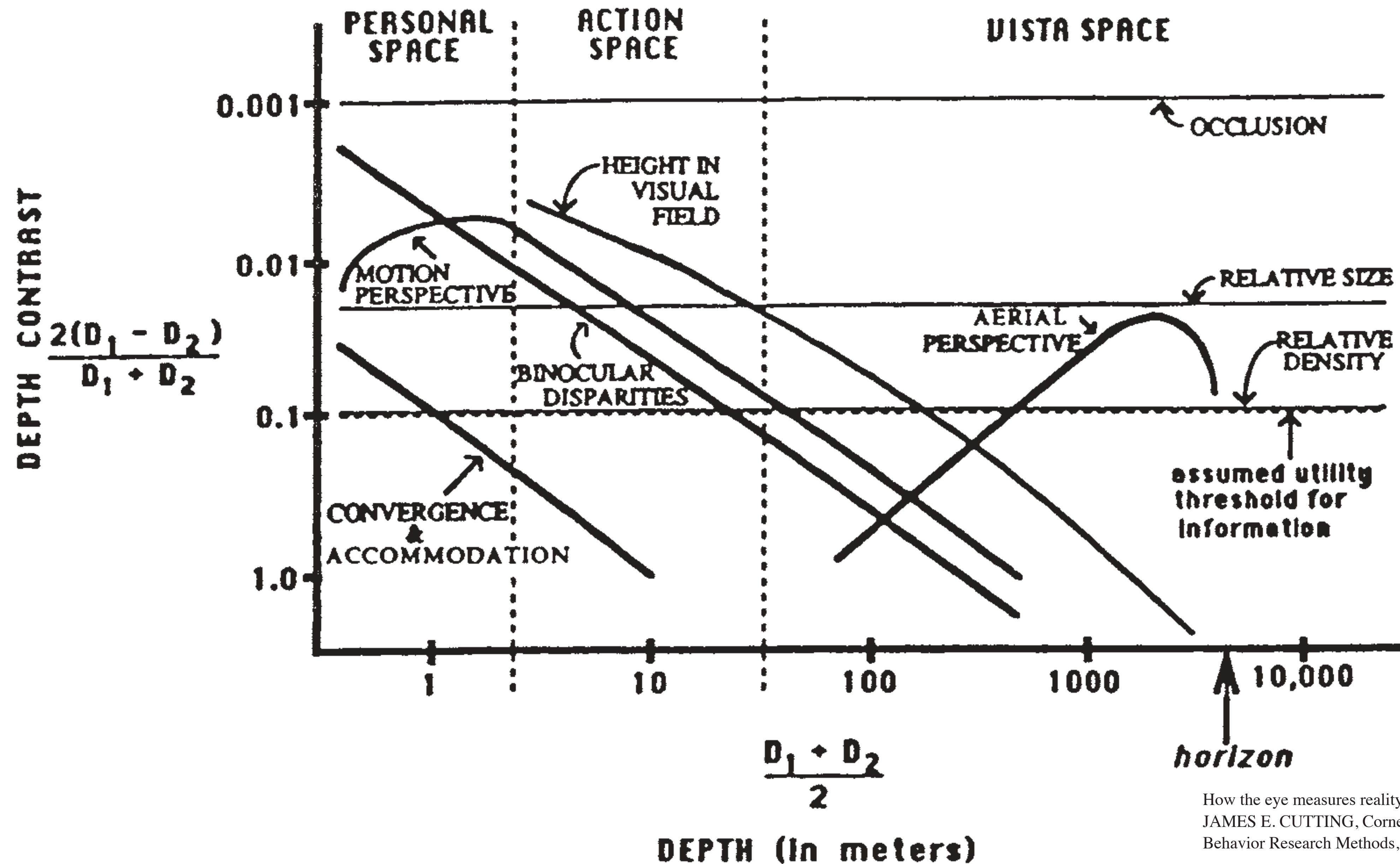


Texture density

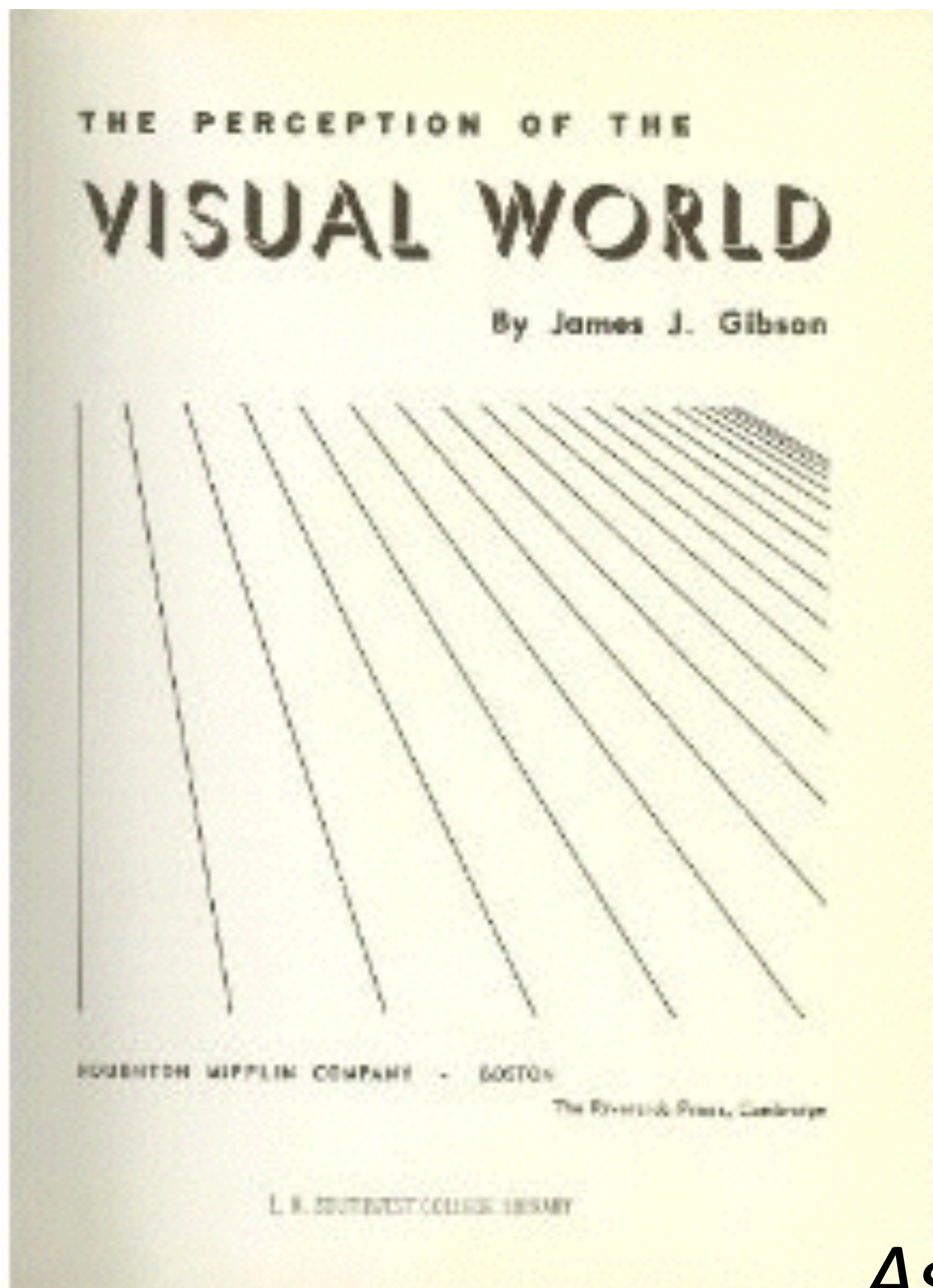


Aerial perspective

Determining distance



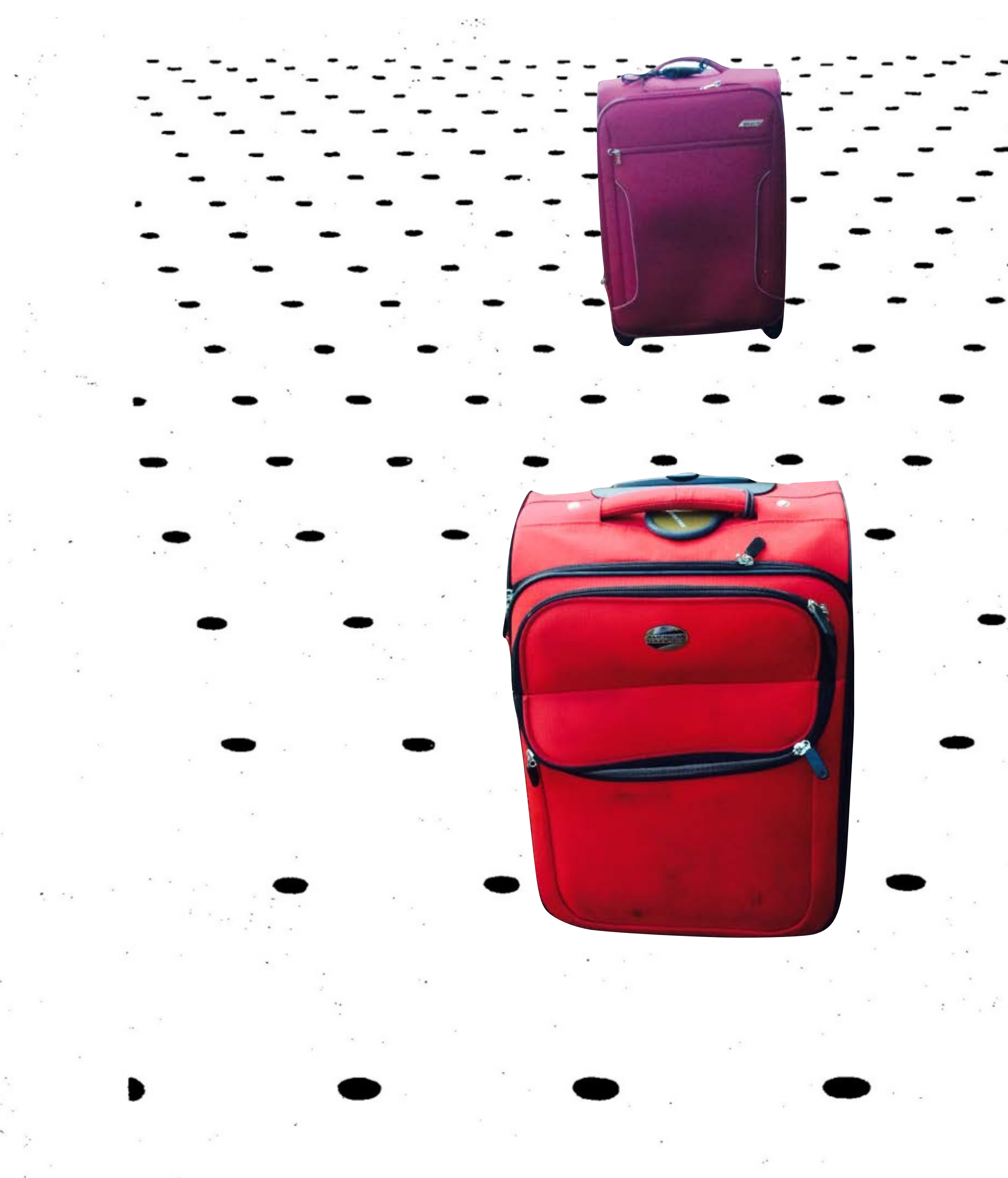
How the eye measures reality and virtual reality.
JAMES E. CUTTING, Cornell University, Ithaca, New York.
Behavior Research Methods, Instruments, & Computers 1997, 29 (1), 2

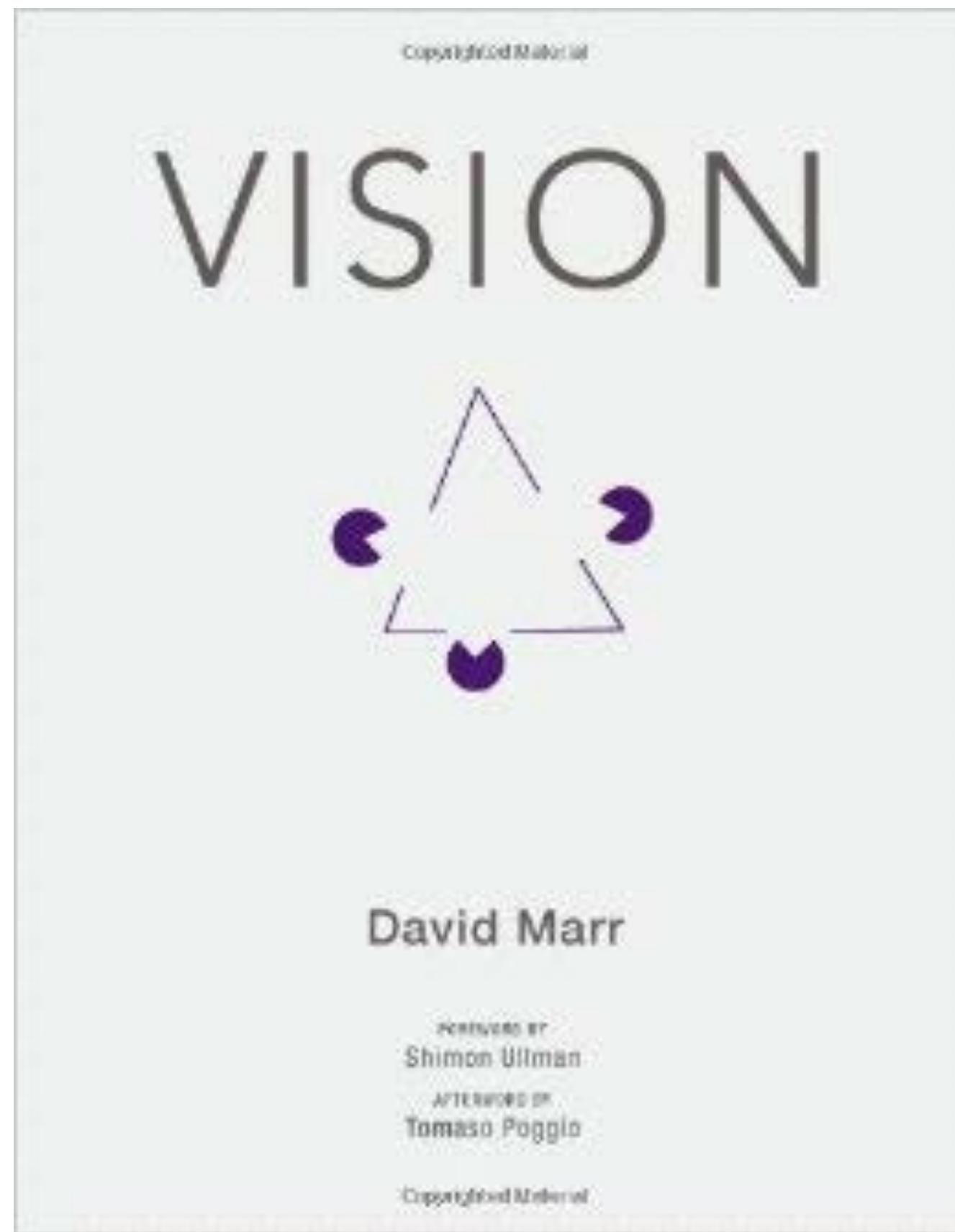


1950

ecological vision

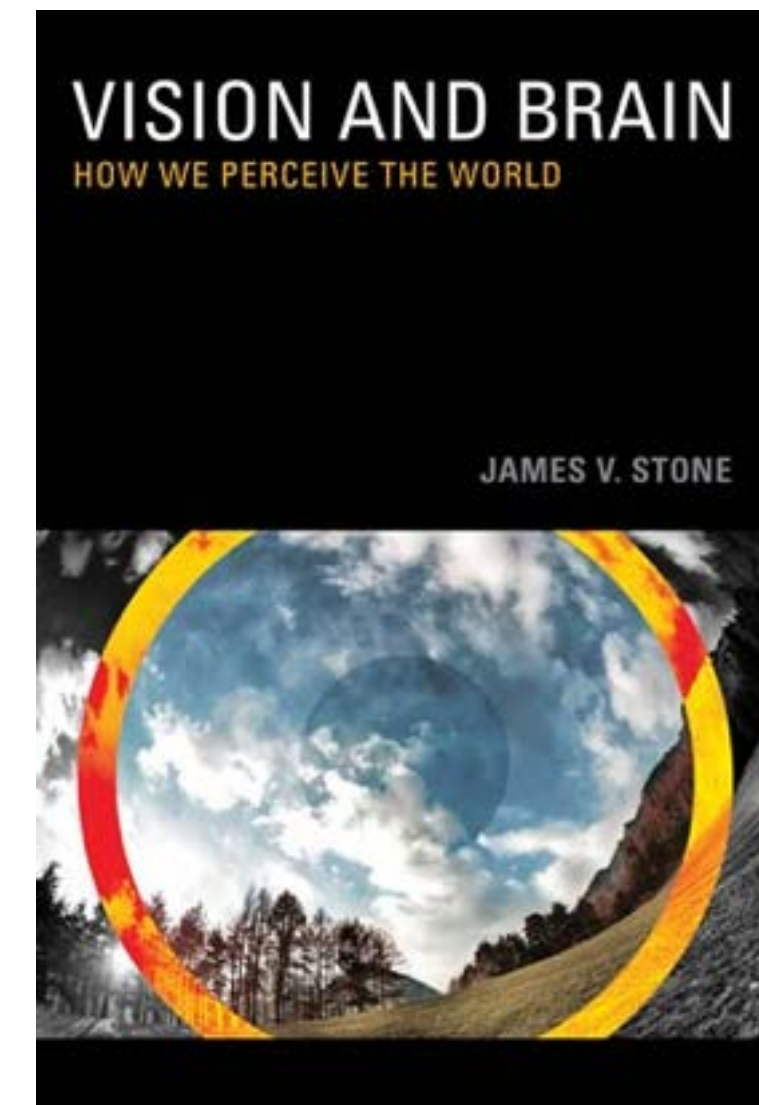
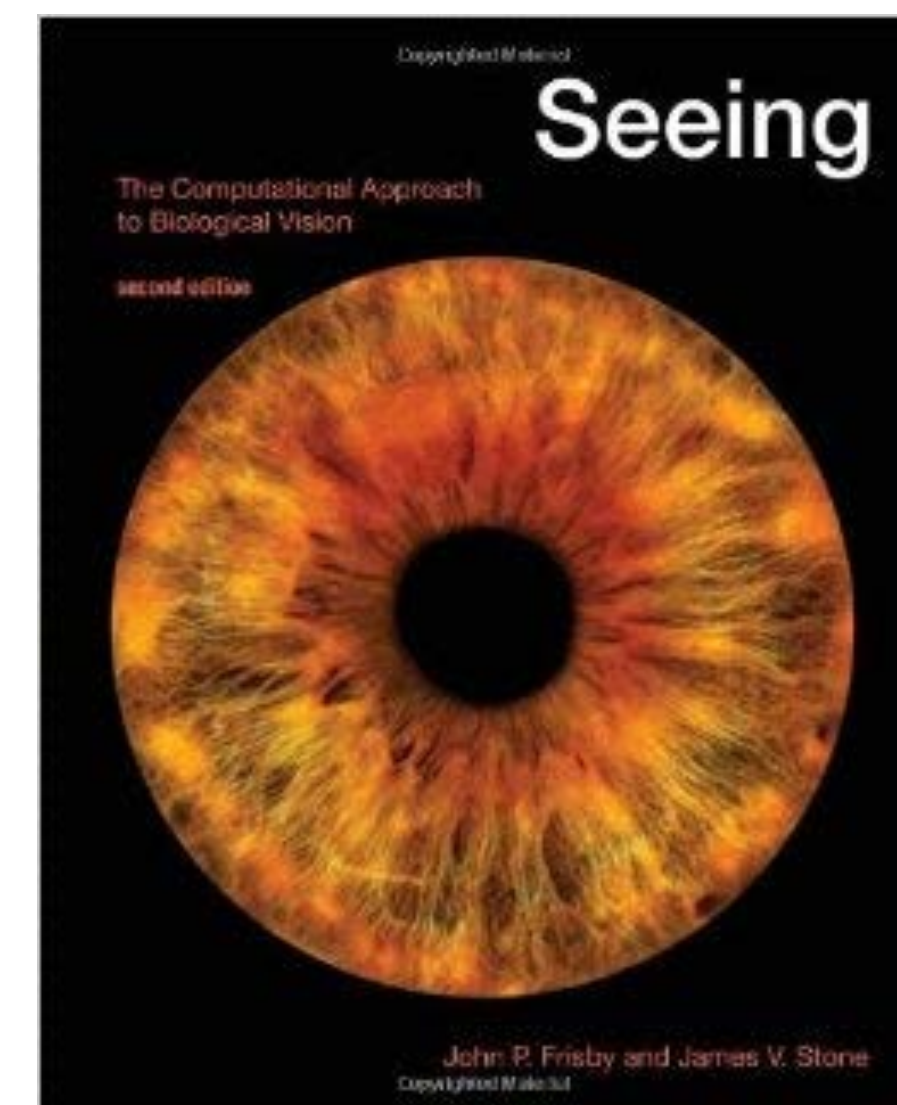
*Ask not what's
inside your head,
but what your
head's inside of.*





1982

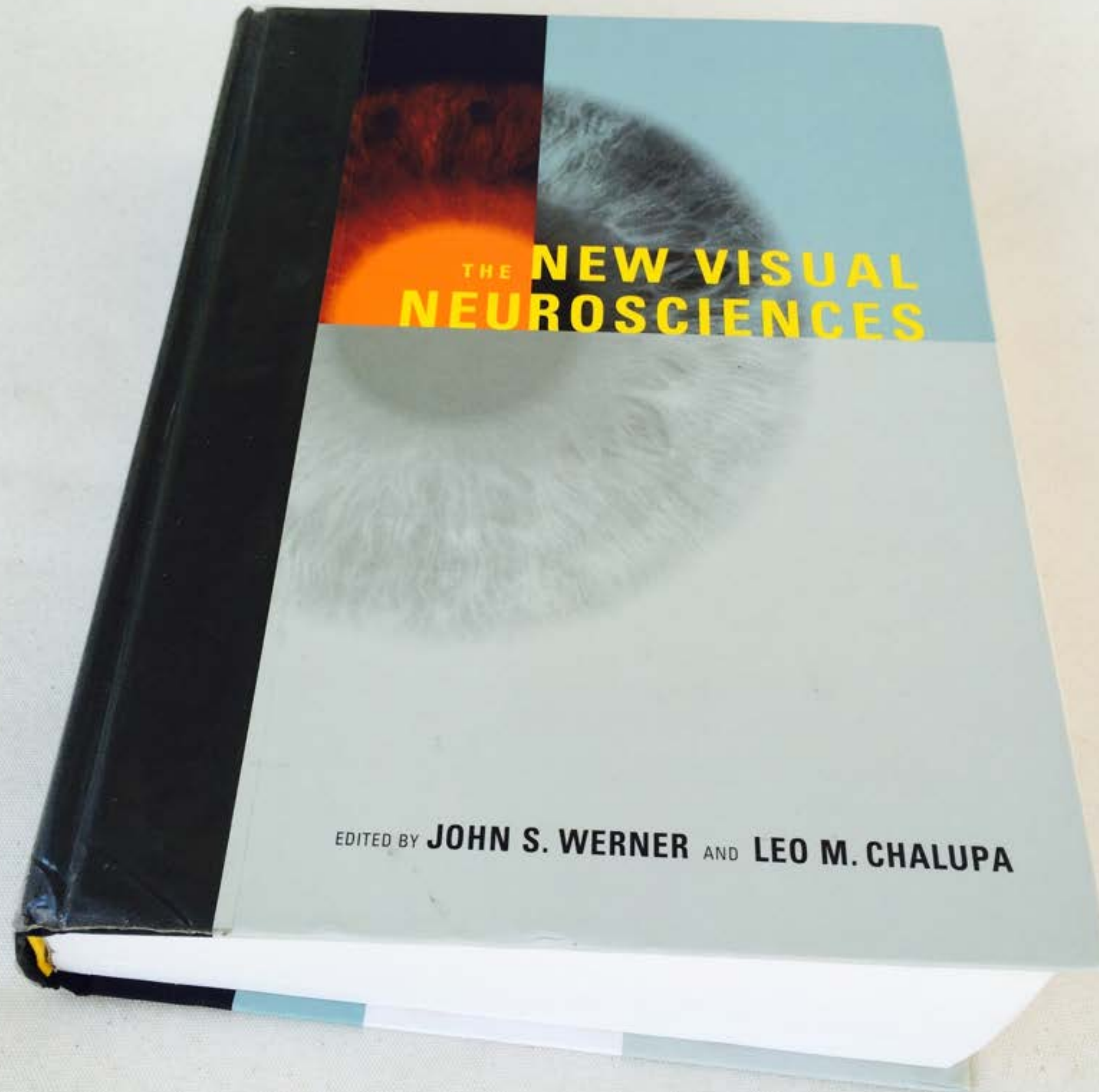
- Computational theories
- Algorithms
- Circuitry





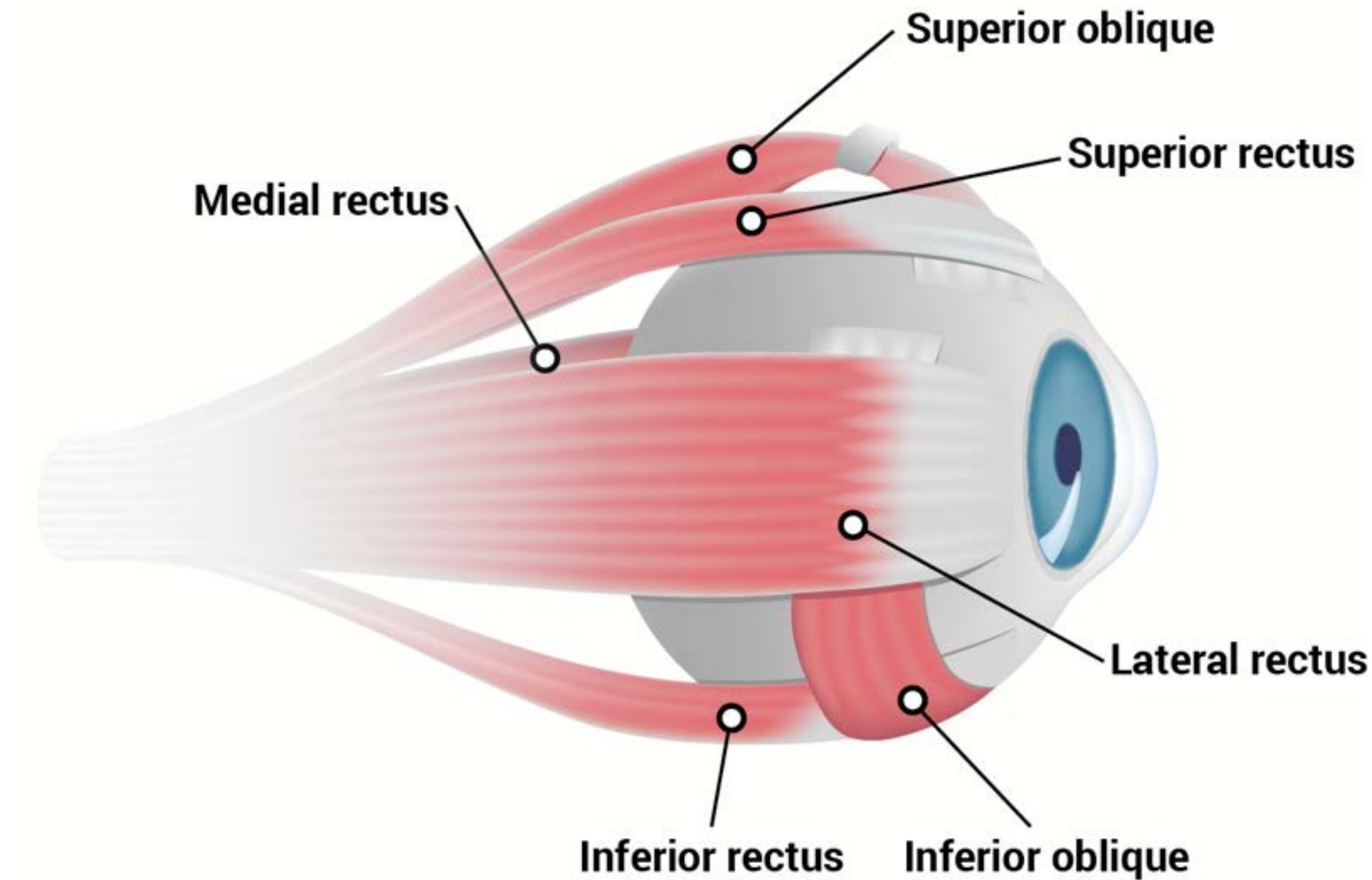
Human Brain Project

The BRAIN initiative

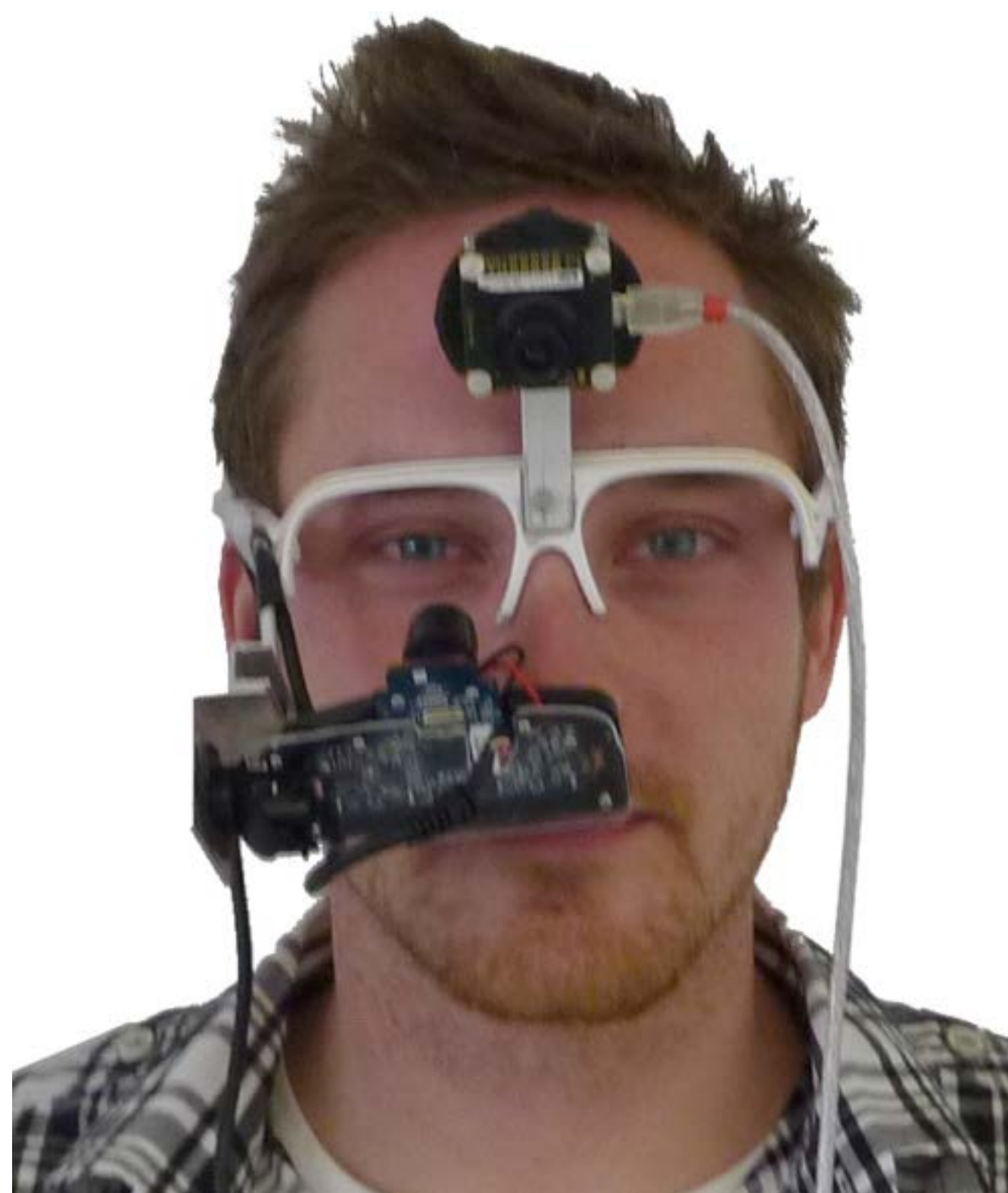


Visual neuroscience growing at a massive rate

Seeing is an **active** process



- Fastest moving muscles in the human body
 - rotational speed: 600deg/s
 - rotational acceleration: 35,000 deg/s²



What are the material
circumstances of the
family?



Figure 109 "Seven Records of eye movements by the same subject" 1967
Yarbus, A. L.

Artwork: *Unexpected visitors*
Ilya Repin | 1884-1888

What age are the
figures in the painting?



Figure 109 "Seven Records of eye movements by the same subject" 1967
Yarbus, A. L.

Artwork: *Unexpected visitors*
Ilya Repin | 1884-1888

What type of clothes
are the family wearing?



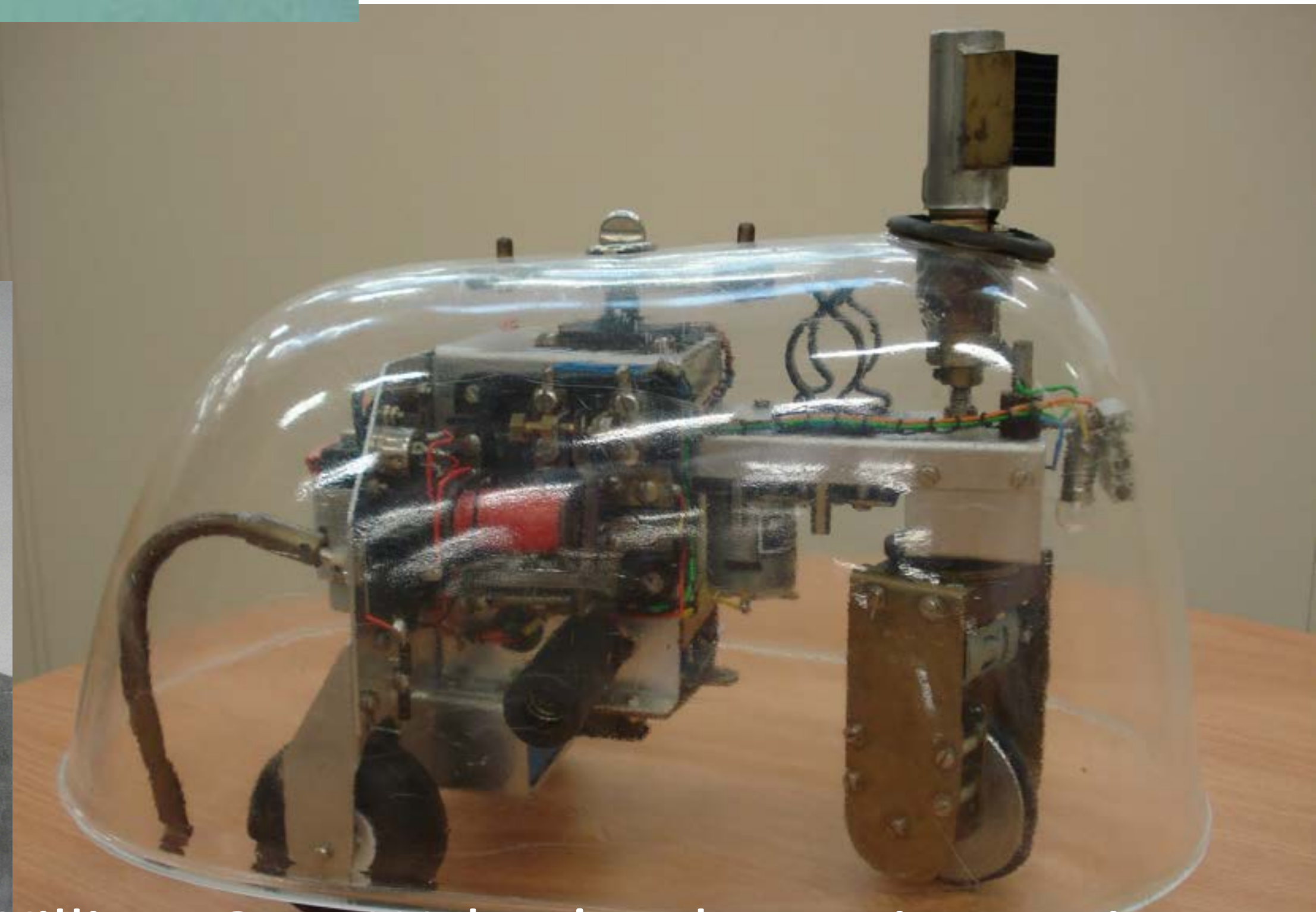
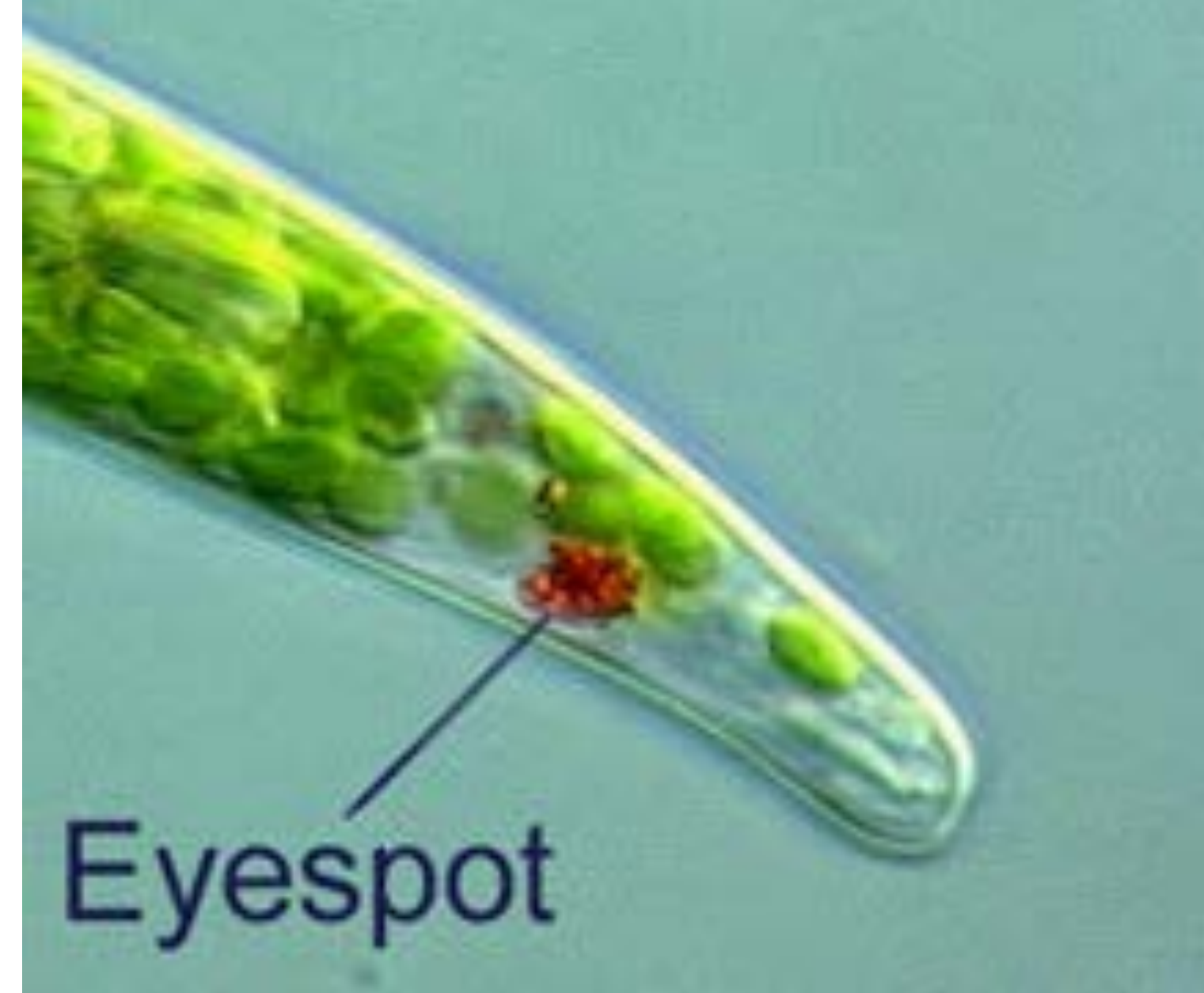
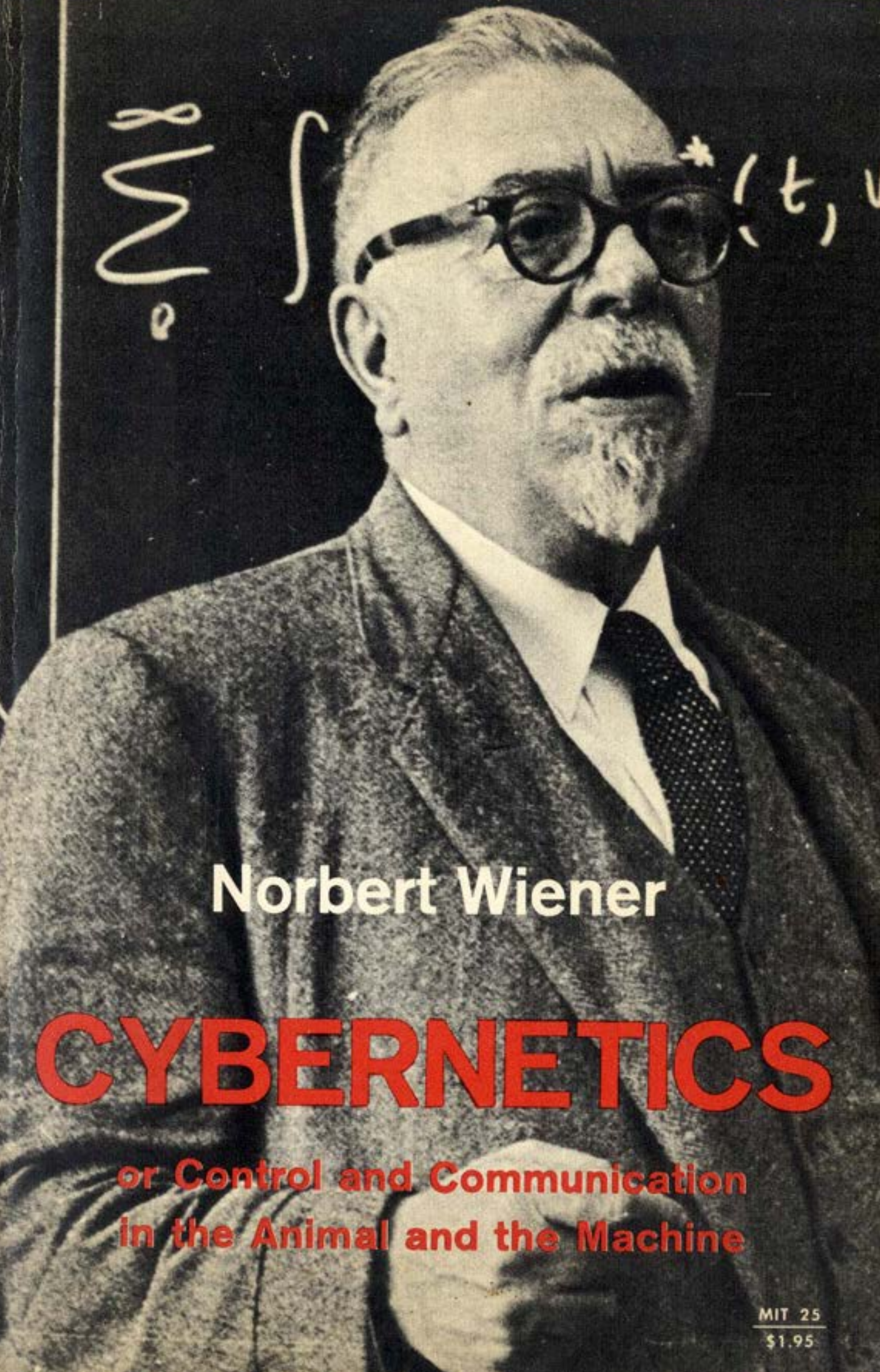
Vision is the process of **discovering** from images what is present in the world and where it is.

David Marr (1982)

Active vision/perception

- Active Vision, Yiannis Aloimonos et al., IJCV, 1988
- Active Perception, Ruzena Bajcsy, Proc IEEE, 76(8) August 1988

robotics and vision



William Grey Walter's cybernetic tortoise

A Proposal for the Dartmouth Summer Research Project on Artificial Intelligence

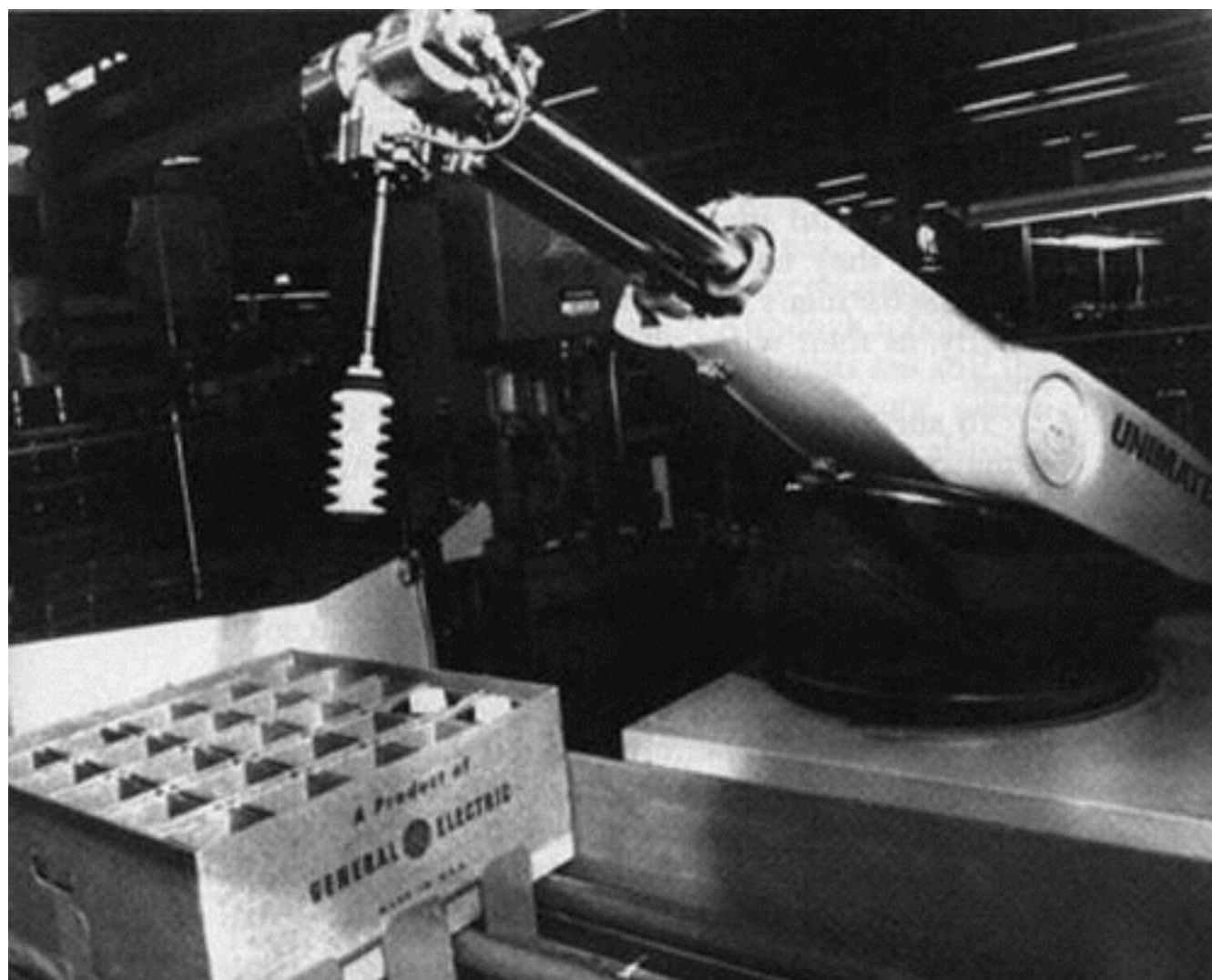
August 31, 1955

*John McCarthy, Marvin L. Minsky,
Nathaniel Rochester,
and Claude E. Shannon*

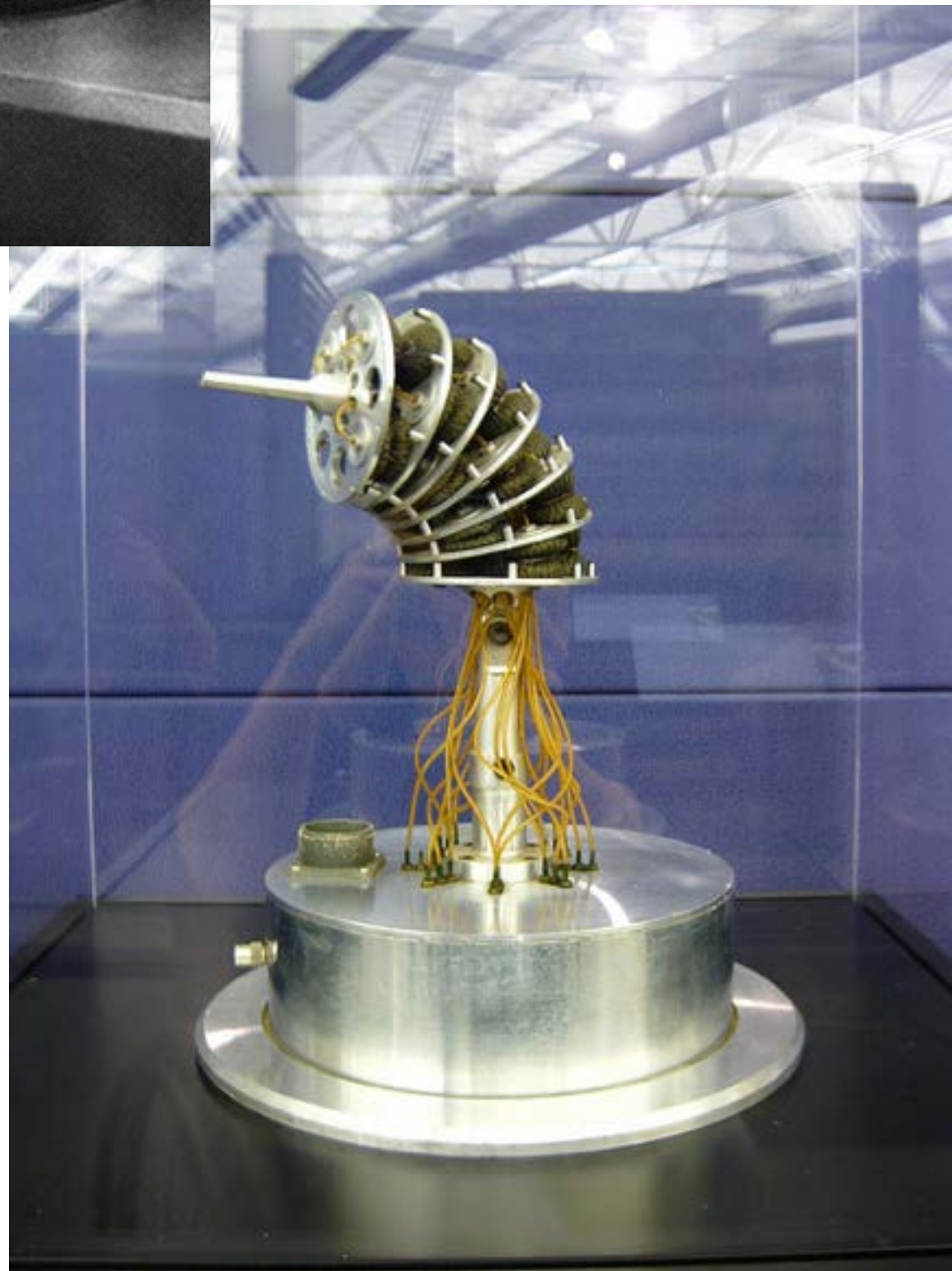


The study is to proceed on the basis of the conjecture that every aspect of learning or any other feature of intelligence can in principle be so precisely described that a machine can be made to simulate it...

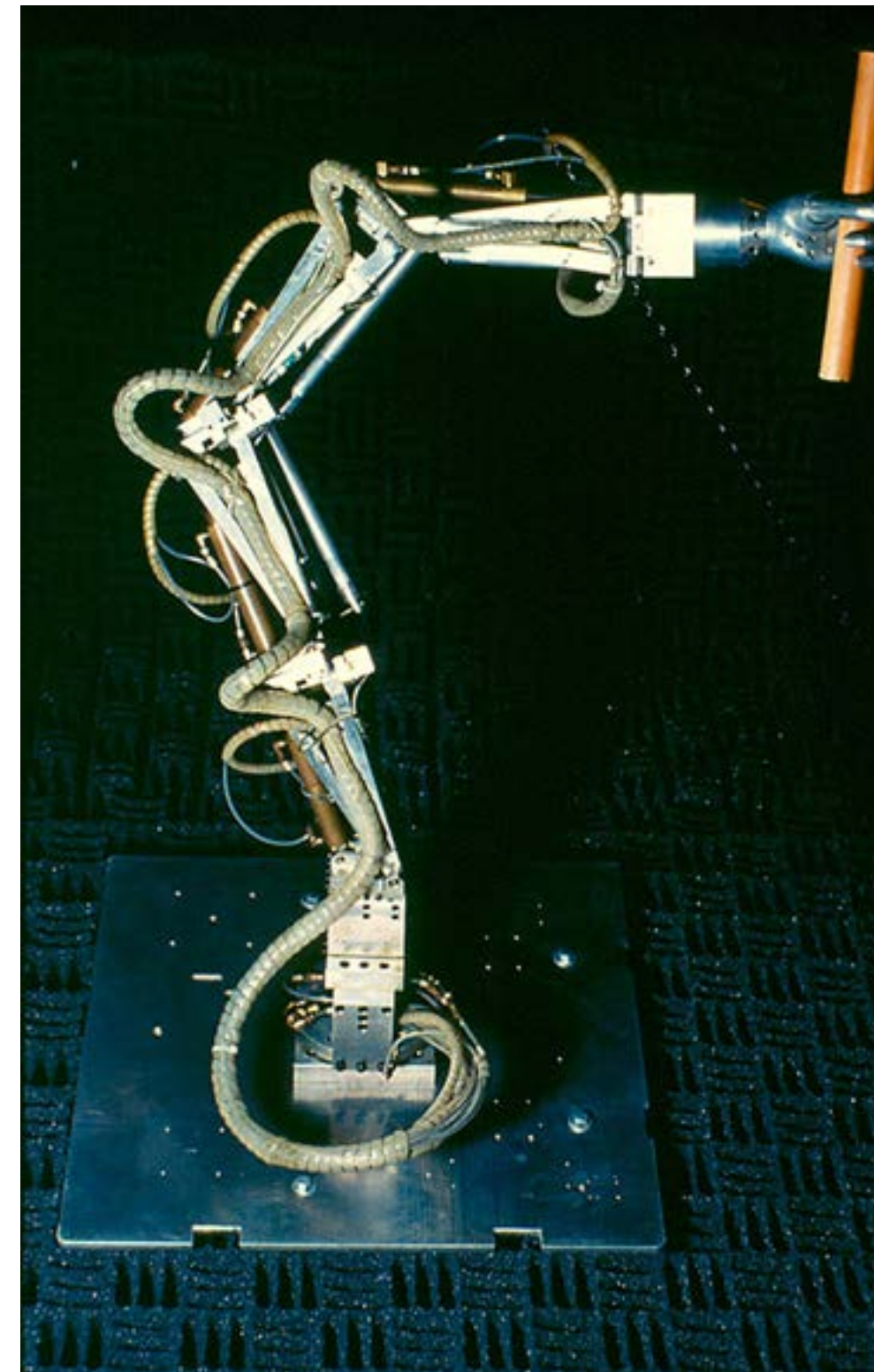
- computers,
- natural language processing,
- neural networks,
- theory of computation,
- abstraction and
- creativity



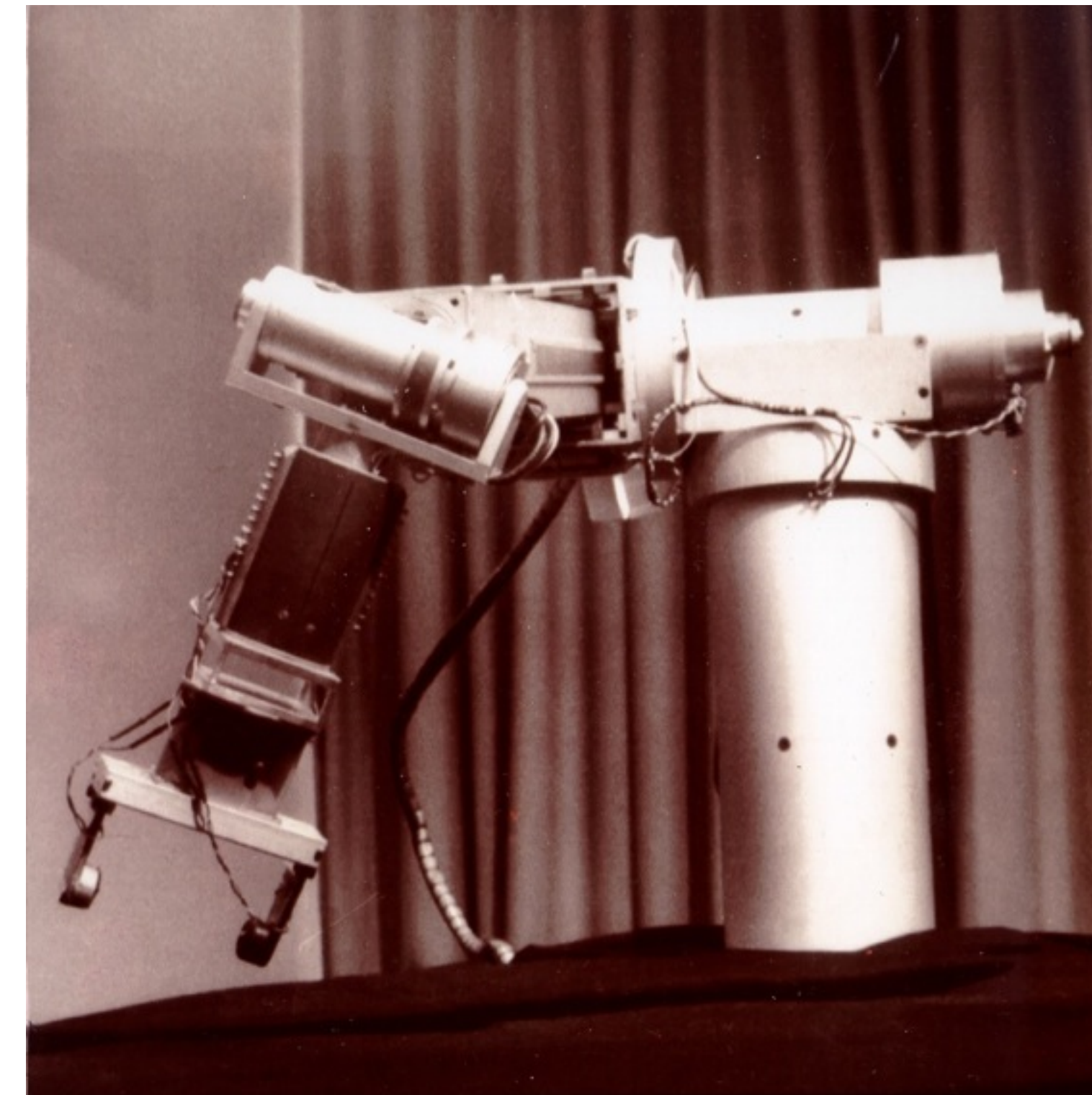
Unimate
 Devol
 Unimation Inc. 1961



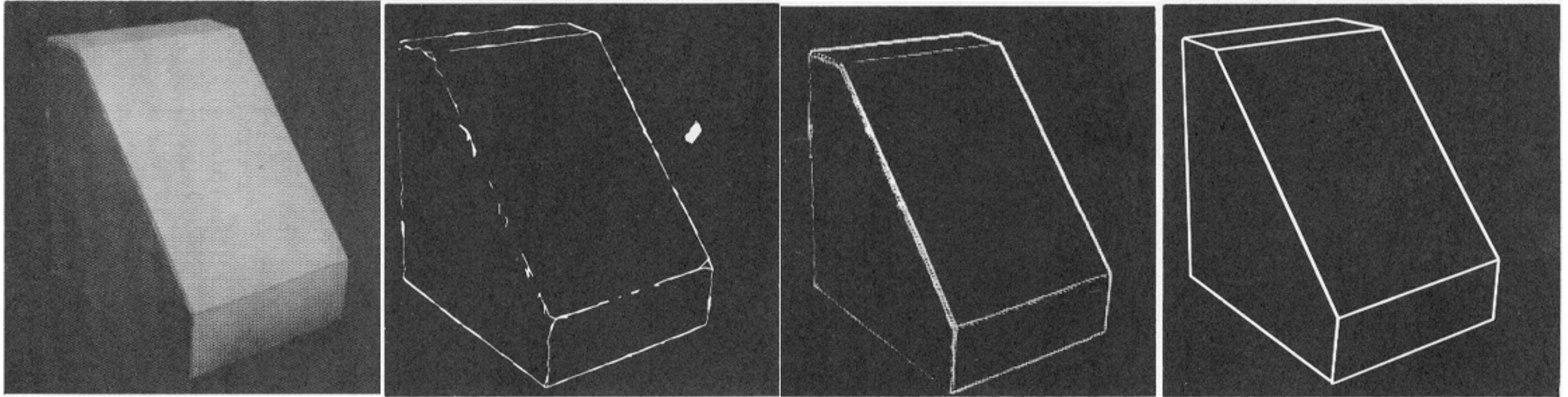
ORM
 Scheinman & Leiffer
 Stanford 1965



Tentacle arm
 Minsk
 MIT 1968





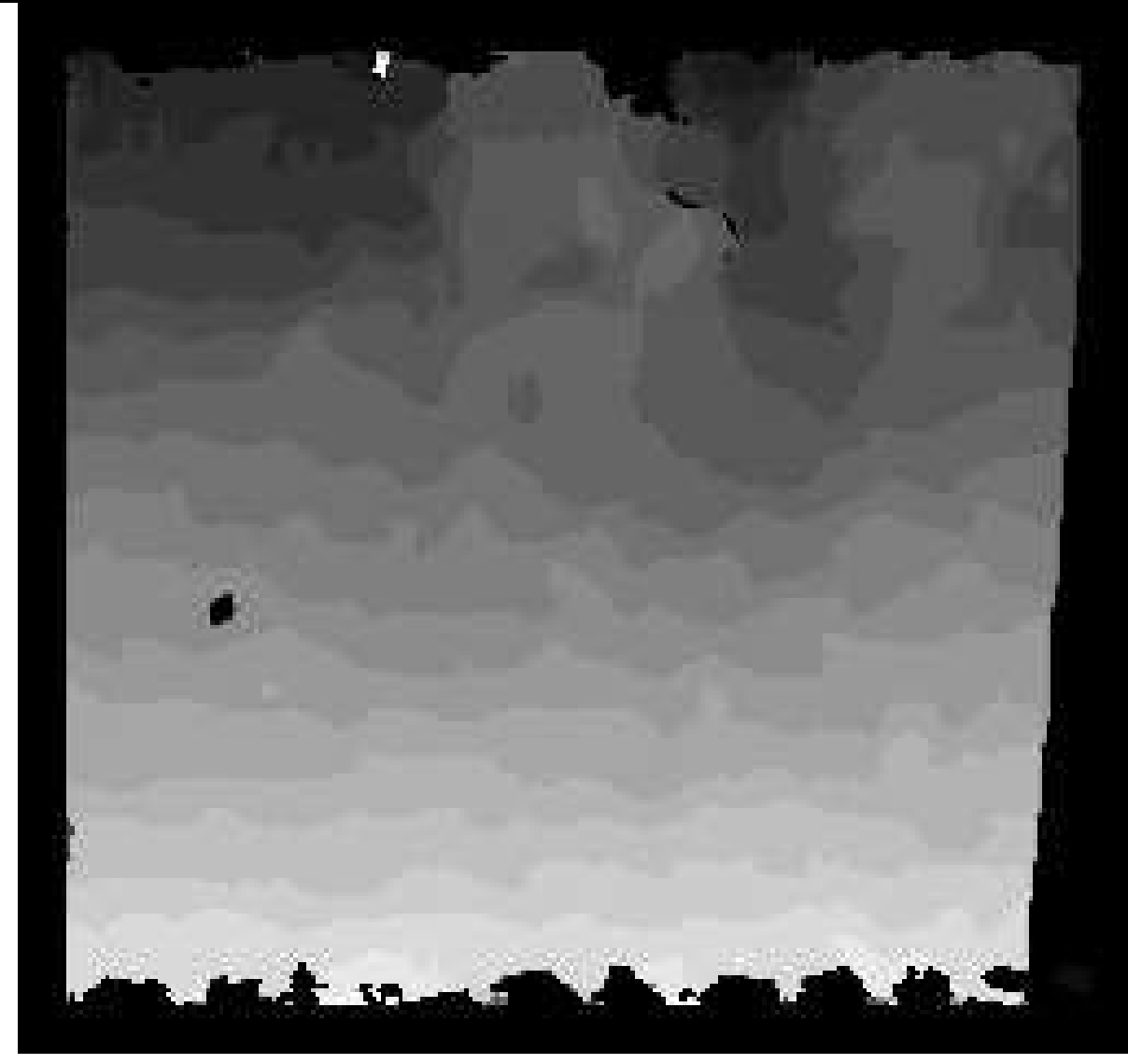


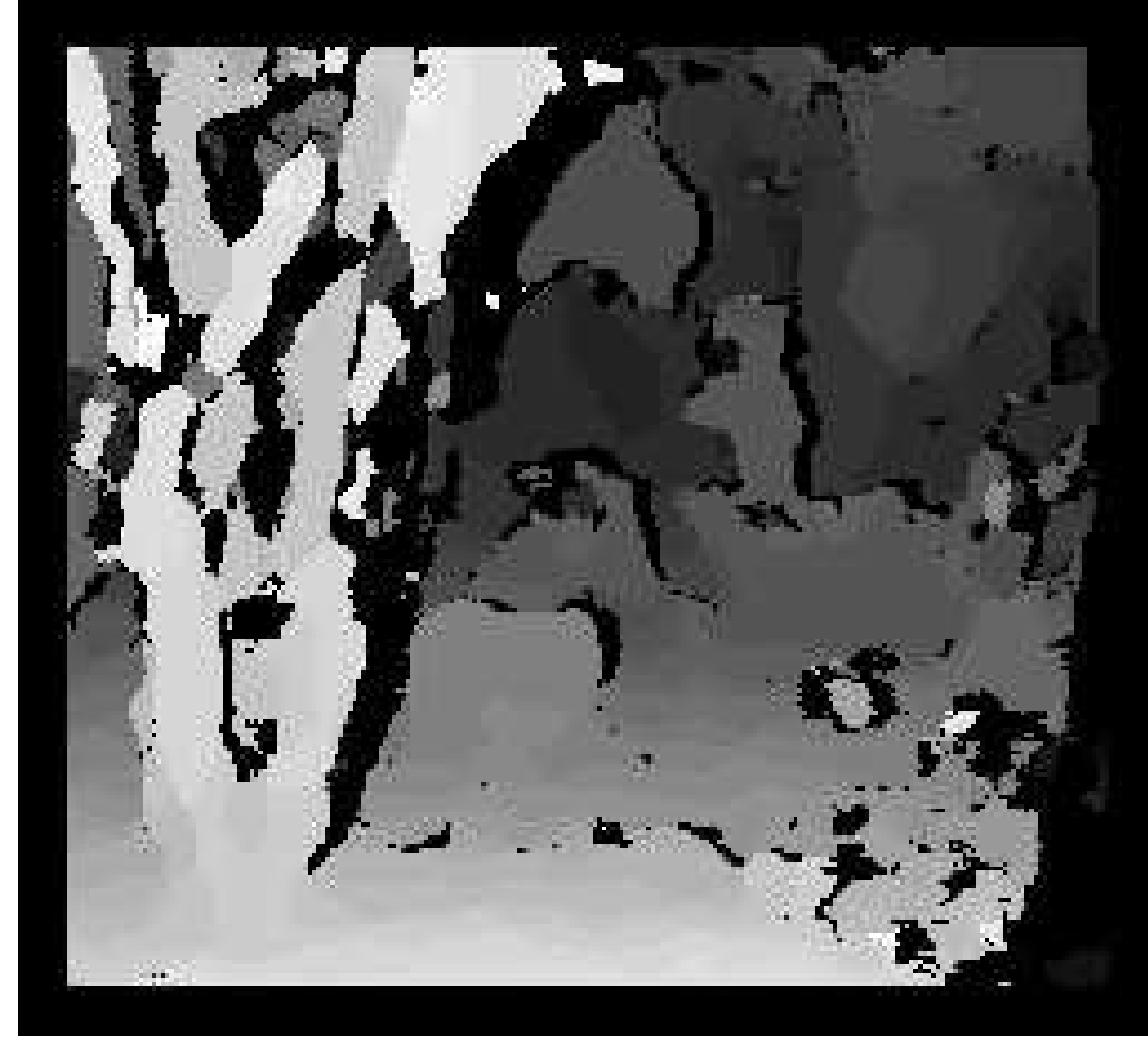
Stanford arm
 Scheinman
 Stanford 1968



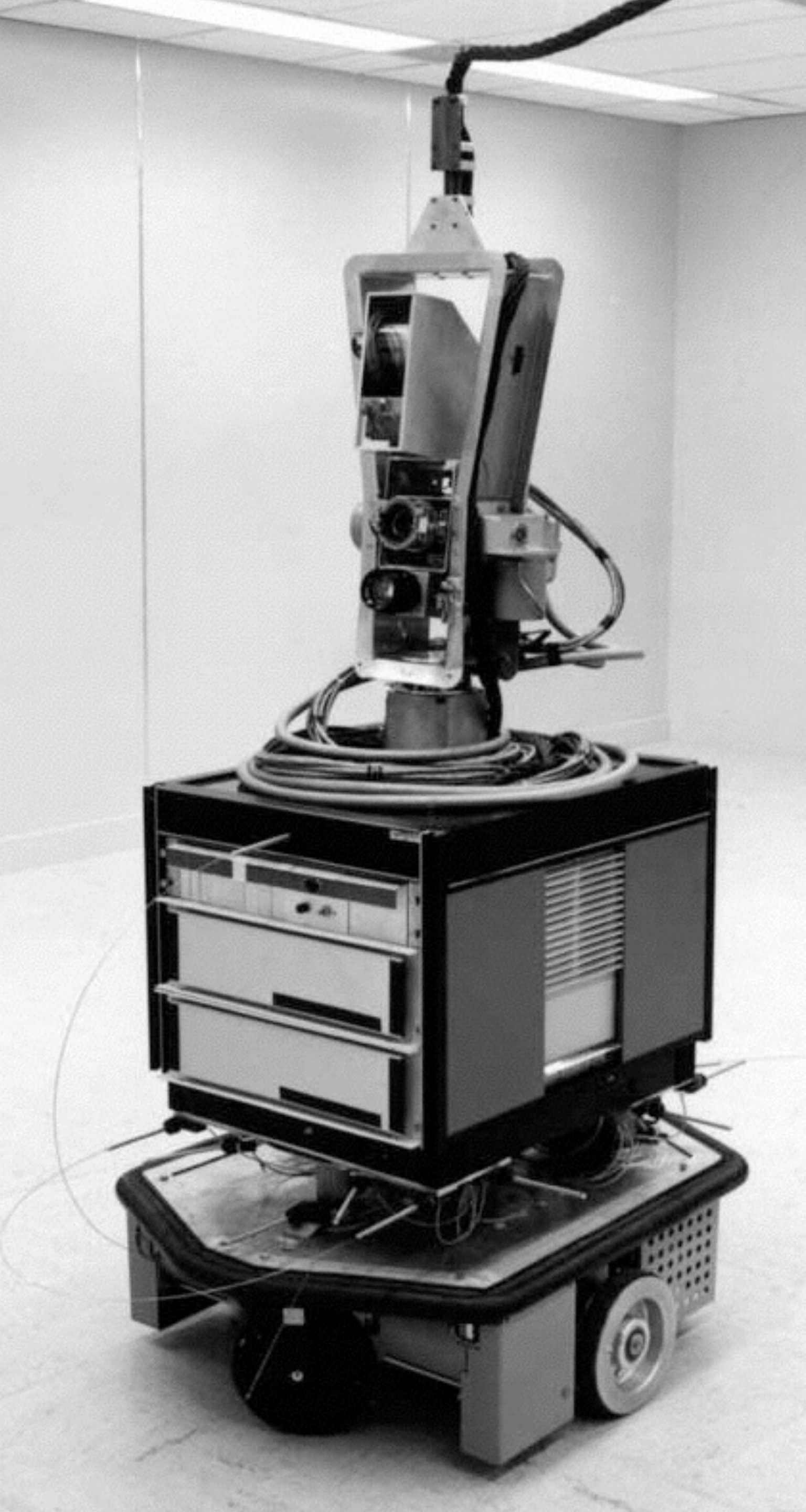
Machine perception of three-dimensional solids

Larry Roberts

MIT 1964

Left	Right	Depth
		
		

JISCT data set and early stereo results
1993



Shakey
SRI
1966-72



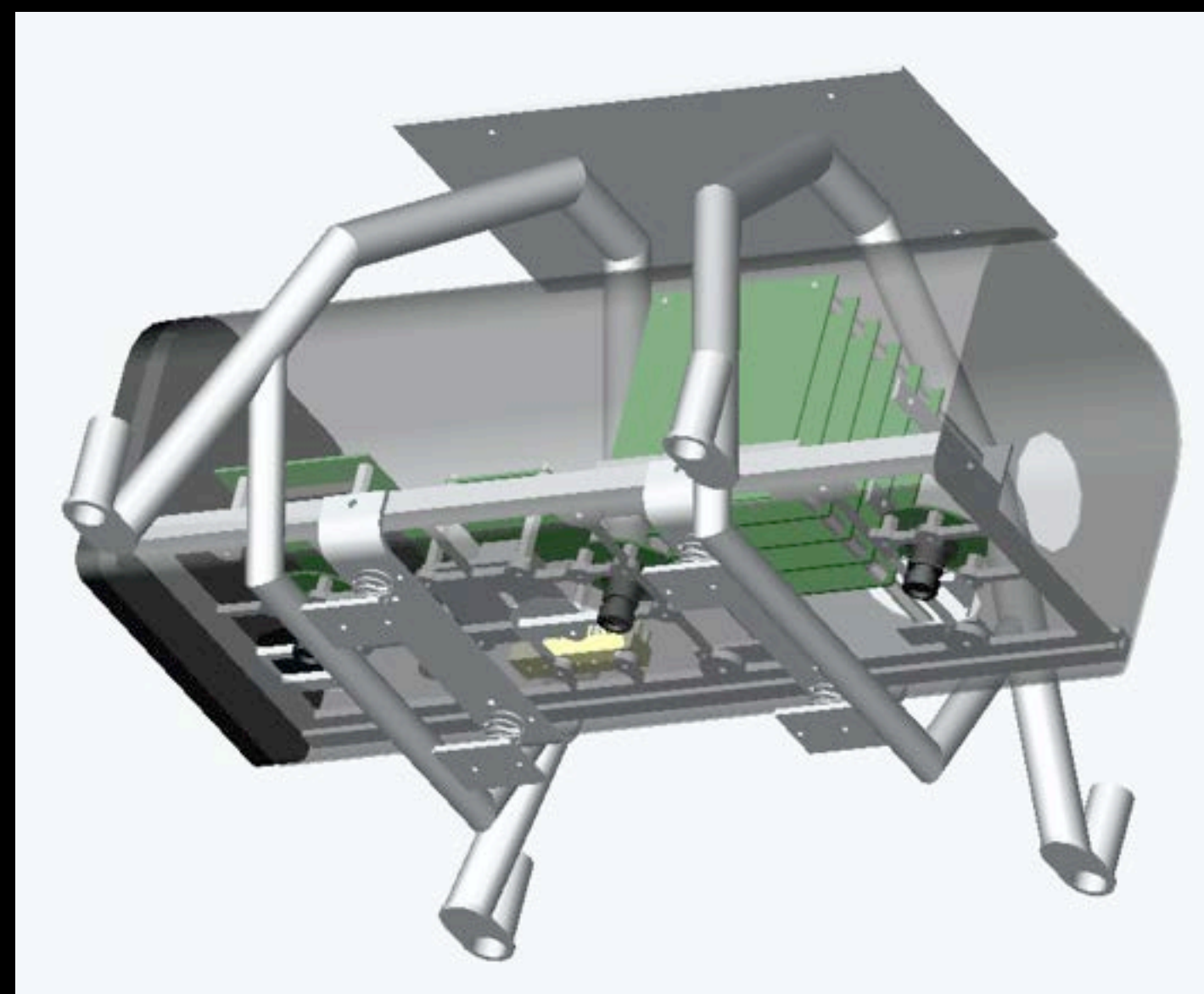
Stanford Cart
Moravec
1971-80



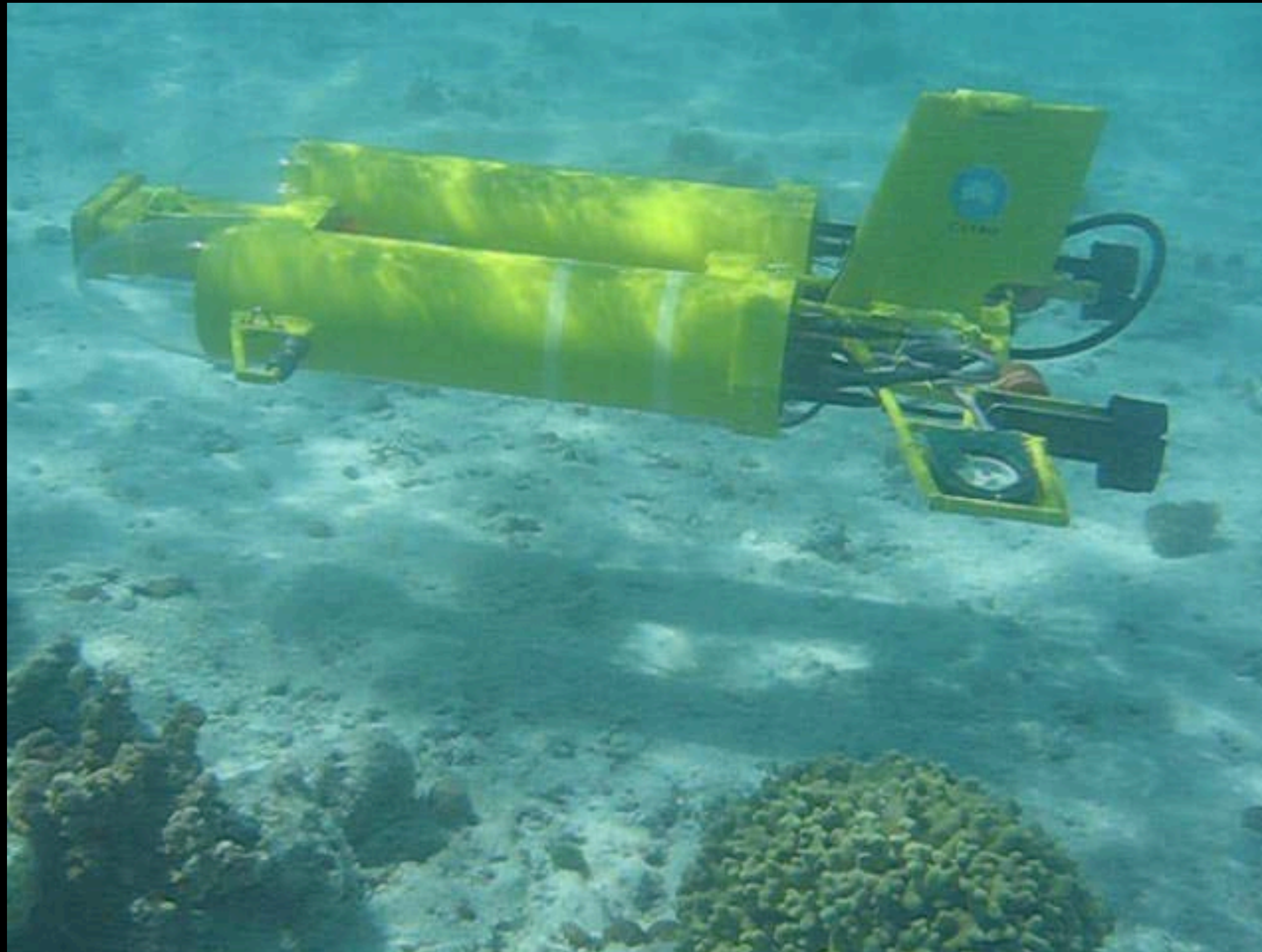
Corke
GRASP lab + CSIRO
1988-92



Steady velocity tests



Buskey, Roberts, Corke
CSIRO
2004

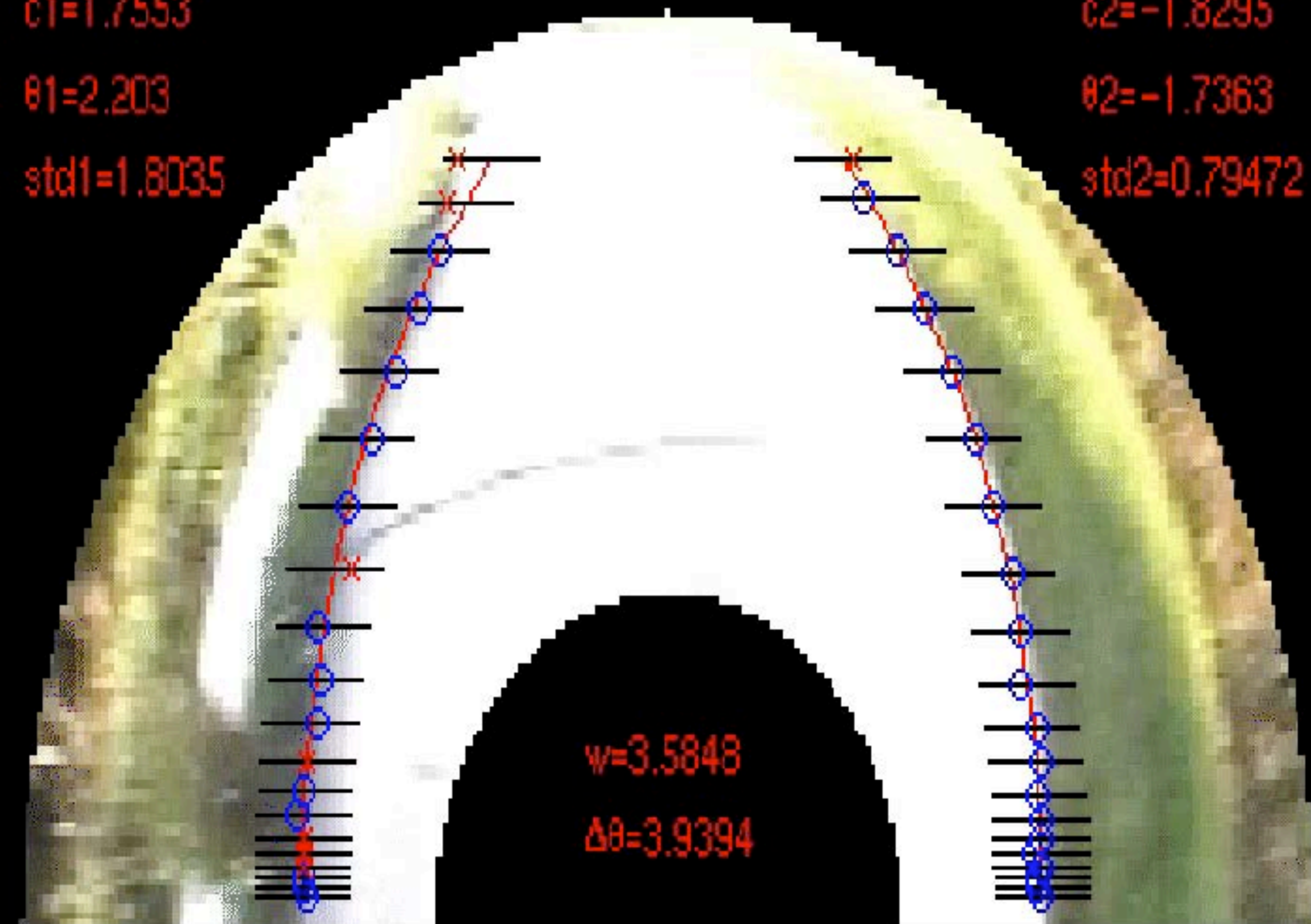


Corke and Dunbabin
CSIRO
2002



$c1=1.7553$
 $\theta1=2.203$
 $std1=1.8035$

$c2=-1.8295$
 $\theta2=-1.7363$
 $std2=0.79472$

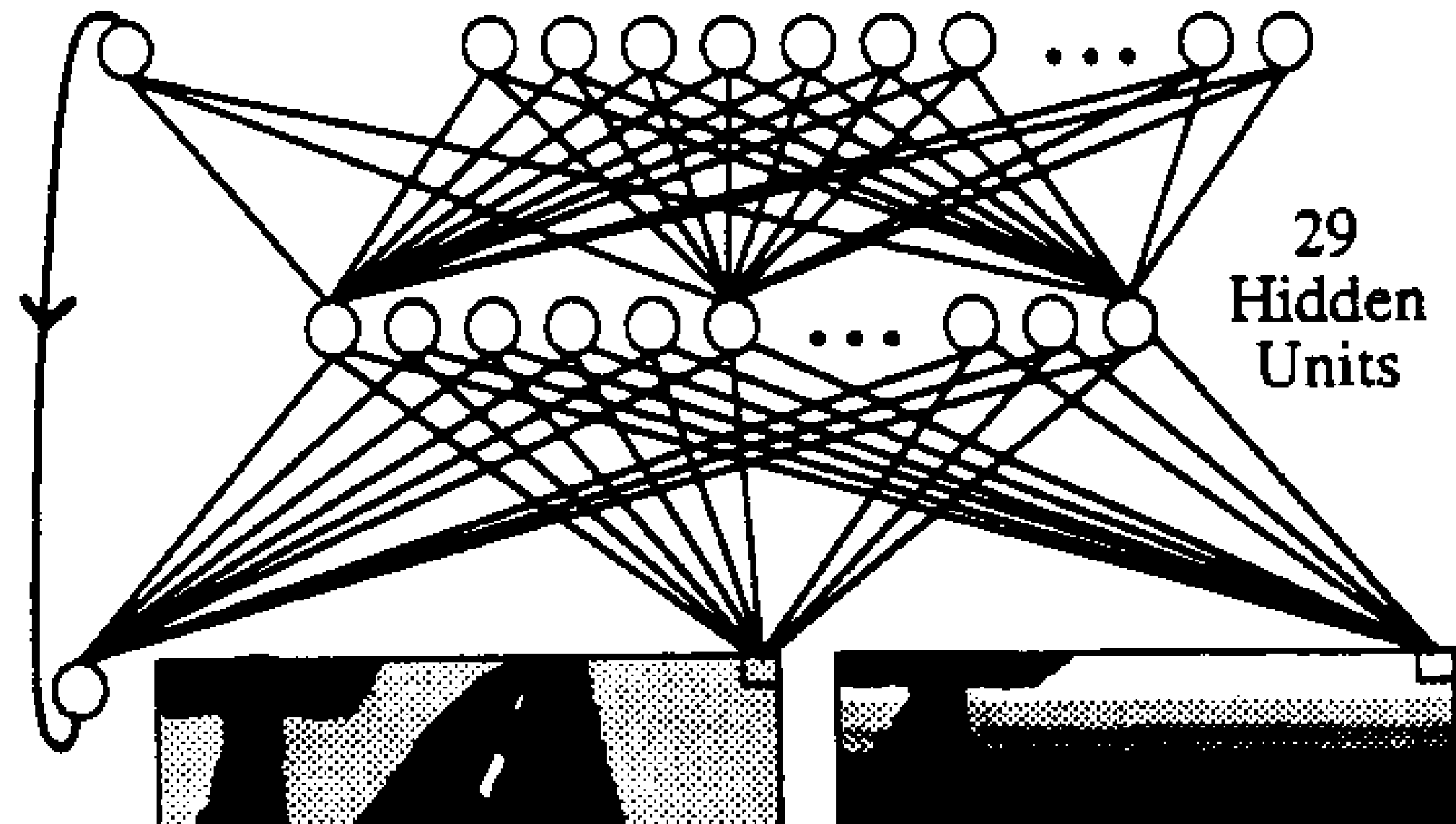


Corke
CSIRO
2002

ALVINN Architecture

Road Intensity
Feedback Unit

45 Direction
Output Units



30x32 Video
Input Retina

8x32 Range Finder
Input Retina

Figure 1: ALVINN Architecture

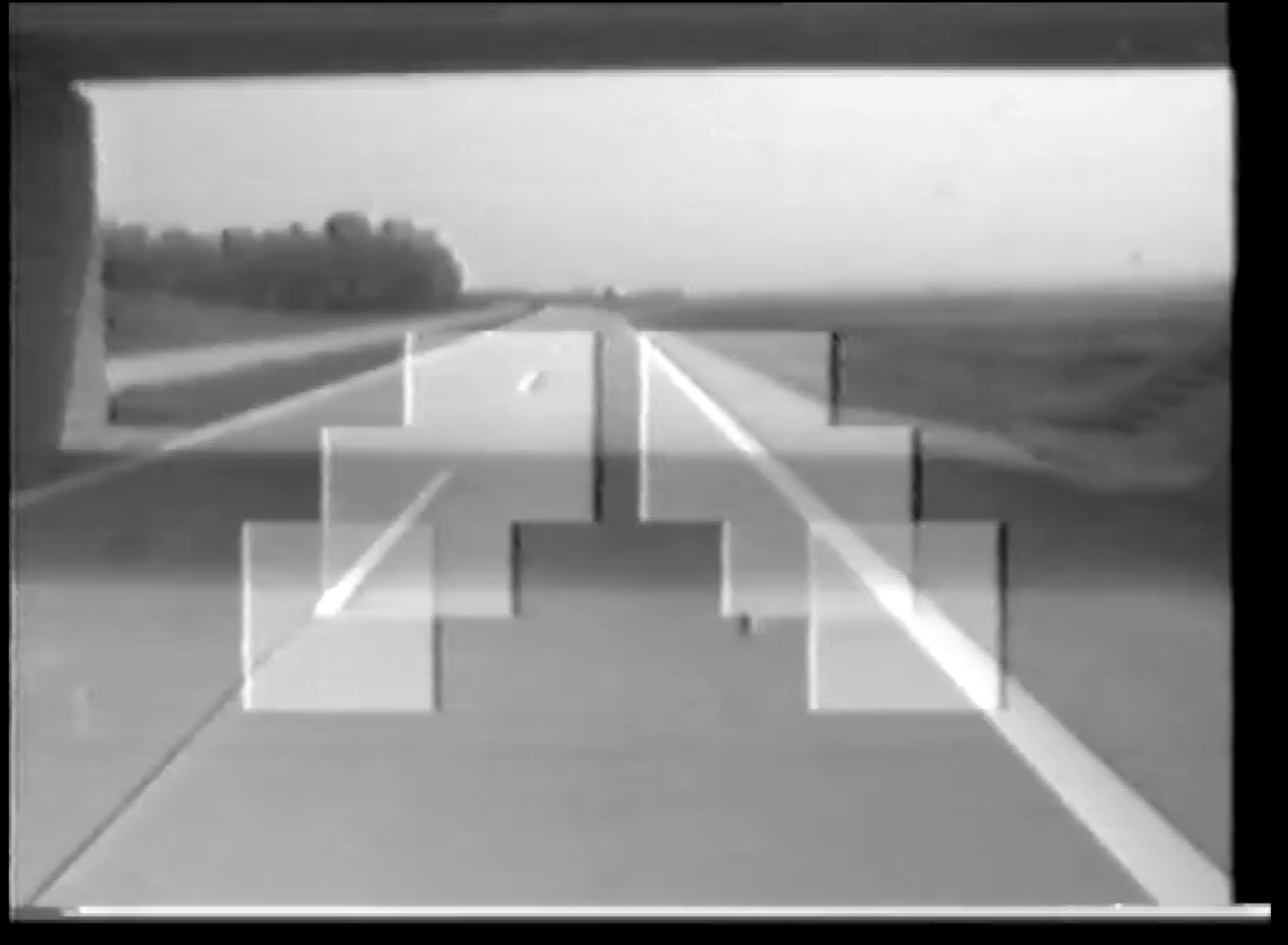


1989

3 x Sun computers
Warp systolic array computer (100 Mflop)



Ernst Dickmanns
Bundeswehr Universität München
1980-2004

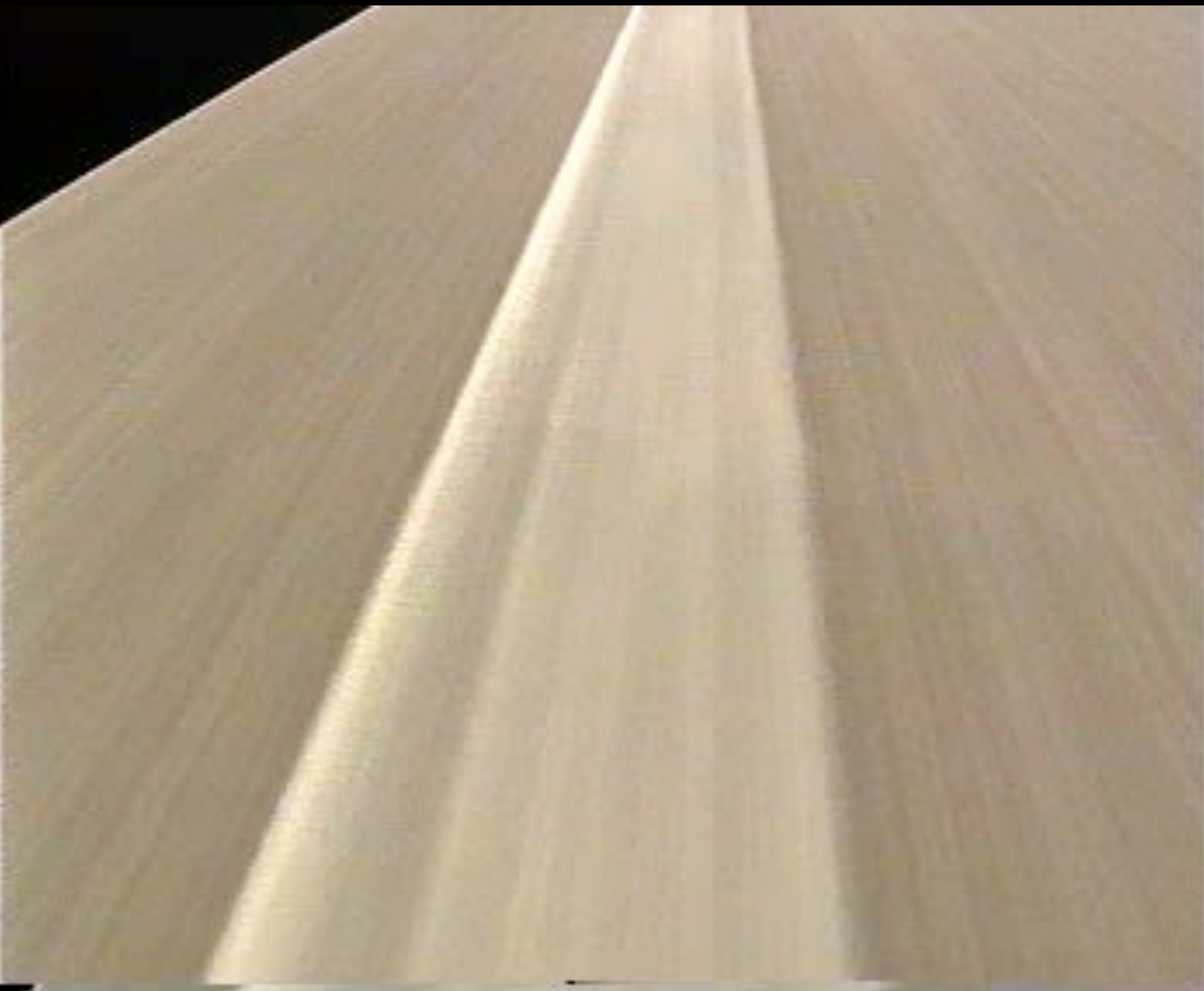


movie courtesy of Joe Wuesche
Bundeswehr Universität München



PROgraMme for a European Traffic of Highest Efficiency and Unprecedented Safety

1987-95



Munich to Copenhagen return
Upto 175km/h
Overtaking
Distance between human interventions:
- Avg 9km
- Max 158km

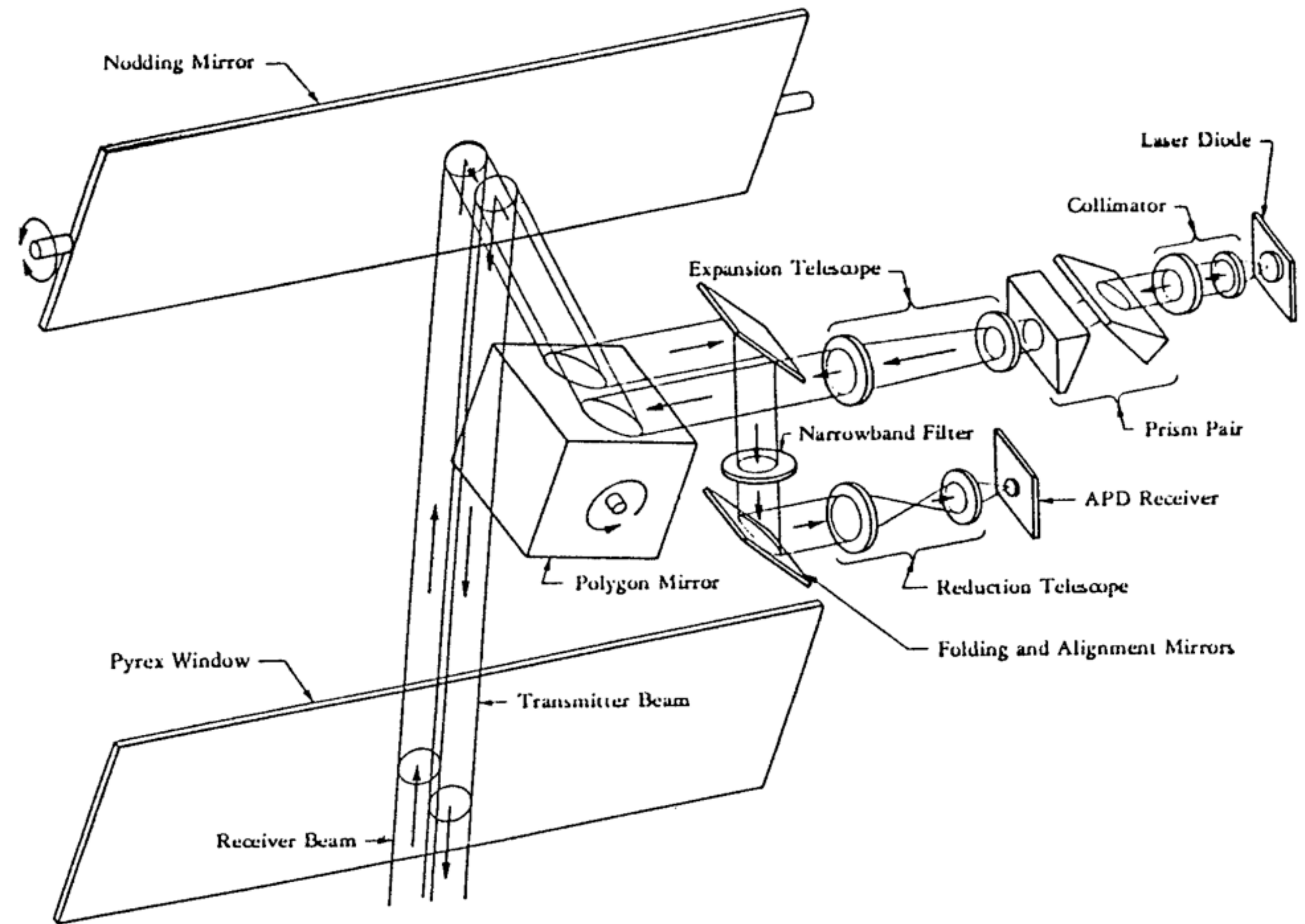
the split



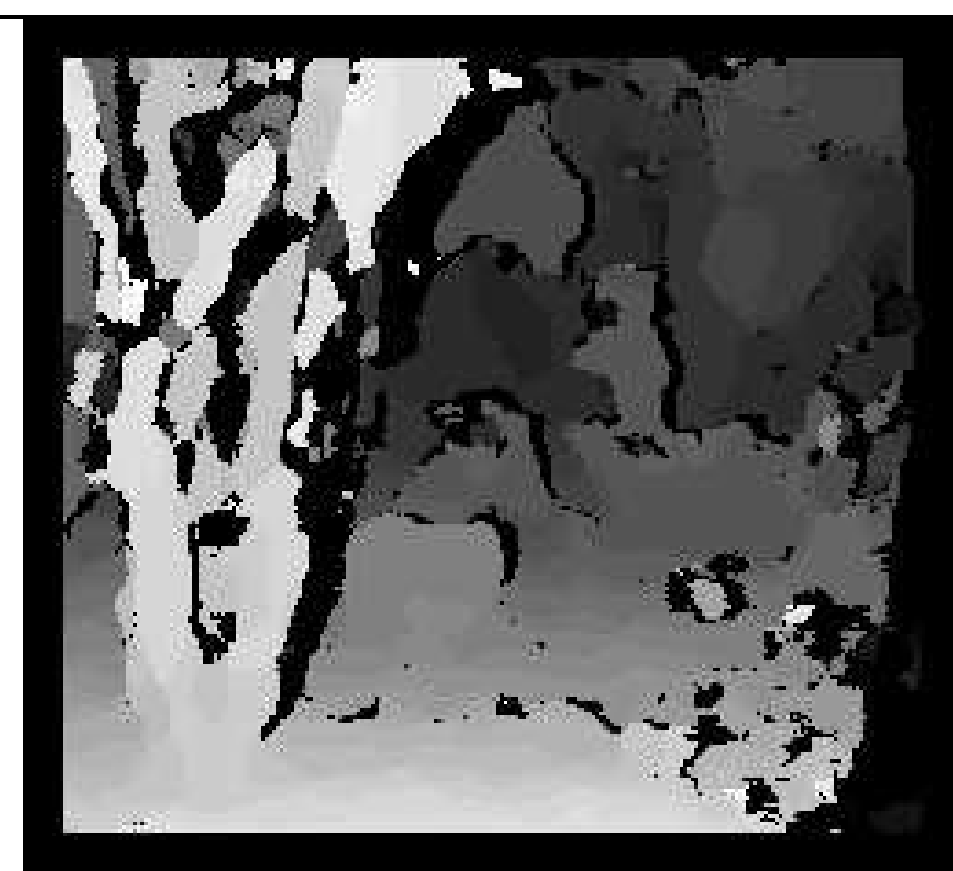
early 1990s



cost < \$10,000



CMU Navlab 1 + ERIM laser scanner
1980s





- Improved sensors appeared over time
- Increase in resolution, rotation rate, reflectance data quality, number of beams
- But maximum range and minimum cost little changed



Jericho, Oxford, August 2014.



RGB-D
cameras



but it's only geometry

we're missing color, texture, object
recognition, face recognition, human
actions, human intent etc.

what have the vision people been doing?



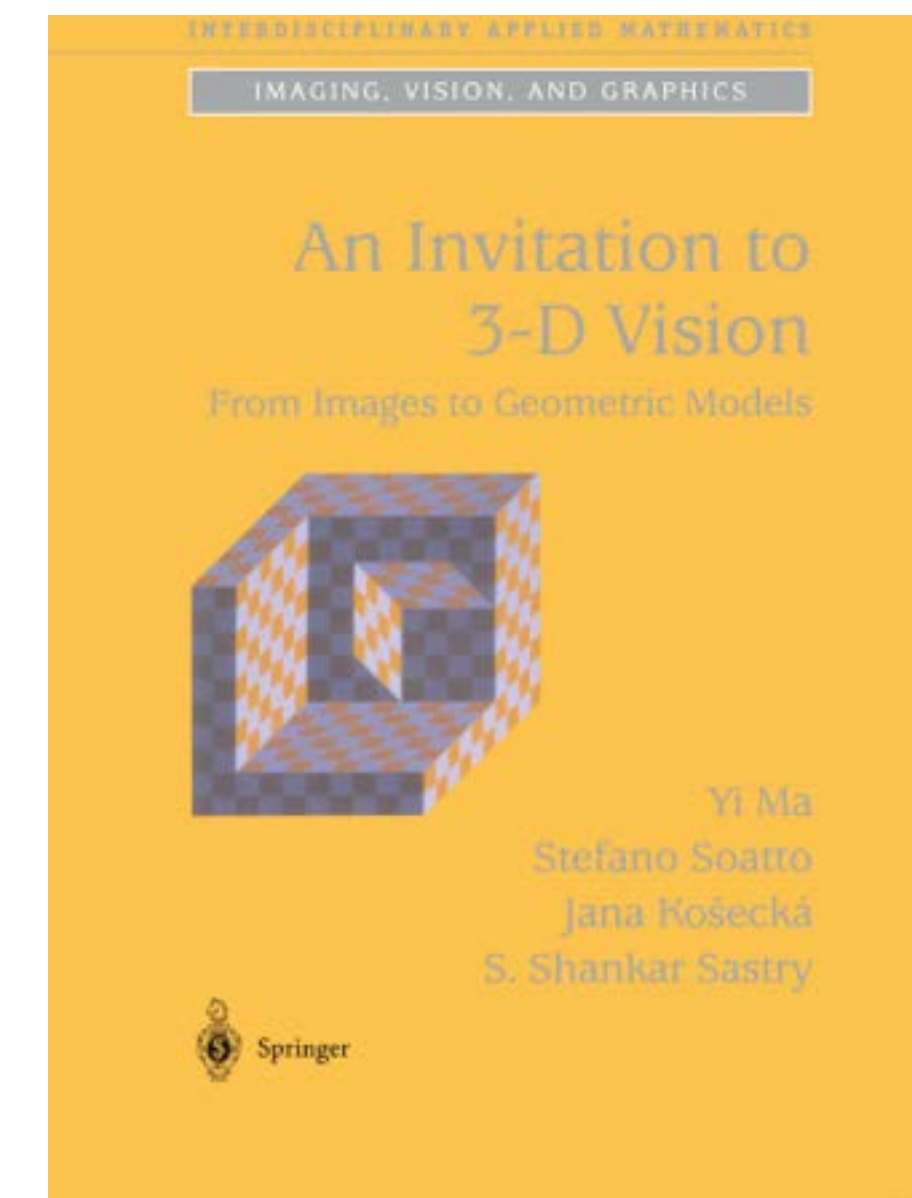
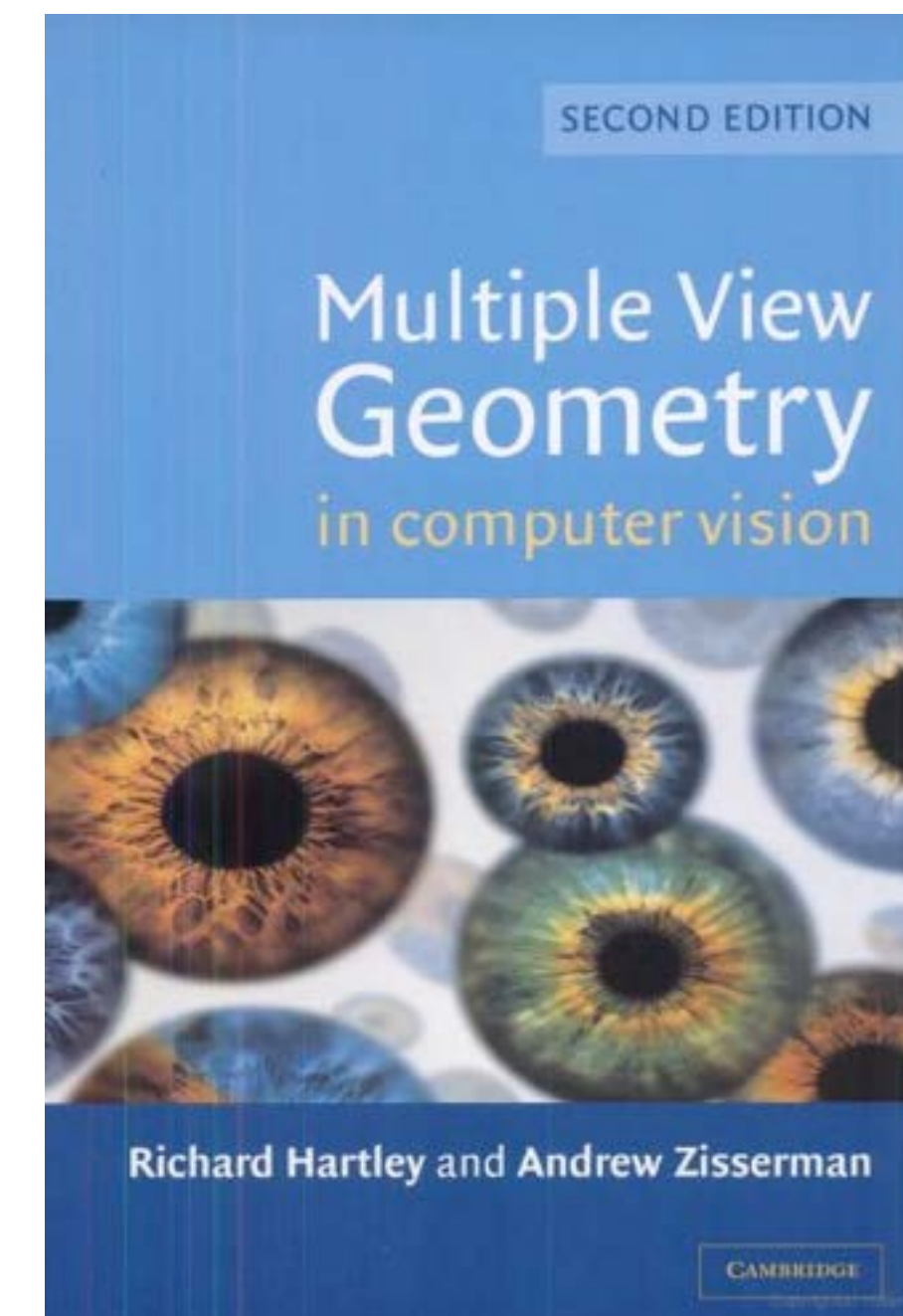
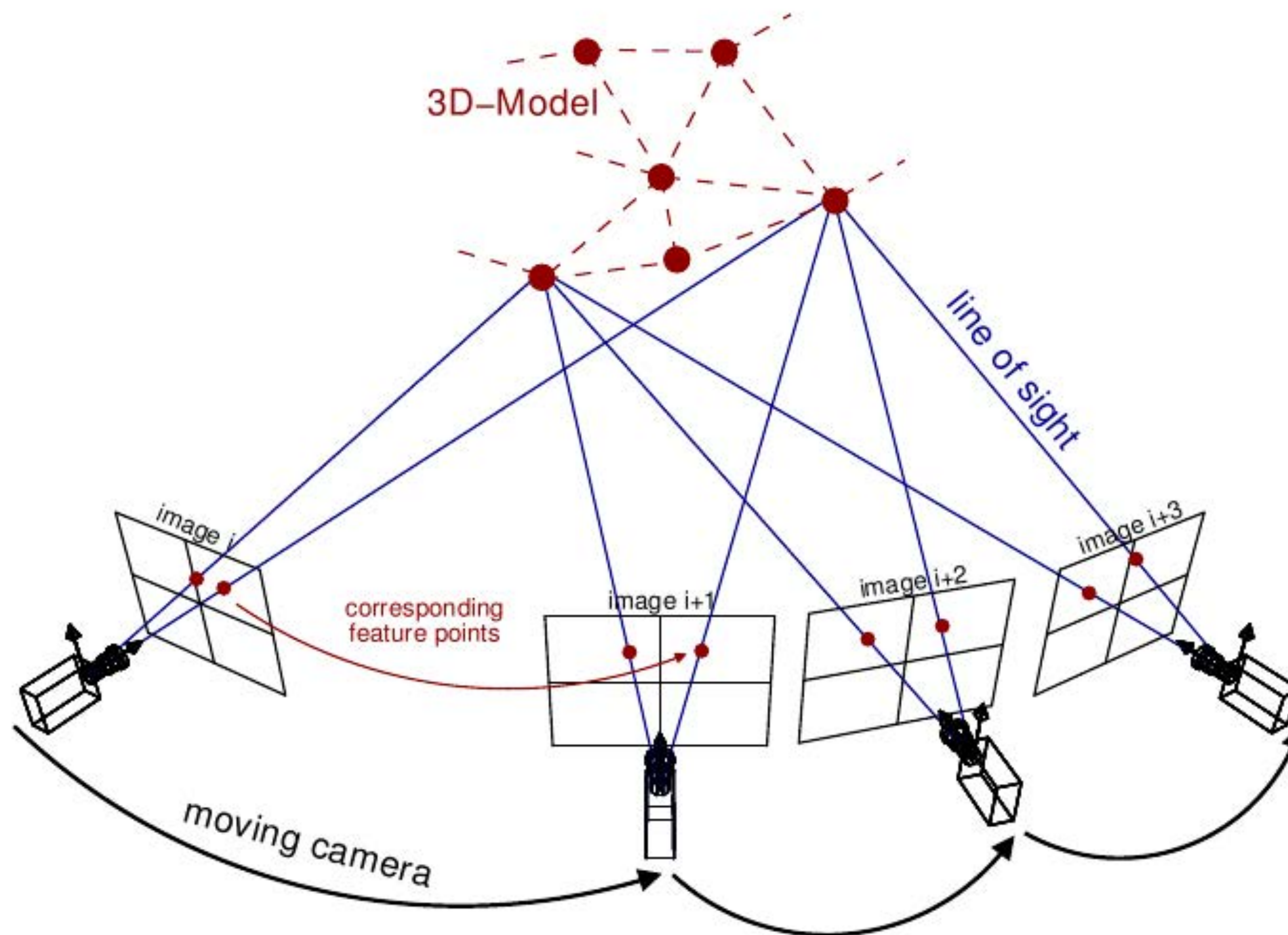
What makes Paris look like Paris?



Built Rome in a day

- Multi-view geometry
 - Structure from motion
 - Stereo
 - Visual odometry
 - Visual SLAM
- Stitching
- Super resolution
- Pose estimation
- Segmentation
- Image retrieval
- Object recognition
- Face recognition
- Action recognition
- Text recognition
- Pedestrian detection
- Calibration
- Feature detectors & descriptors: SIFT, SURF, FAST, BRIEF, BRISK, MSER, FREAK, HOG, CenSureE







Middlebury stereo dataset



Stereo reconstruction





Monocular camera reconstruction of Scott Reef
He, McKinnon, Upcroft
QUT

image sequence courtesy of U. Sydney



Alessandro Bissacco, Mark Cummins, Yuval Netzer, and Hartmut Neven.

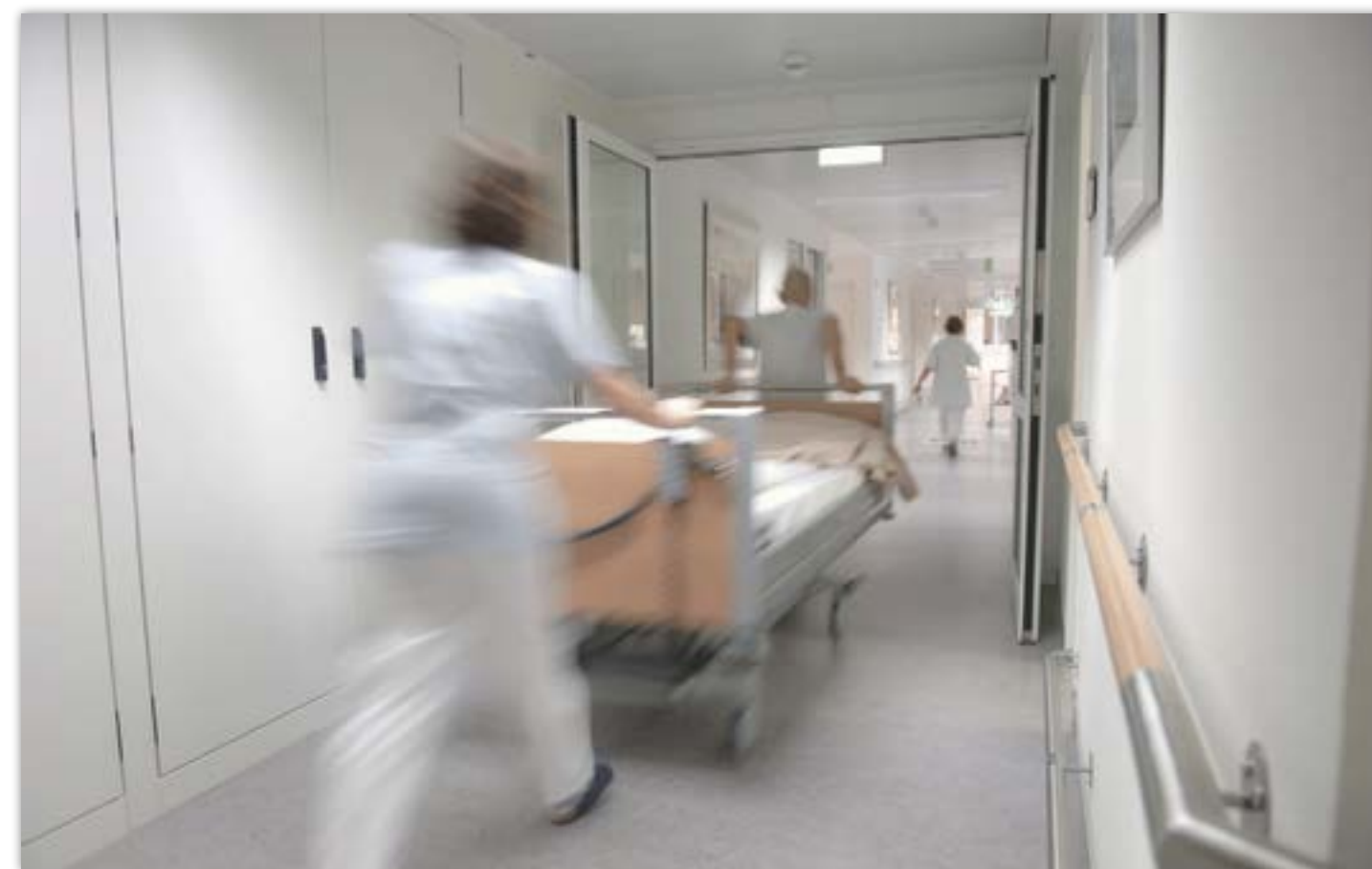
Photoocr: Reading text in uncontrolled conditions. In Computer Vision (ICCV), 2013 IEEE International Conference on, pages 785–792. IEEE, 2013.

lots of awesome (and useful)
stuff roboticists need to be
aware of

but the camera is
generally passive

robots → computer vision

the future



where are all
the robots?



Seeing

recognising objects & stuff

recognising places

detecting motion

paying attention

recognizing humans, their activities and intent

seeing creates memories
memory helps seeing

context for seeing
seeing for context

see to move
move to see

So why is it hard?

- Vision is a great sensor but
 - The rich visual information is encoded
 - 3D world is projected to 2D
 - To recover the “lost” information we need
 - assumptions (context, world knowledge)
 - to interact with the environment (move)
 - Many distractors
 - shadows, lighting change, seasons etc.
 - information is ambiguous

but we know it is possible

- 120Mpixel
- 20bit dynamic range
- 3 colors

× 2

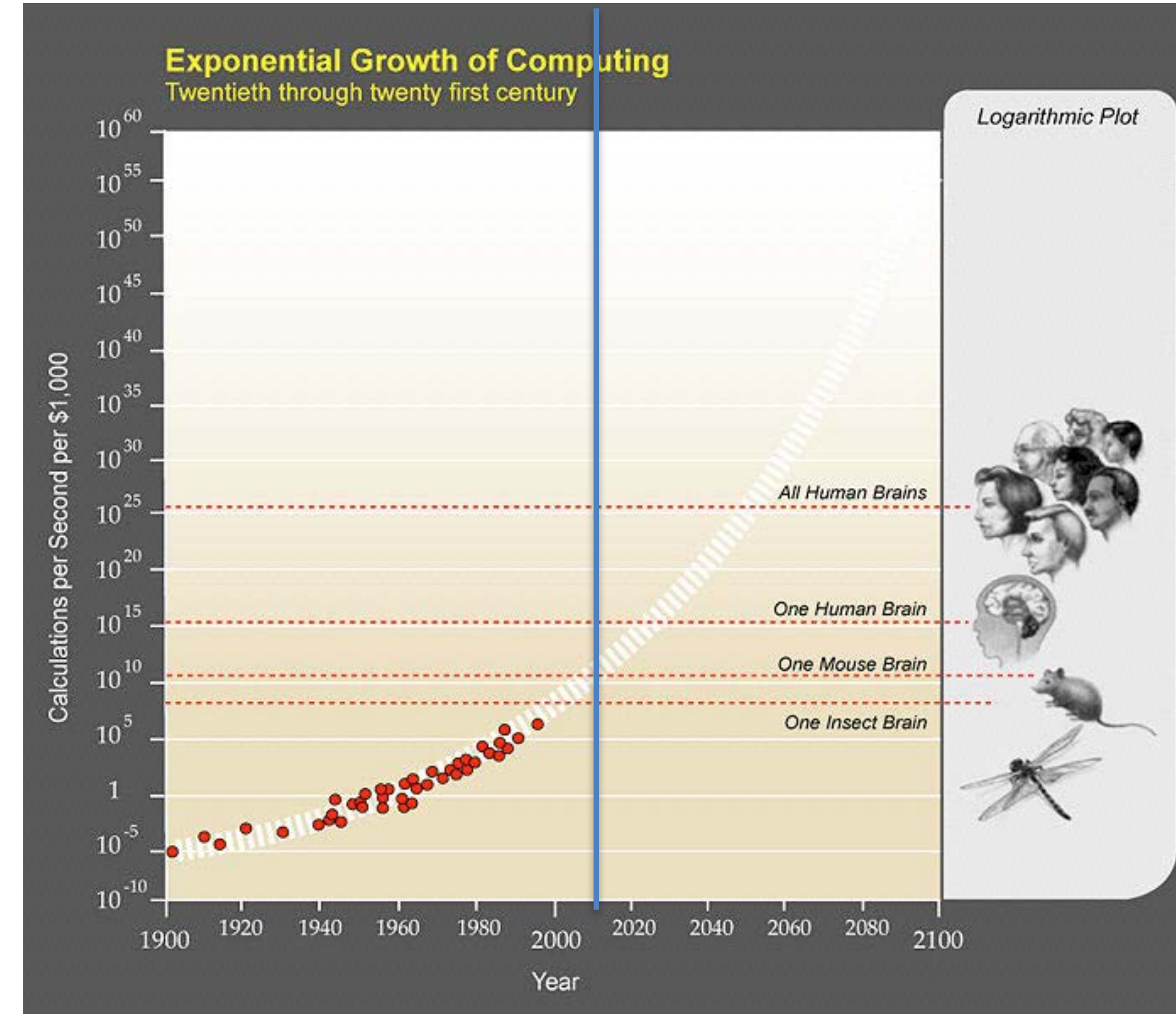
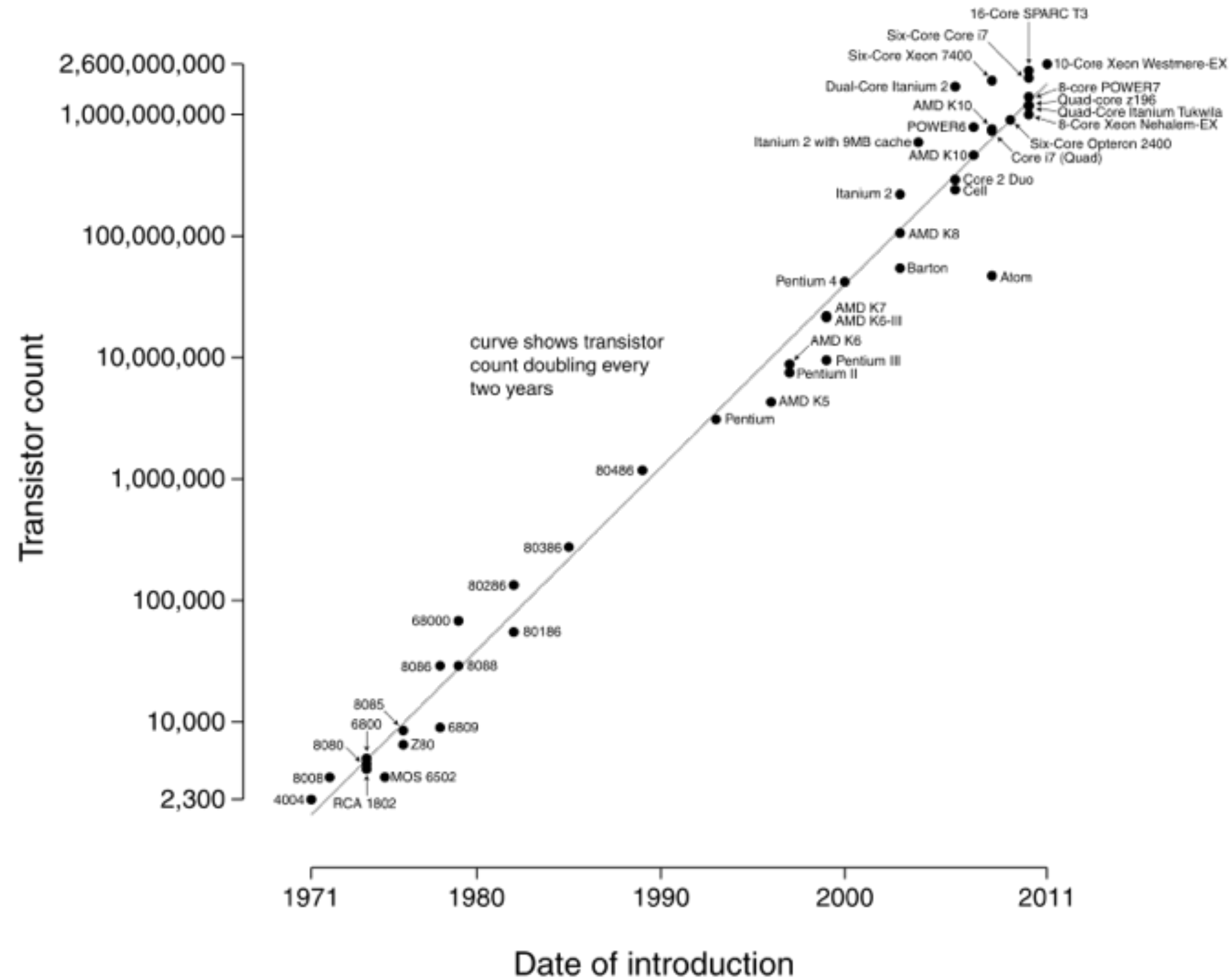
- 3 gyroscopes
- 2 accelerometers

× 2

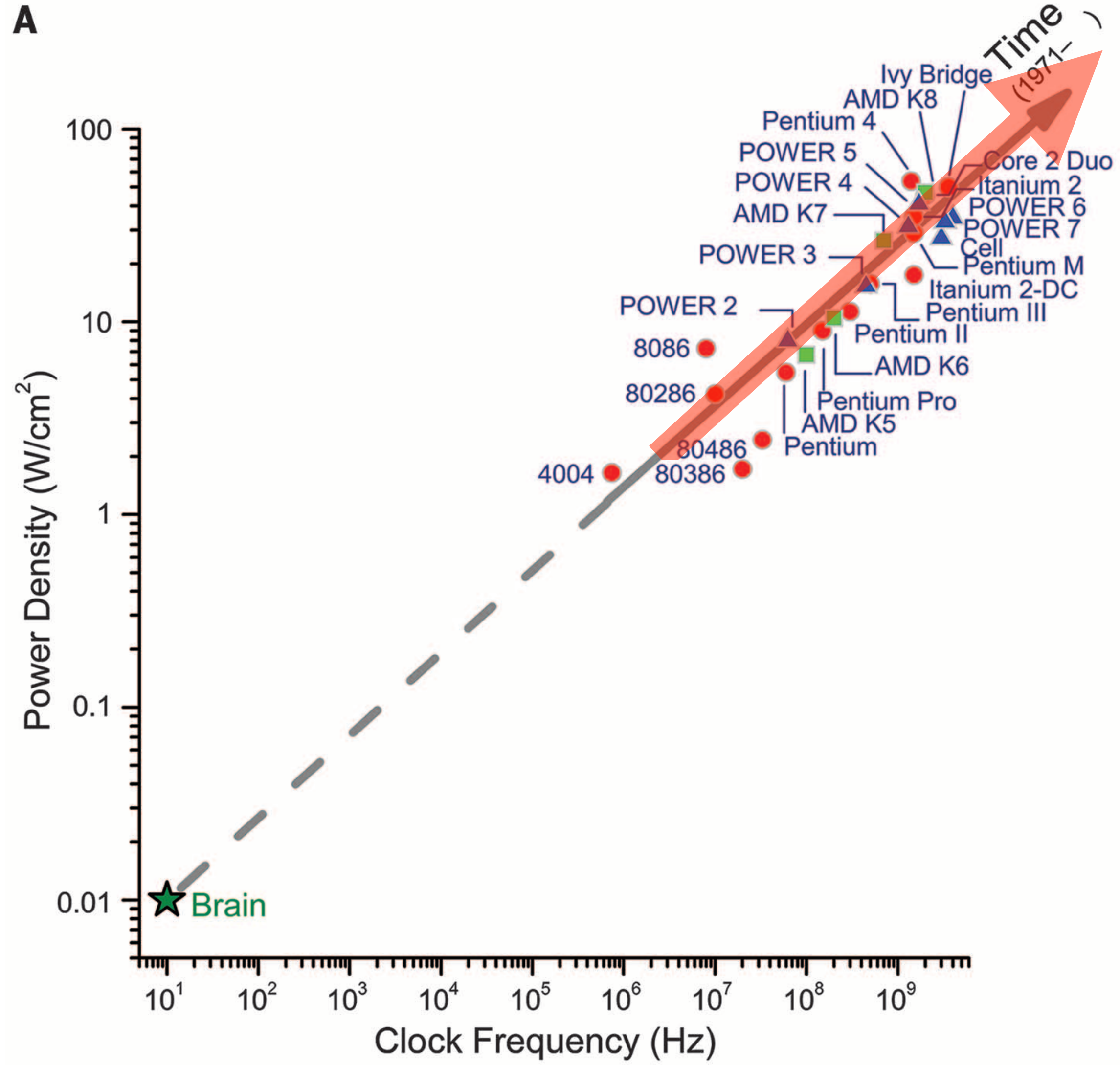
- Vision engine
 - 3×10^{10} neurons
 - 500g
 - 6W

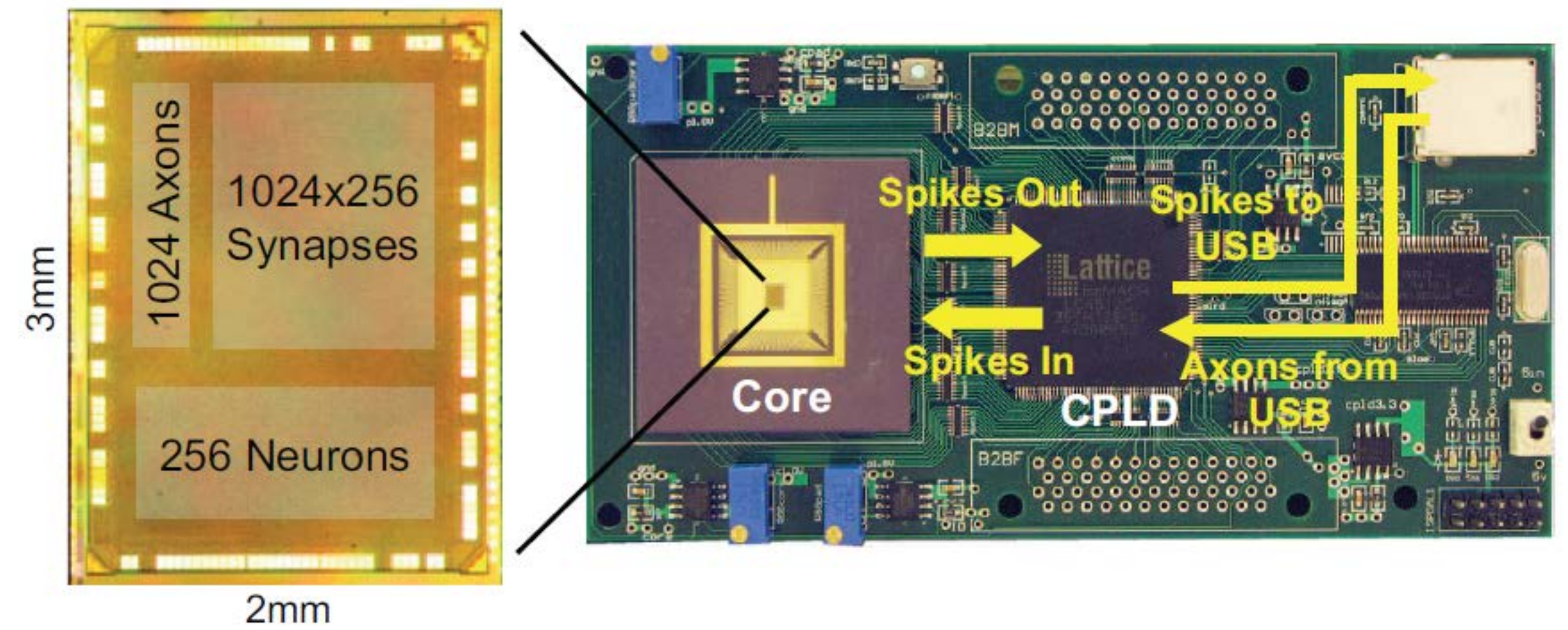
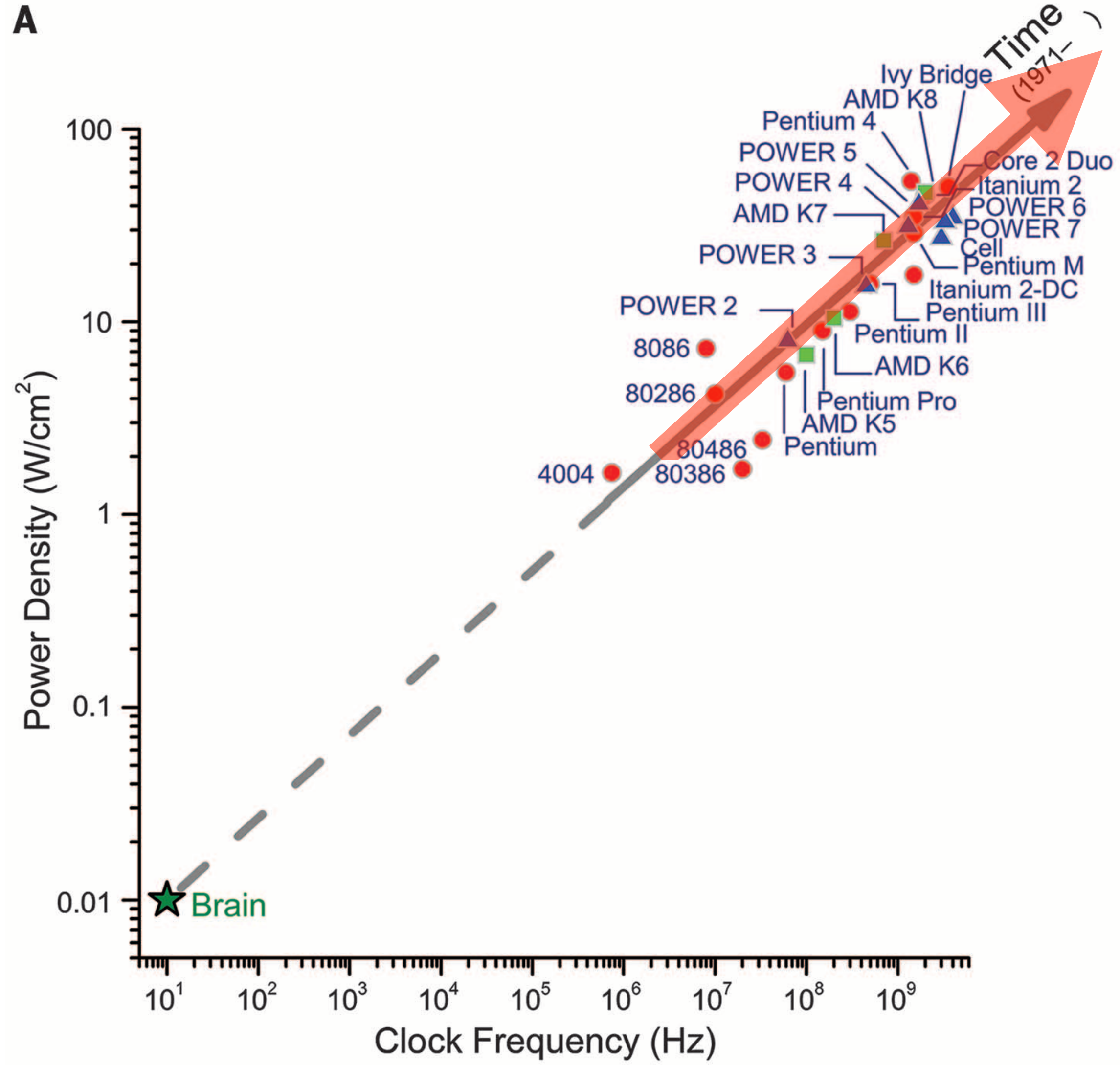


Microprocessor Transistor Counts 1971-2011 & Moore's Law



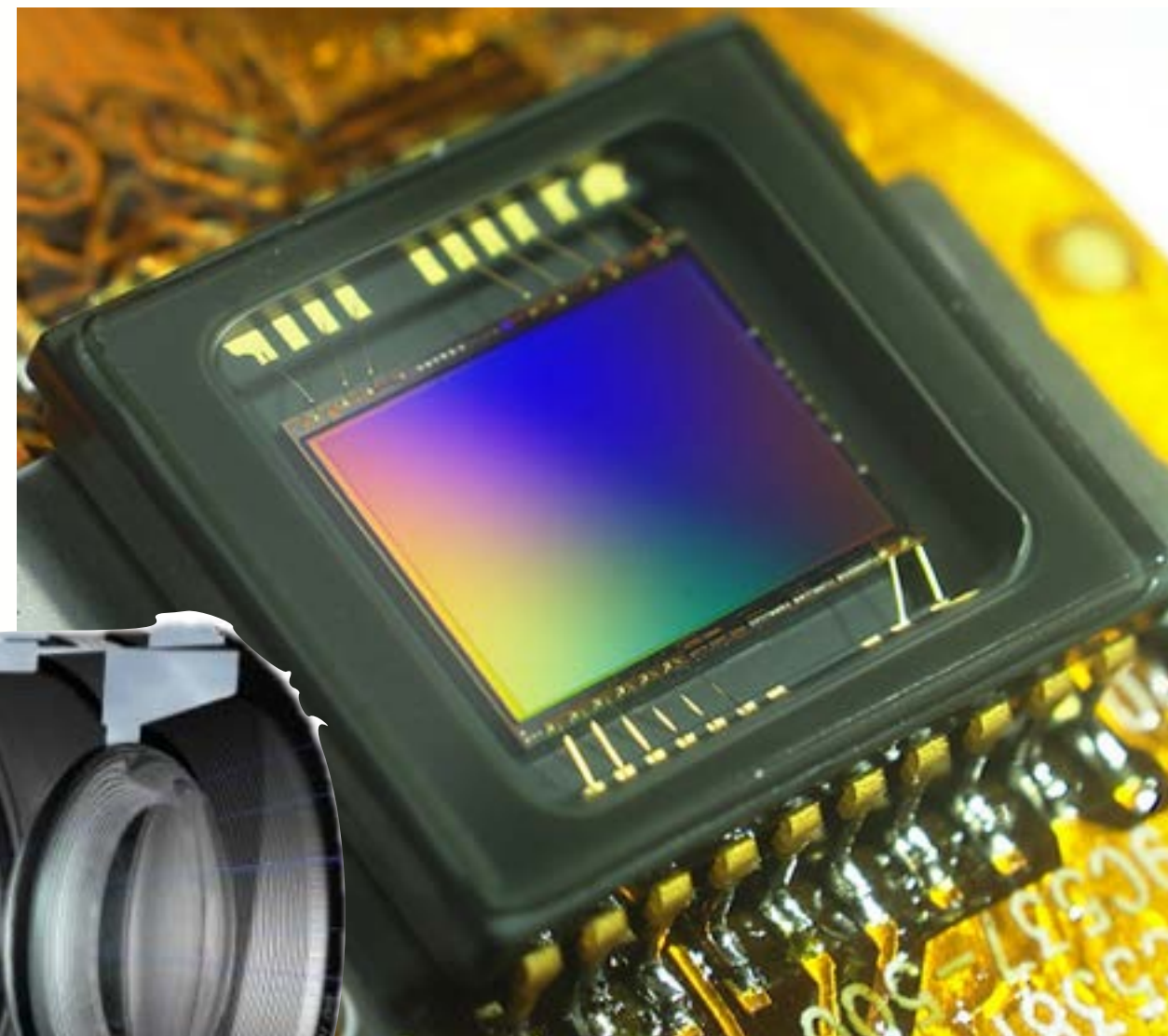
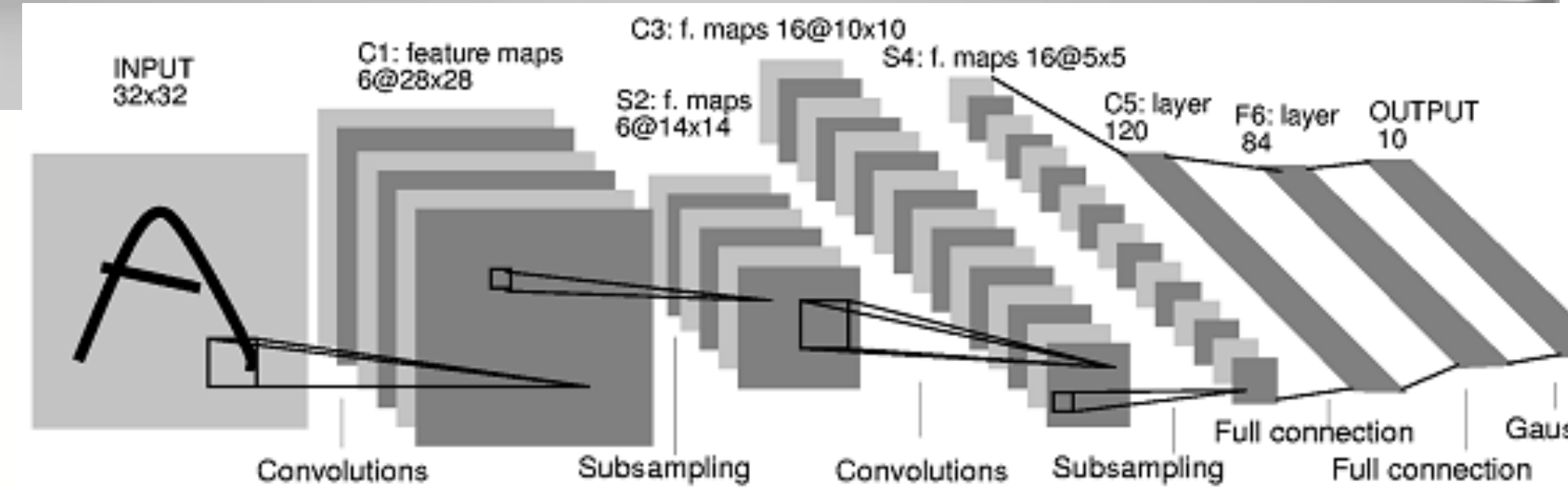
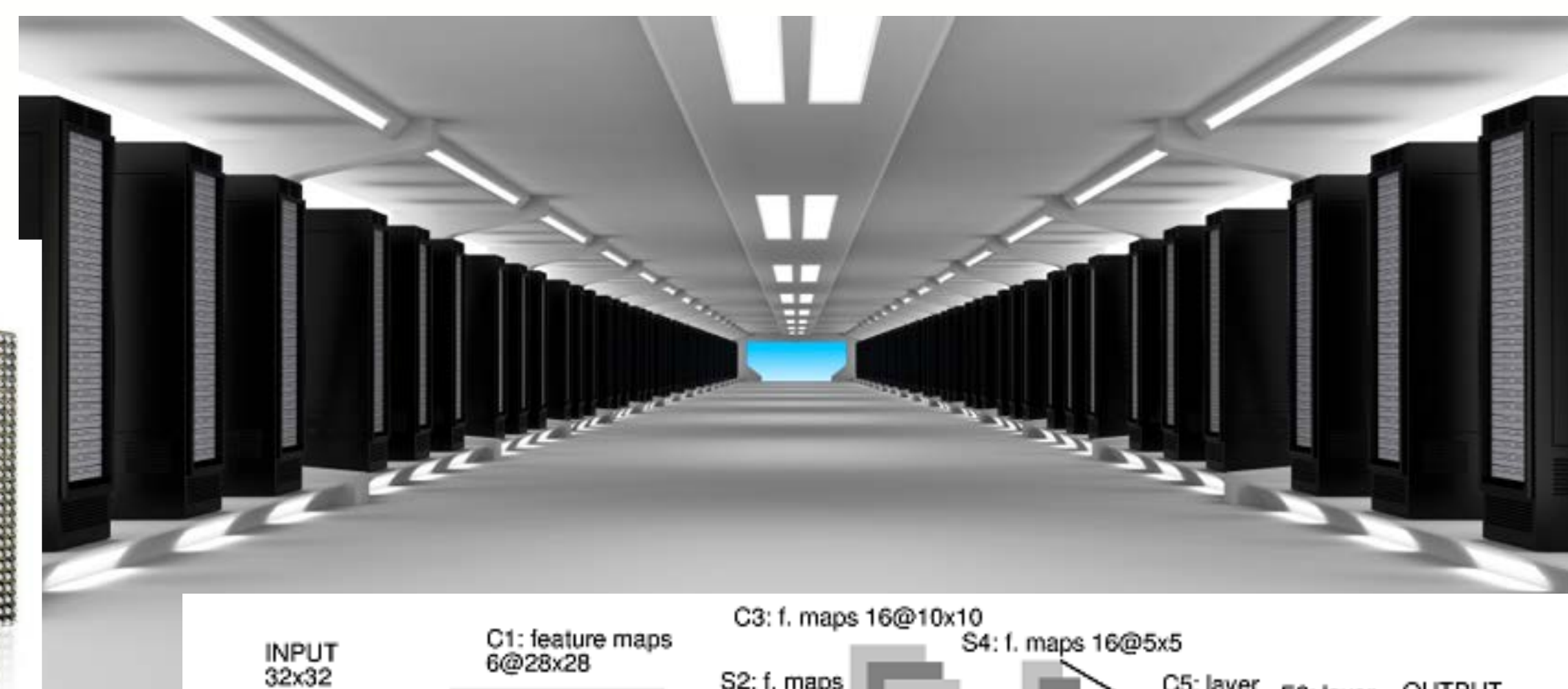
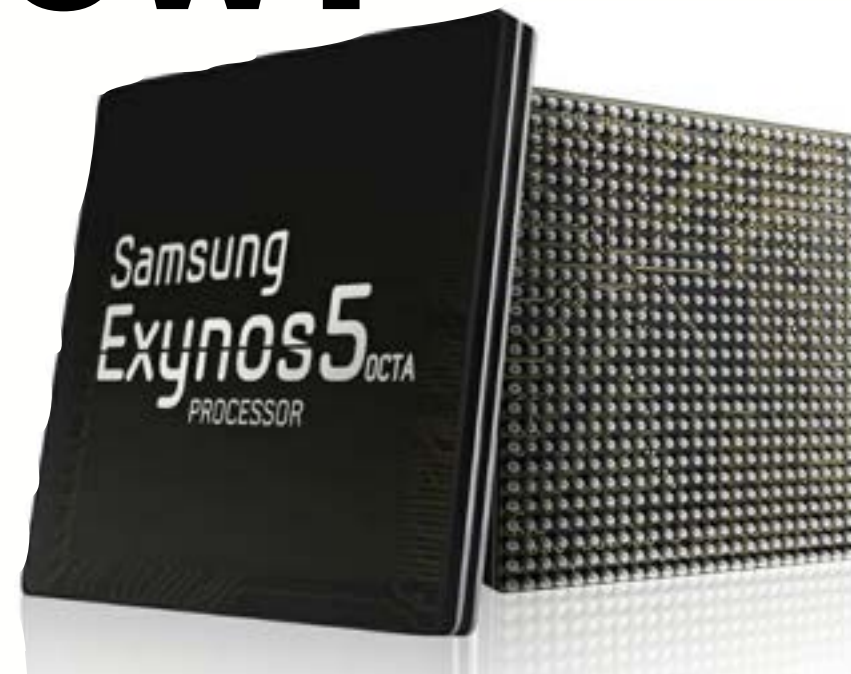
Ray Kurzweil

A

A

What's different now?

- Computation (ops/\$/W)
- Algorithms
 - all that computer vision stuff, particularly semantics
 - machine learning, CNNs etc.
- (Visual) neuroscience
- Sensors
 - light field
 - low light
 - high dynamic range



An **ambition** for a sensor

- Cost < \$1,000
- Works in a usefully wide range of lighting conditions
- Provides geometric and semantic description of all objects in the scene
 - particularly people: activity, intent, etc.
- Low power
- Learns, exploits attention & context

ARC Centre of Excellence for **ROBOTIC VISION**

\$25M
7 years
13 CIs

16 new postdocs
50 PhD students
8 projects

Robust vision



- better sensors, comp. photography
- robust algorithms for poor images
- contextual priming

Semantic vision



- understanding from images
- lifelong learning

Vision & action



- seeing to move
- moving to see

Algorithms & architectures



- real-time, energy efficient
- new architectures
- local + cloud computing



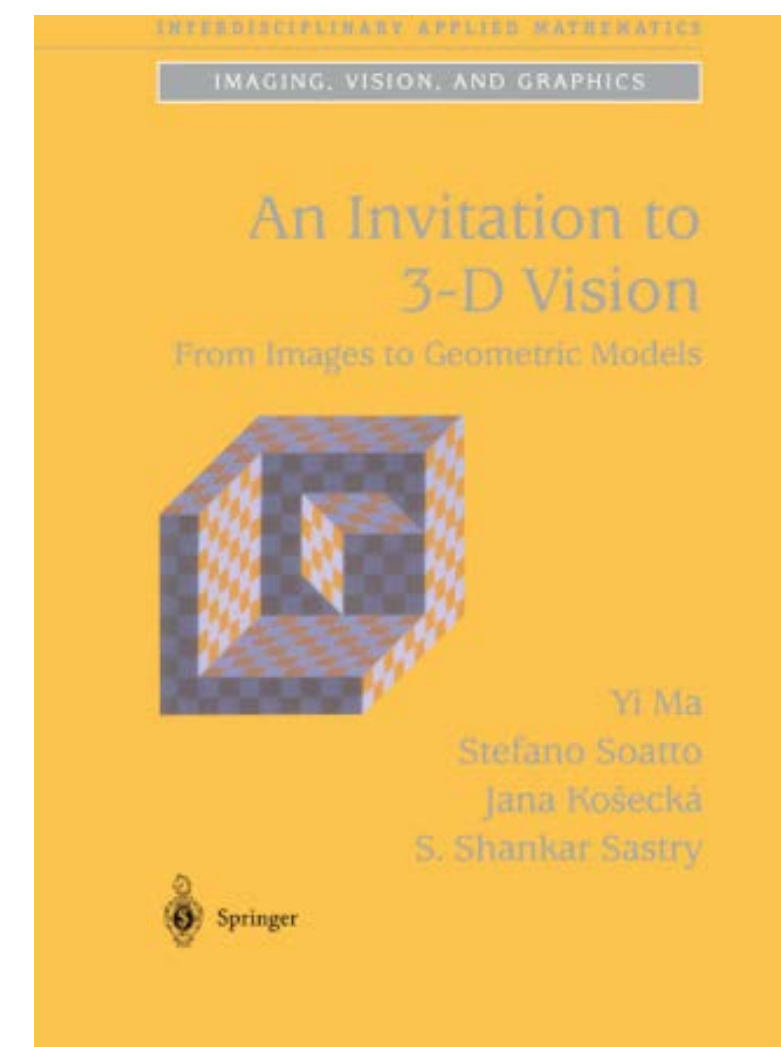
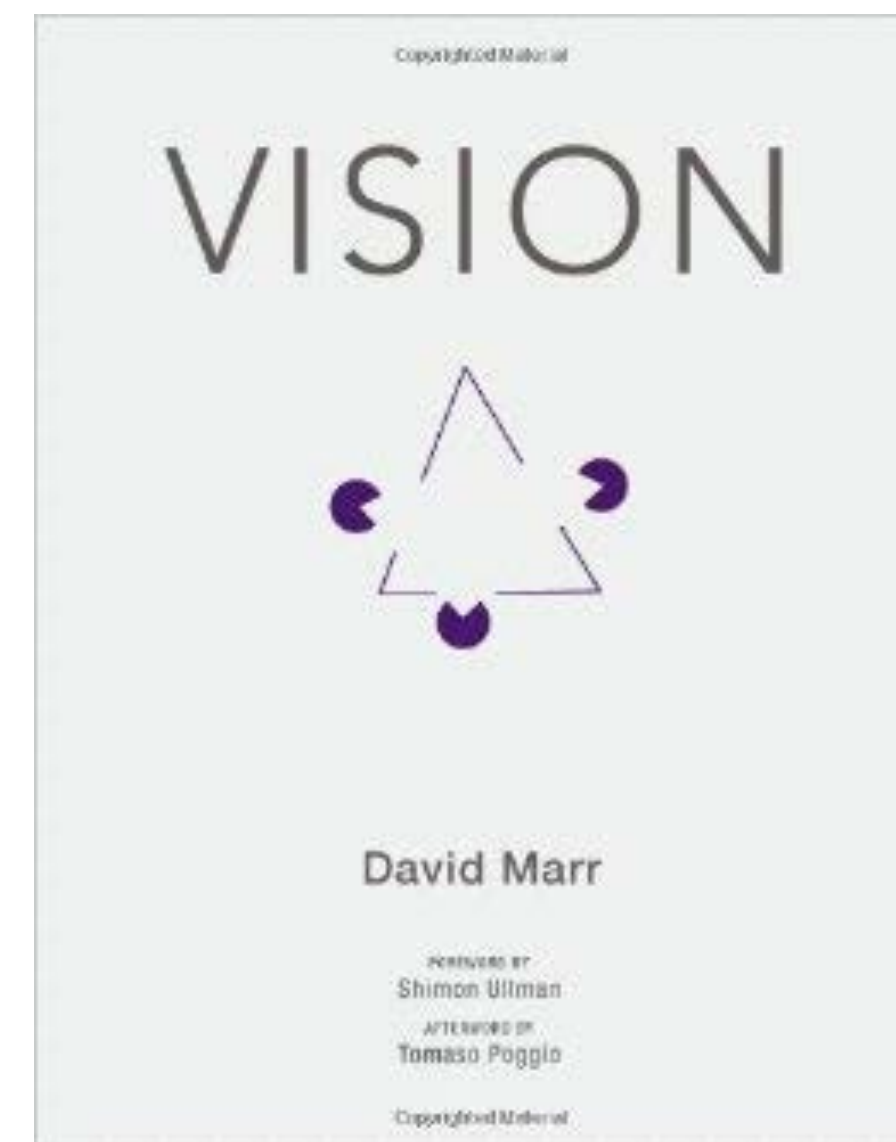
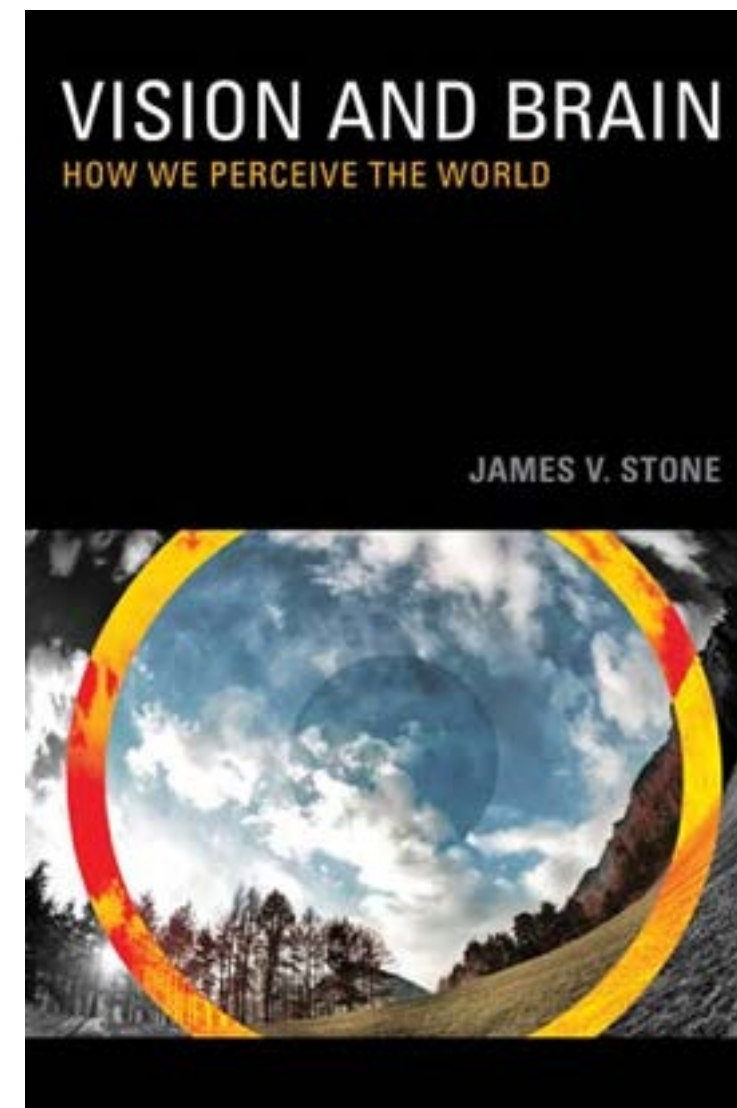
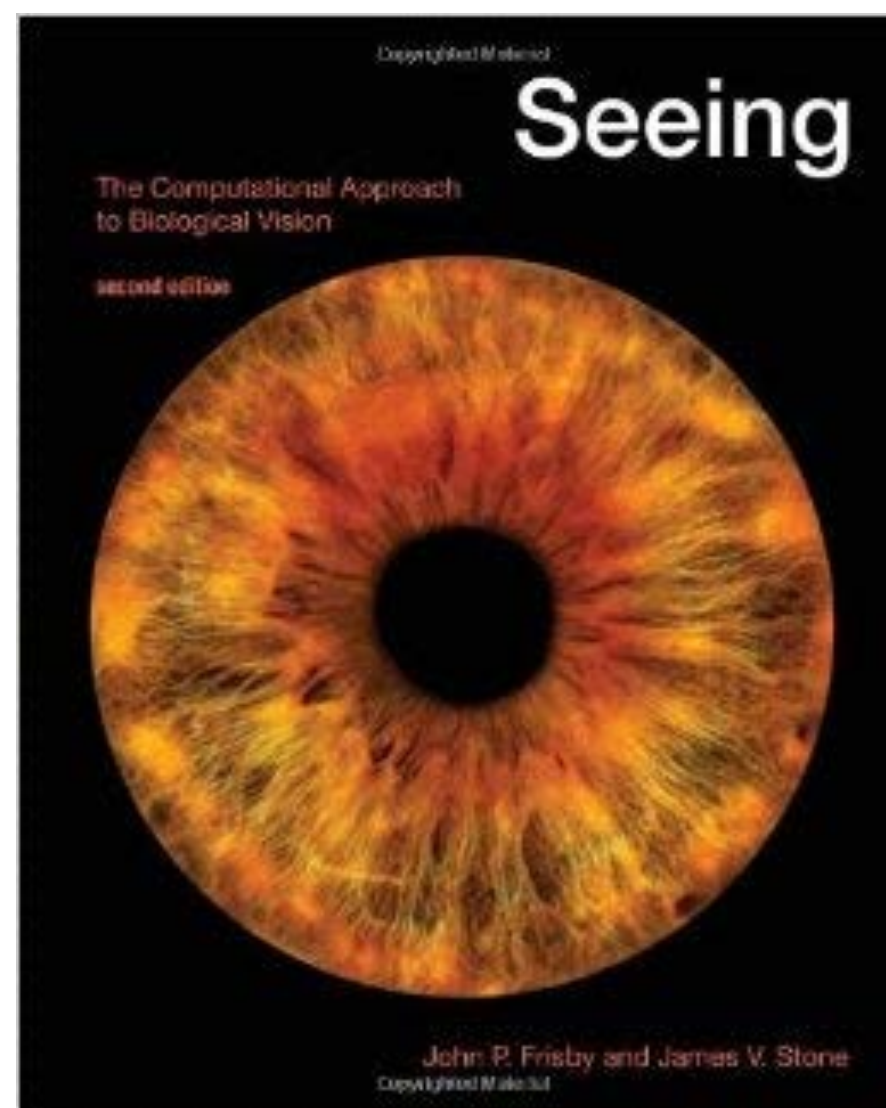
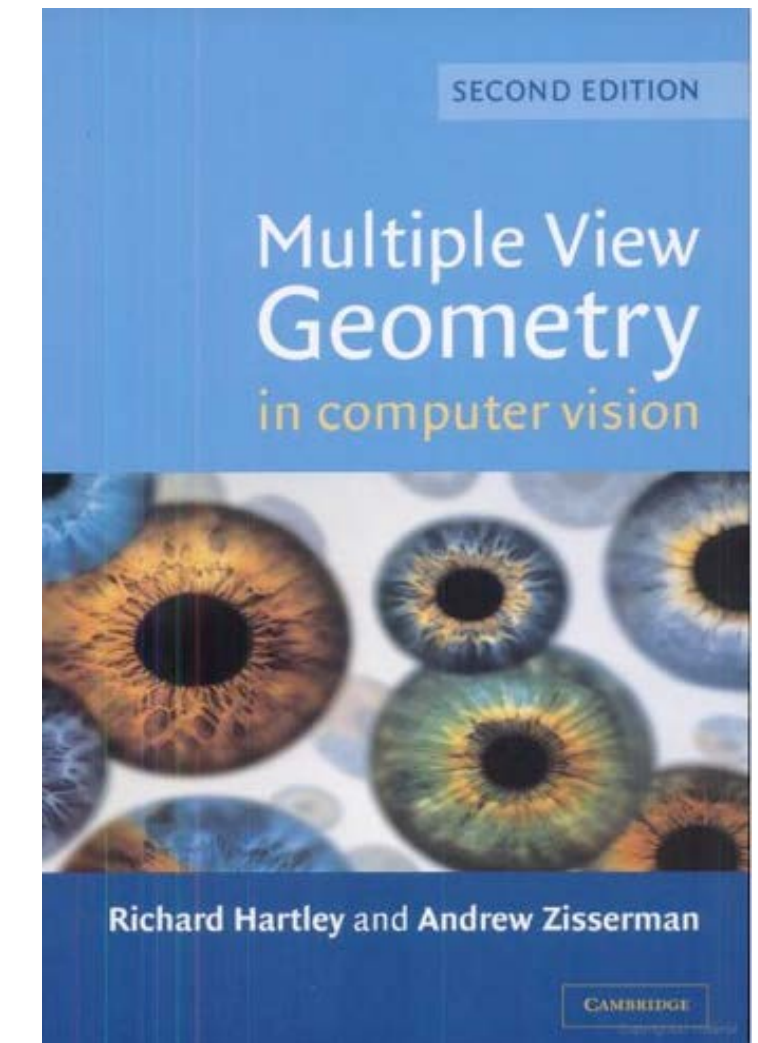
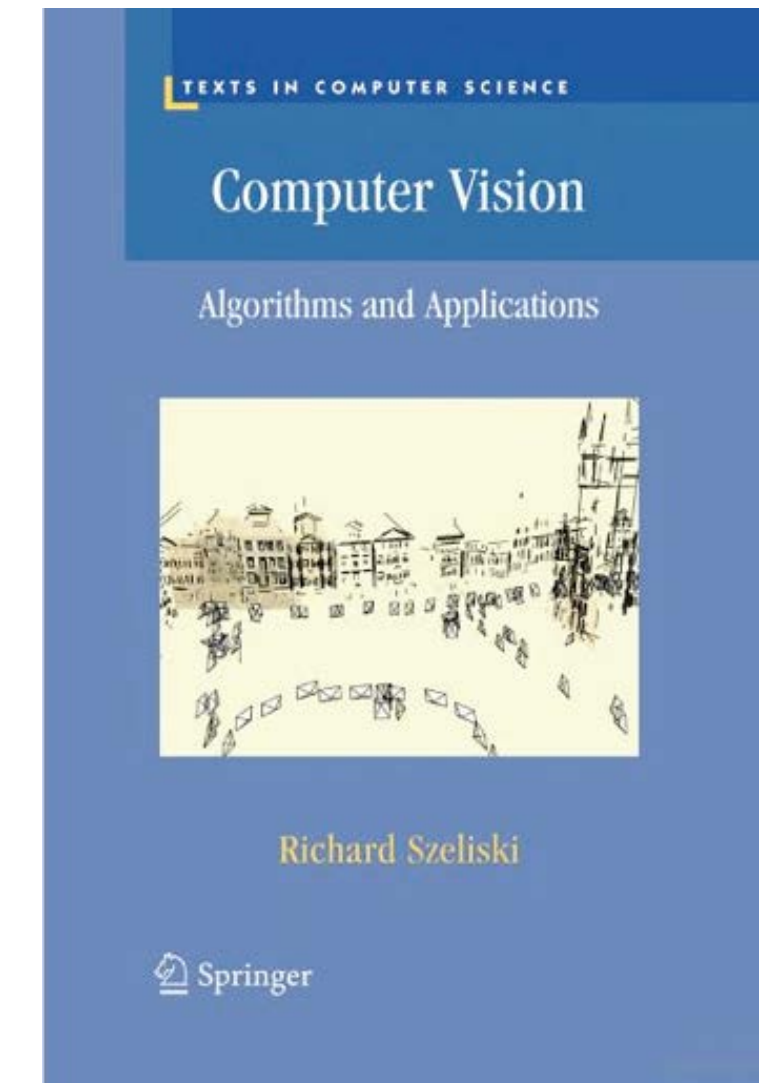
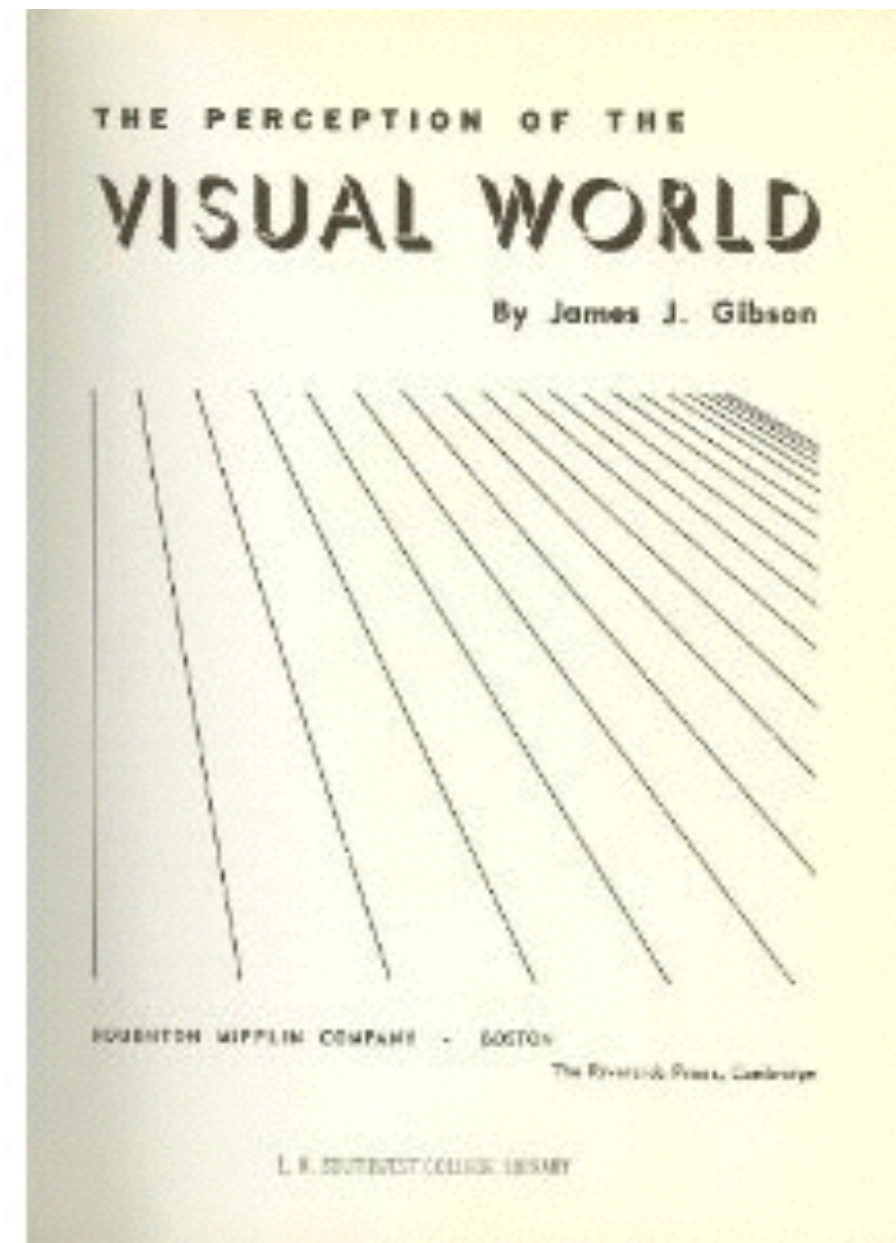
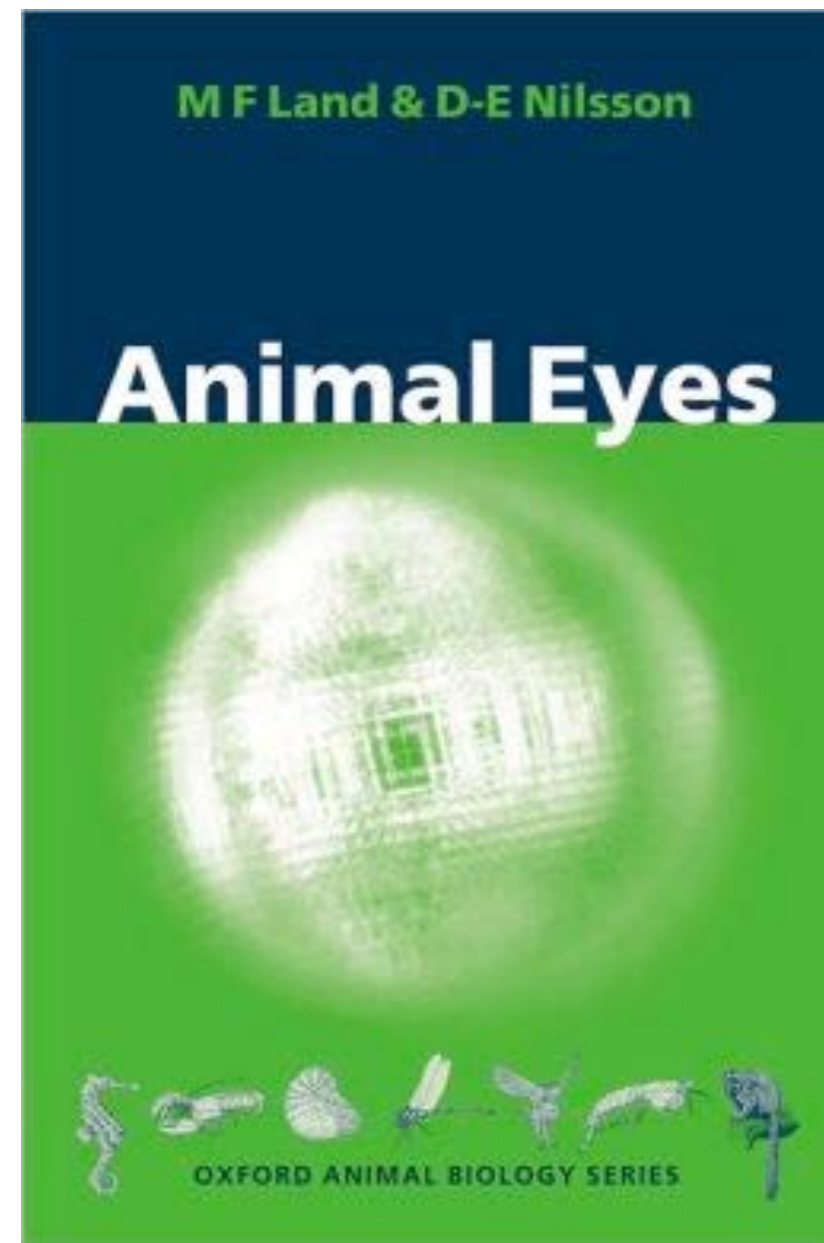
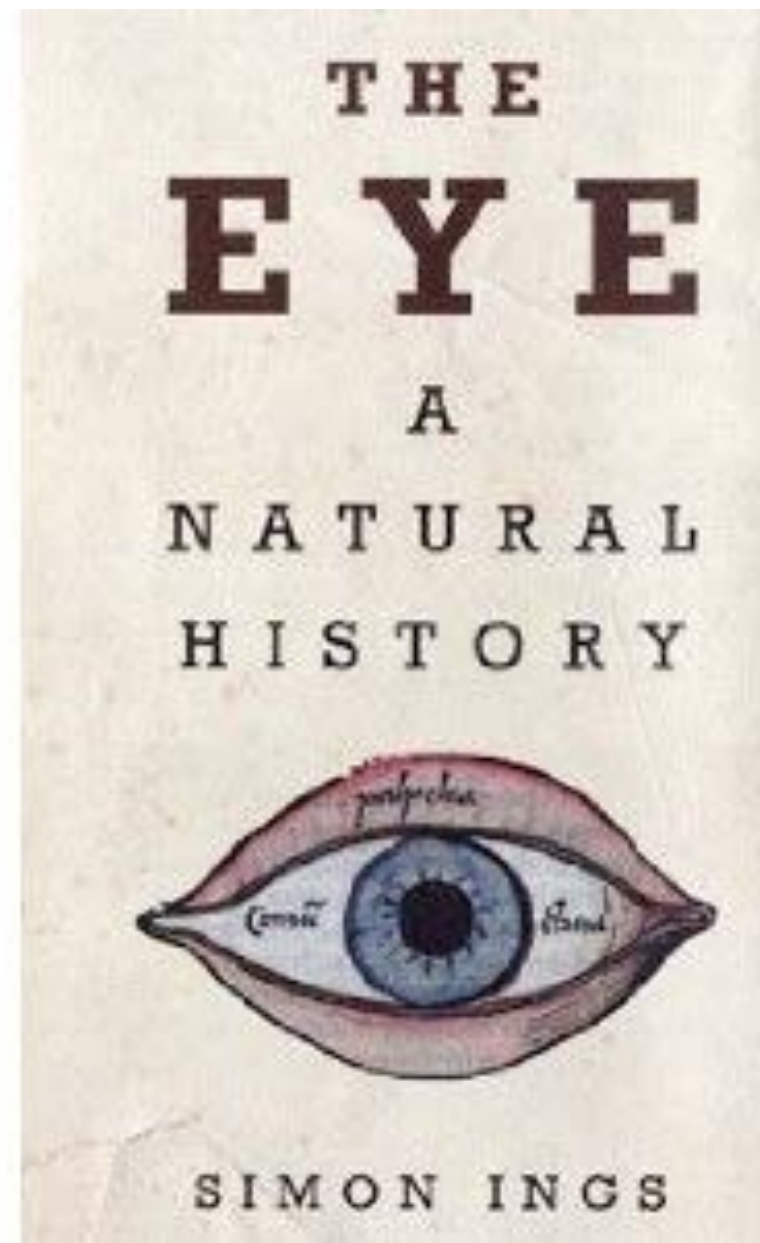
Vision is a **540M** year old
technology

Radically transformed life on planet Earth

It is sufficient to enable almost all useful tasks

It's time to bring vision and
robotics back together!

Further reading



Vision is a **540M** year old
technology

Radically transformed life on planet Earth

It is sufficient to enable almost all useful tasks

It's time to bring vision and
robotics back together!





Cave Paintings ~40,000 years ago



Ideal City (1470)

Piero della Francesca (1415–1492)

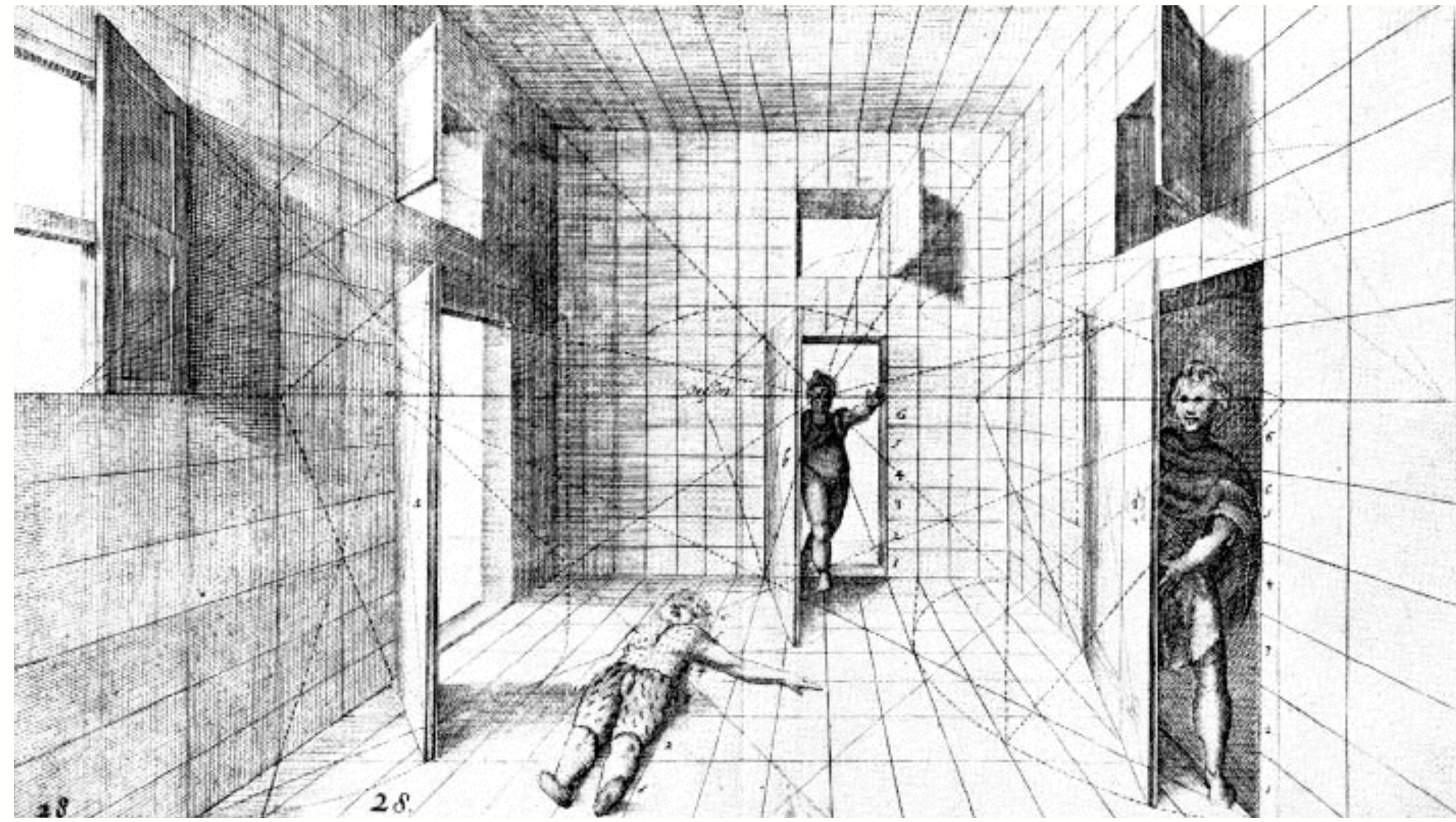


Figure 28 (Jan Vredeman de Vries, 1604).

Used with permission from Perspective, Dover Publications, 1964.



People are actually avoiding walking in the "hole" 2007
Joe Beever | CC A2.0



Stunning 3D chalk drawing from Zebit stops Liverpool shoppers in their tracks on Bold Street. 2012

Bill Hunt Original art: Zebit | CC A2.0



Edgar Meuller <http://www.metanamorph.com>
Edgar Mueller | CC-BY-SA-3.0, via Wikimedia Commons