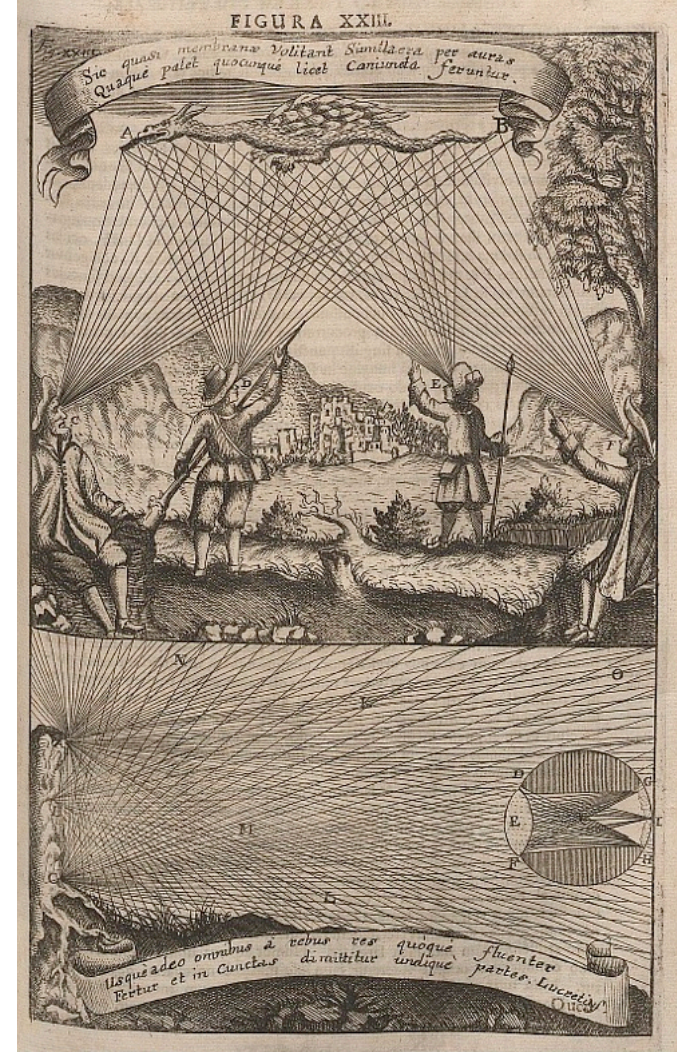


CS-503 Visual Intelligence: Machines and Minds

Amir Zamir

Lecture 1

Vision: understanding the world through light.



Vision: understanding the world through light.

This Course:

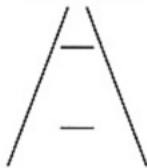
- **Basics:** biological vision, theories, etc.

Which horizontal line is longer?



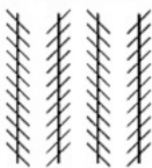
A

Which horizontal line is longer?



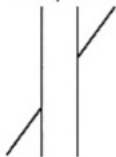
B

Are the long lines parallel or tilted?



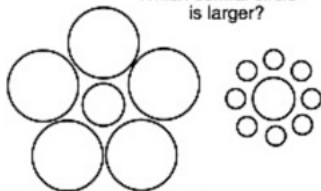
C

Do the diagonal lines line up or not?



D

Which central circle is larger?



E

- **Foundation Models & Modern Tools:** visual and multimodal models, etc.



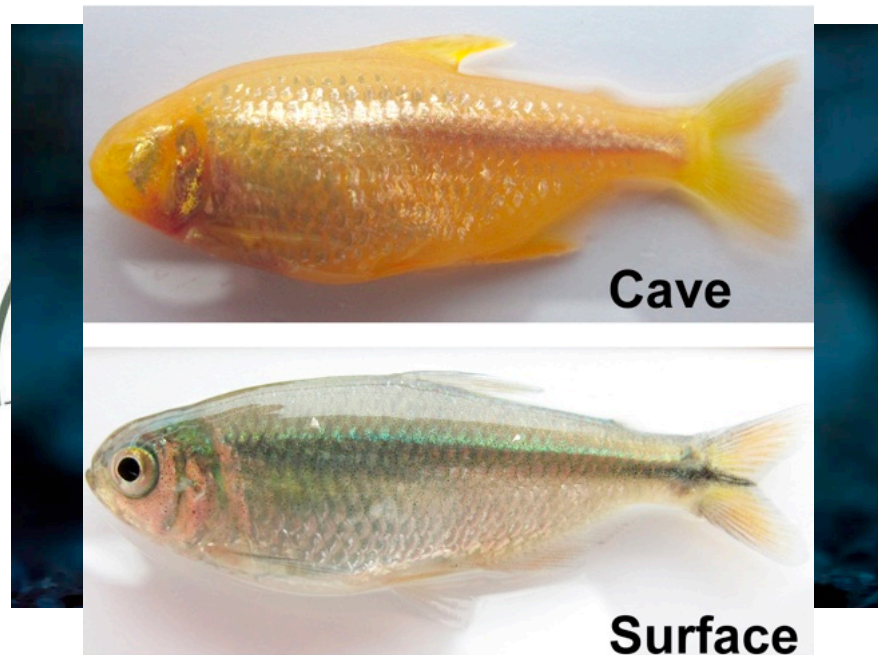
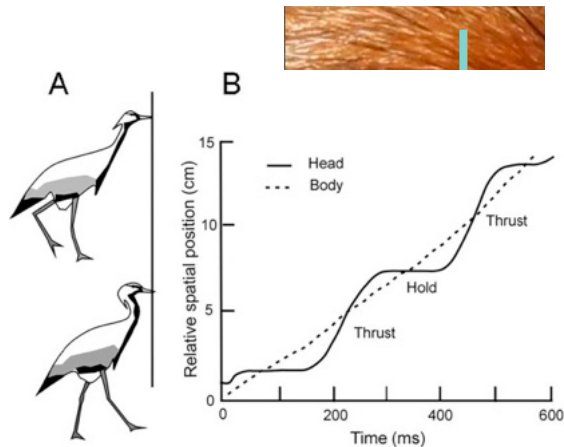
Gemini



Qwen

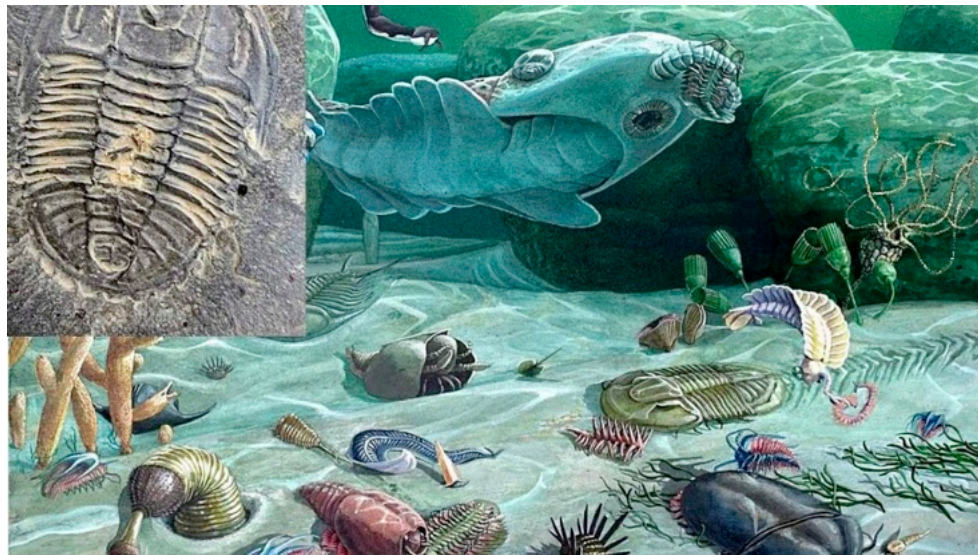


- **Visual perception** is closely linked to agent's actions, body, and surrounding environment.



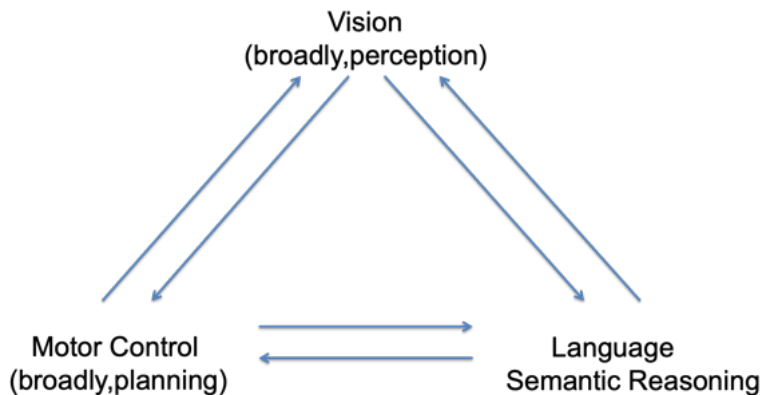
• M. Bank et al., 2015. M. Land, 2002

- Cambrian explosion: ~500 million years ago.
- Light switch theory, A. Parker.



Why vision+action?

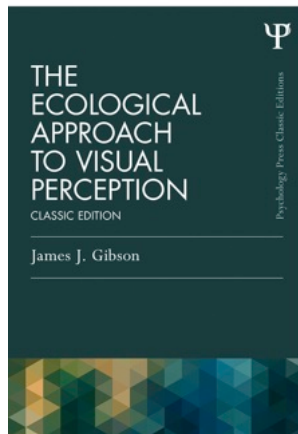
- ~500 million years ago: light switch. Mostly locomotion.
- ~5 million years: bipedalism appeared. Freed hands for tool making/use/ manipulation.
- 50k-150k: modern homo sapiens. Language appears.



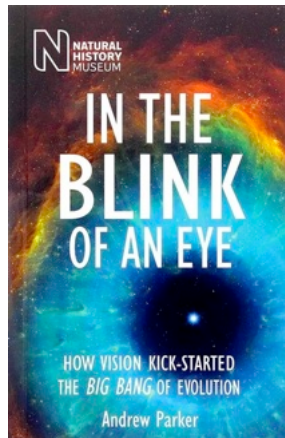
Why vision+action?

- Abundant evolutionary, biological, practical, and computational arguments.

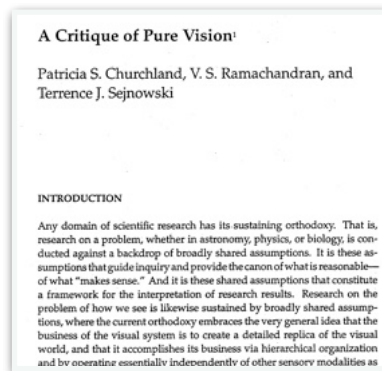
The “ecological vision” argument



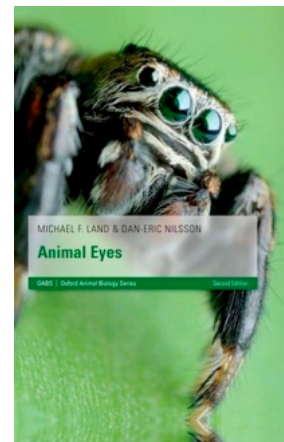
Light switch theory, A. Parker



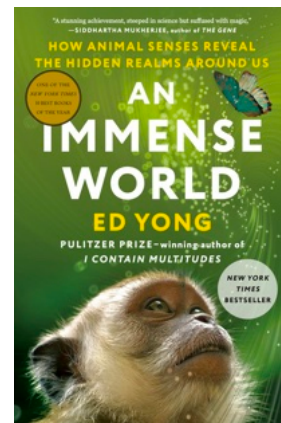
A Critique of Pure Vision, Churchland et al.



Animal Eye, M. Land



“Umwelt”

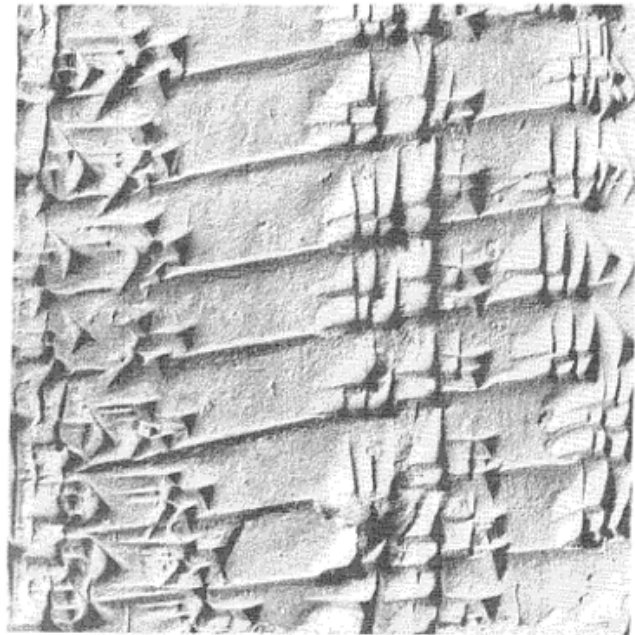


Why the real world/ecology?

Why the real world/ecology?



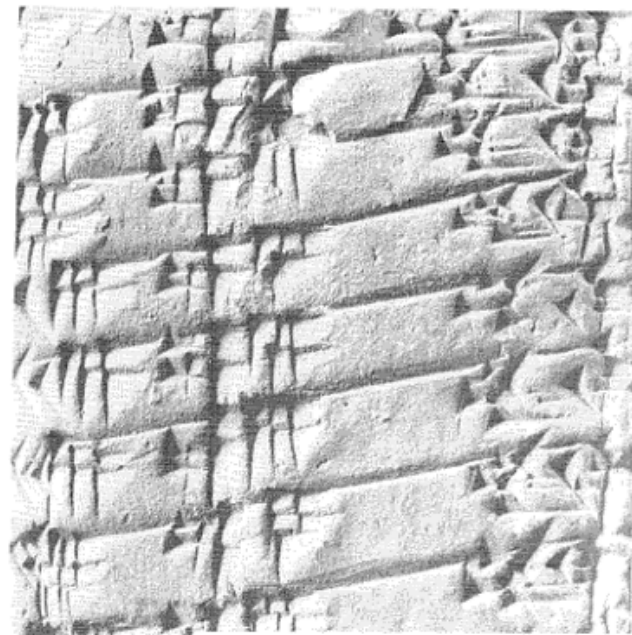
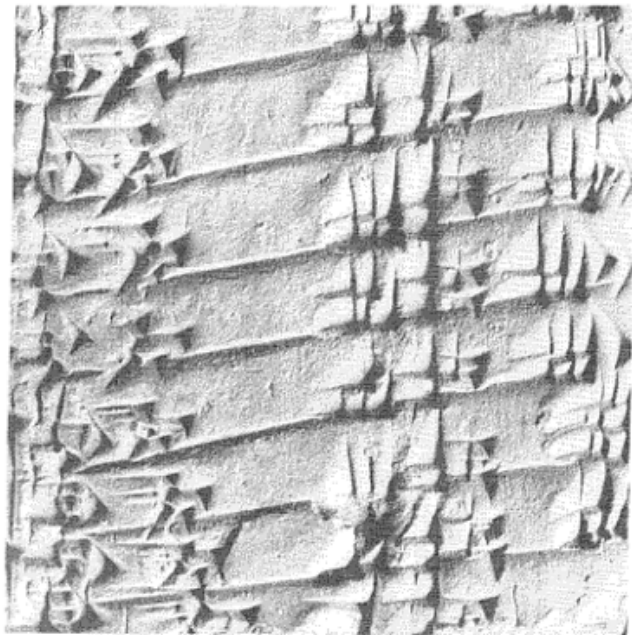
Why the real world/ecology?



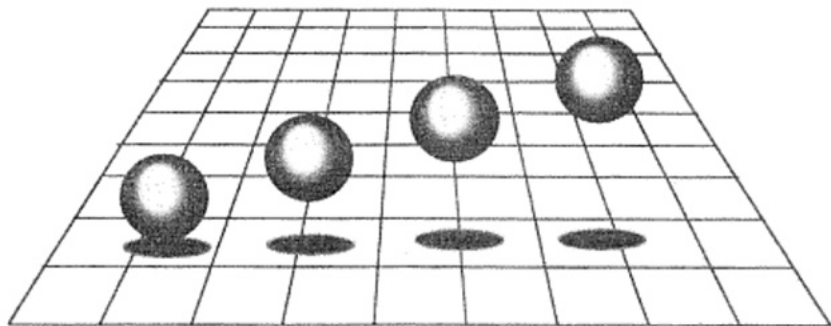
Why the real world/ecology?



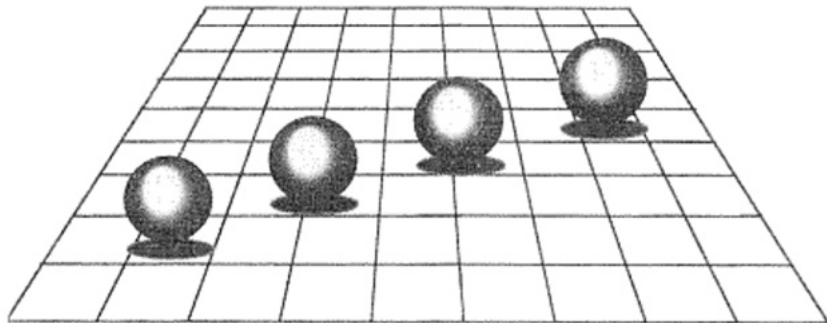
Why the real world/ecology?



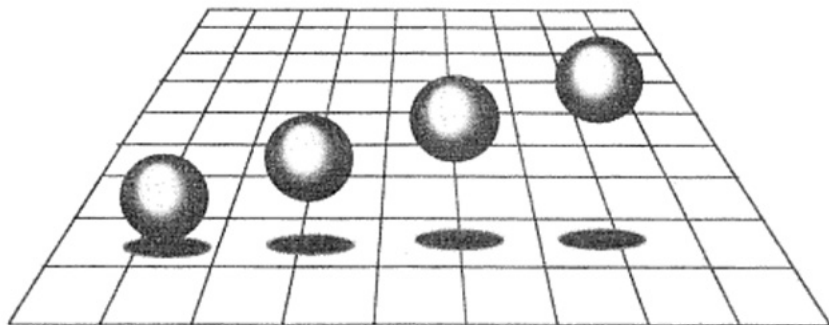
Why the real world/ecology?



Why the real world/ecology?

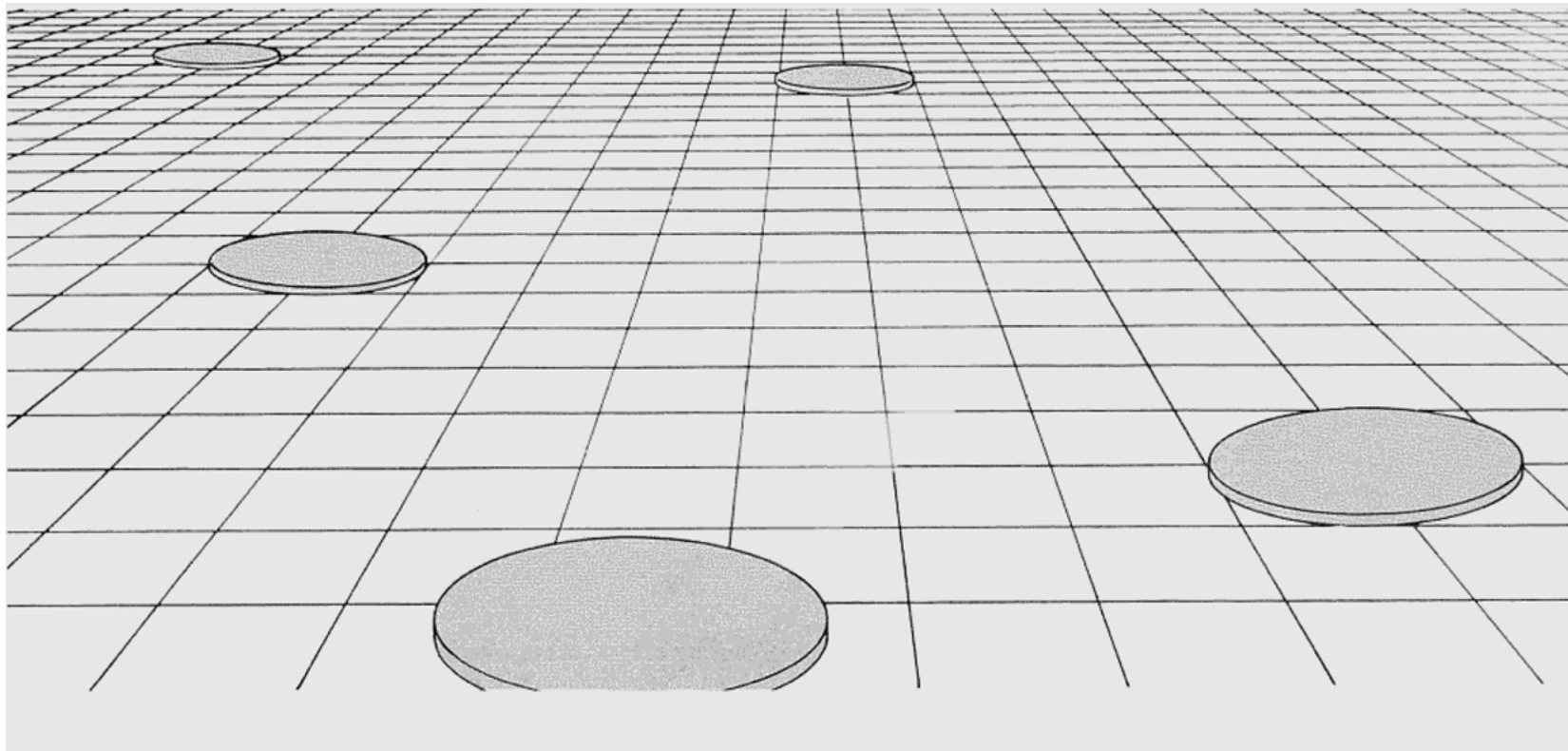


A



B

Why the real world/ecology?



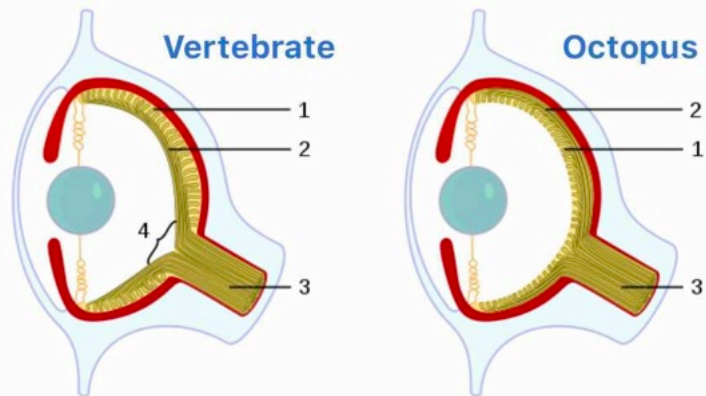
Why natural vision?

- Because most working examples of sophisticated perceptual intelligence are there!



But only looking at nature isn't enough.

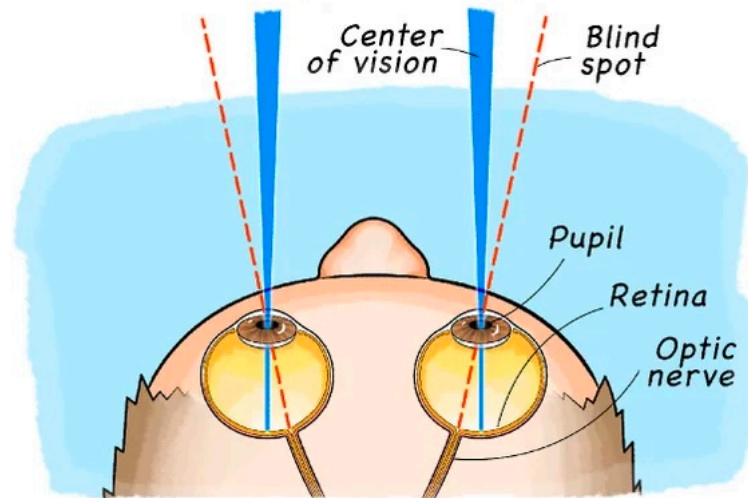
- Copying evolutionary bugs: blind spot (scotoma)



In **vertebrate** eyes, the nerve fibers route *before* the **retina**, blocking some light and creating a blind spot where the fibers pass through the retina and out of the eye. In **octopus eyes**, the nerve fibers route *behind* the retina, and do not block light or disrupt the retina. In the example, **4** denotes the vertebrate blind spot, which is notably absent in the octopus eye. In both images, **1** denotes the retina and **2** the nerve fibers, including the **optic nerve (3)**.

But only looking at nature isn't enough.

- Copying evolutionary bugs: blind spot (scotoma)
 - Check it out: look straight ahead. Close one eye. extend your arm and thumb out. Move slowly from center to side. Your thumb will disappear around ~ 15 degree from center.



But only looking at nature isn't enough.



- Cover your left eye and look at the plus sign, move closer or further away to find where the circle disappears. (Try on your screen, not the projector, since you have to move your head a lot)
 - That location is where your optic nerve attaches to your eyeball and is where no visual information is processed due to the lack of rods and cones, otherwise known as the blind spot

But only looking at nature isn't enough.

■ Wright Brothers

- *“Learning the secret of flight from a bird was a good deal like learning the secret of magic from a magician. After you once know the trick and know what to look for you see things that you did not notice when you did not know exactly what to look for.”*



ORVILLE WRIGHT
DAYTON, OHIO

December 27, 1941.

Mr. Horace Lytle, President,
The J. Horace Lytle Company,
Dayton, Ohio.

Dear Mr. Lytle:-

Your letter of November 26th was duly received, but having become buried among other papers, it has just come to my attention again.

I can not think of any part bird flight had in the development of human flight excepting as an inspiration. Although we intently watched birds fly in a hope of learning something from them I can not think of anything that was first learned in that way. After we had thought out certain principles, we then watched the bird to see whether it used the same principles. In a few cases we did detect the same thing in the bird's flight.

Learning the secret of flight from a bird was a good deal like learning the secret of magic from a magician. After you once know the trick and know what to look for you see things that you did not notice when you did not know exactly what to look for.

Sincerely yours,

Orville Wright

But only looking at nature
isn't enough.





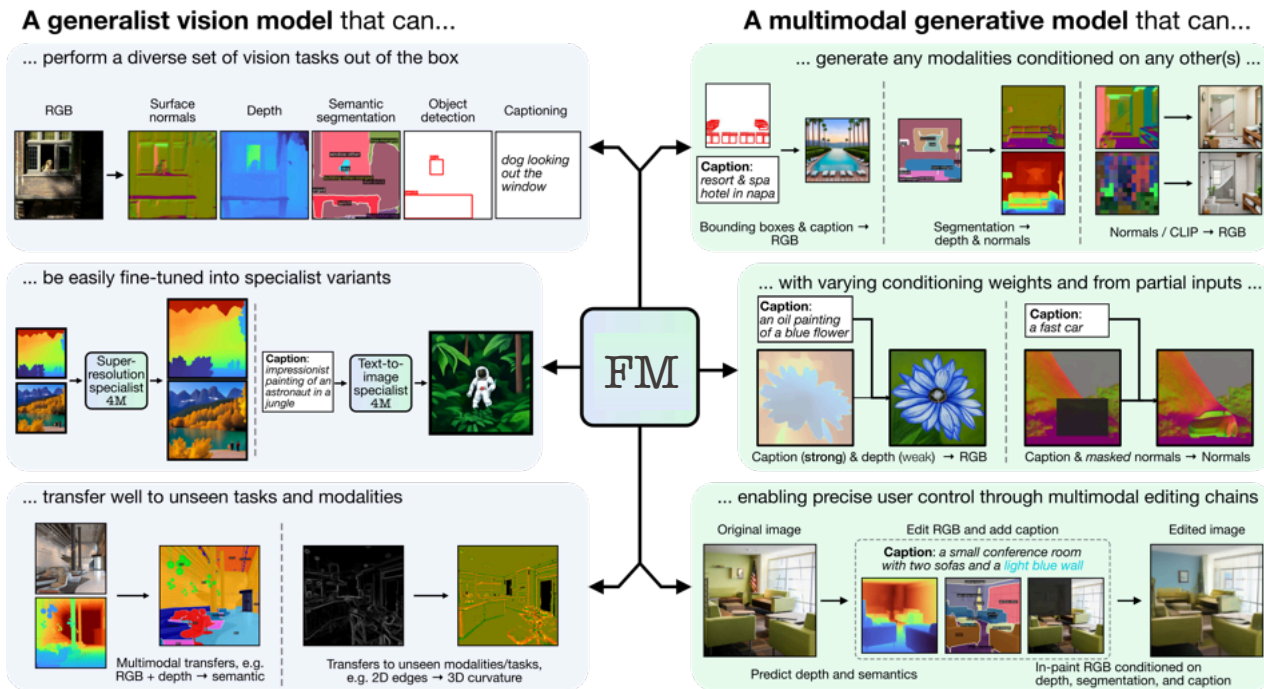
Qwen



EPFL Goal: develop a “Foundation Model” from scratch

24

Zamir



- develop and enhance a foundation (multimodal large vision-language) model end-to-end

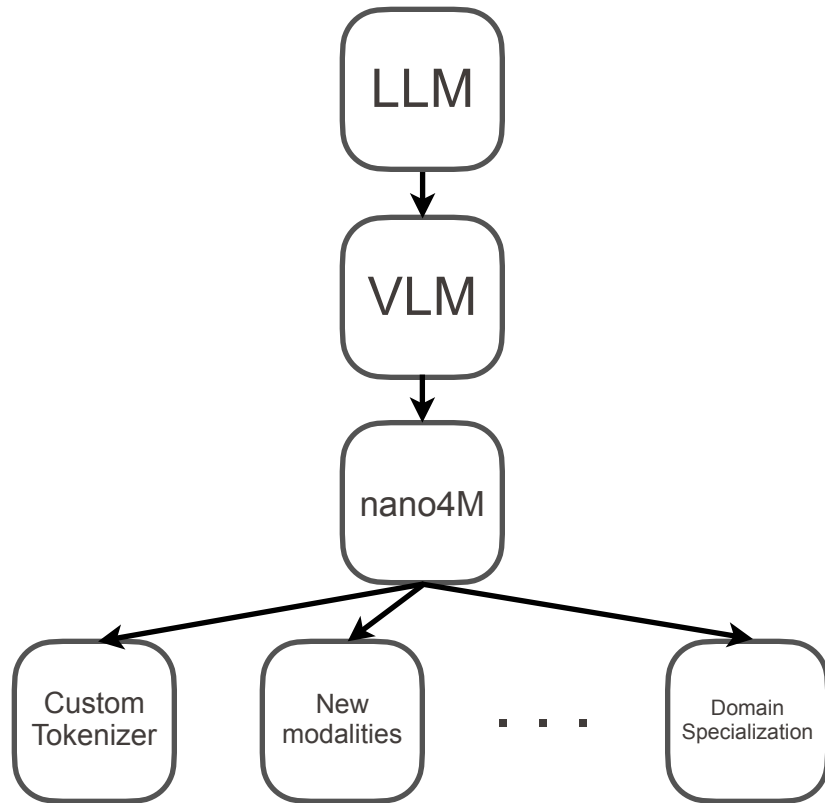
EPFL Goal: develop a “Foundation Model” from scratch

25

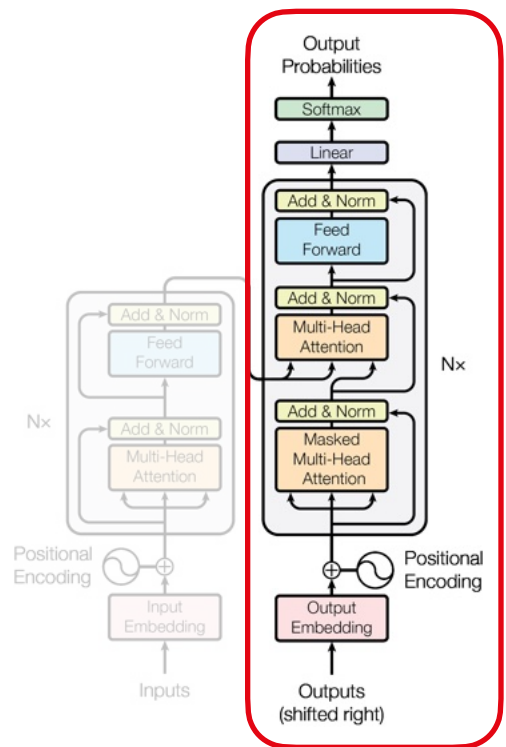
Zamir

- **Basic Task:** develop a nano FM
 1. LLM
 2. Text-to-Image generation
 3. nano4M
- **Extensions:** propose your own, e.g.
 - Custom tokenizers
 - New modalities
 - Domain specialization
 - Reasoning
 - etc.

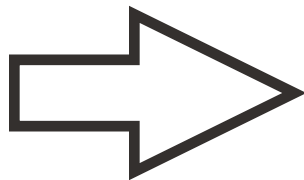
Project model development flow



EPFL 1. Develop: LLM



Transformer Decoder



Generating text

Max prompt length = 5

the robots will bring

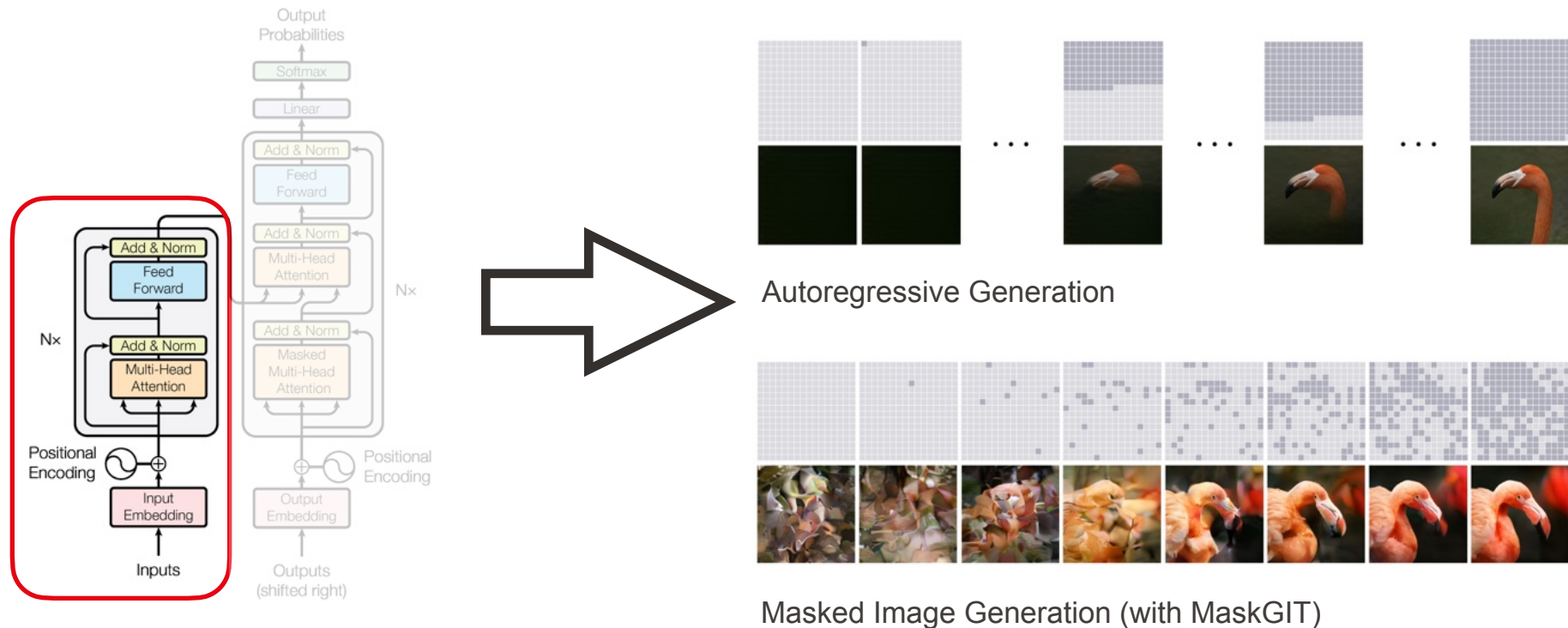
Transformer



Qwen



EPFL 2. Develop: Text-to-Image Generator

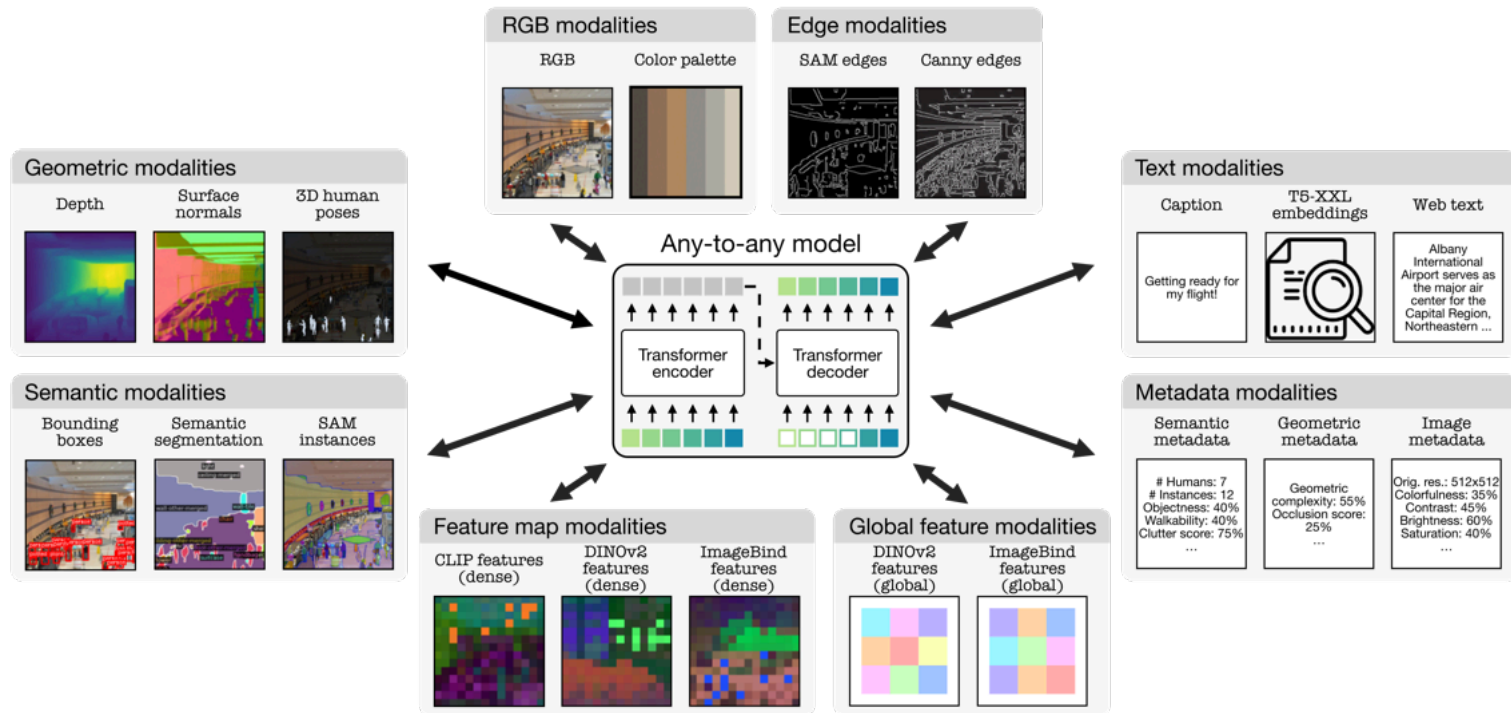


Transformer Encoder

Autoregressive Generation

Masked Image Generation (with MaskGIT)

3. Develop: nano4M



Multi-modal Encoder-Decoder

3. Develop: nano4M

Tokenization

Bounding boxes

xmin=0.30 ymin=0.51
xmax=0.68 ymax=0.99
horse

4M chained multimodal generation

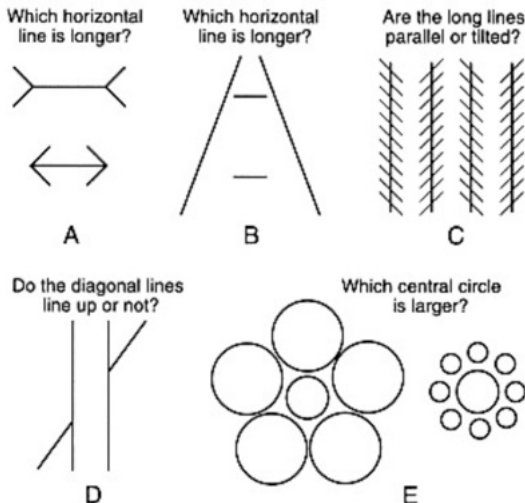
Iteration 1 2 3 4 5 6 7 8 9 10

Transformer
encoder

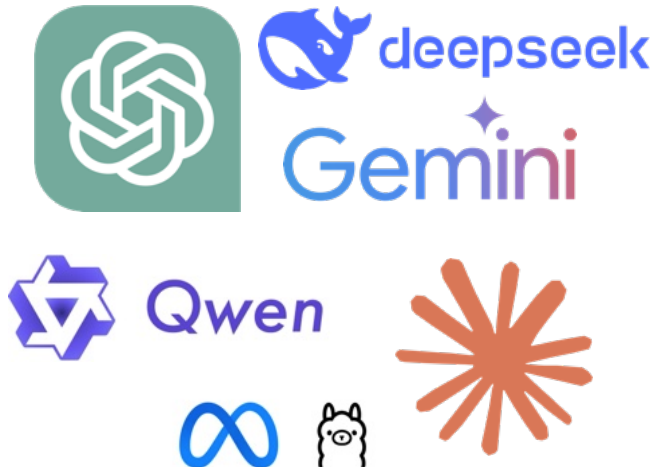
Transformer
decoder

Course Topics and Lectures

- **Basics:** biological vision, theories, etc.



- **Foundation Models:** visual and multimodal models.



Course Topics and Lectures

- **Foundation/Concept Lectures**
 - e.g. theories of vision, sensory umwelt, perception-action loop, multimodality, etc.
- **SOTA/Tools Lectures**
 - e.g., Architectural/transformer advances, FMs, etc.
- **Domain Expert Lectures**
 - e.g., neuroscience, RL, etc.

- Credit for a fraction of visuals/slides: James J. Gibson, Michael Land, Bruno Olshausen, Jitendra Malik, Stella Yu, Alexei Efros, Dan-Eric Nilsson.

course logistics and policies

- Instructor: Amir Zamir (amirzamir.com)
- Doctoral Teaching Assistants: (vimm-ta@groupes.epfl.ch)
 - Rishubh Singh (head TA)
 - Roman Bachmann
 - Zhitong Gao
- Group: <https://vilab.epfl.ch/>

- Anyone interested in open-ended learning and research in vision and perceptual AI.
 - In particular **active & multimodal** intelligent agents in the **real world**.
 - Emphasis on exploring and what is **not** possible today, than what it is.
- Part of the course is about defining **problems** than providing **solutions**.
 - A vast sea between the two. Eg:
 - Solution: we know how to train a DNN to predict 1000 predefined classes given an image
 - Problem: how to use DNNs to have an agent that interacts with an environment, continually learns a model of it, and performs in it.
 - Often there is no clear-cut solution.

- Machine Learning
 - E.g. Machine Learning (CS-433) or Introduction to Machine Learning (CS-233) or equivalent course on the basics of machine learning.
- Deep Learning
 - E.g. Deep Learning (EE-559) or Artificial Neural Networks (CS-456) or equivalent course on the basics of deep learning.
- For projects with an active agent component, Reinforcement Learning or control theory.
- Expertise at one of the common Deep Learning frameworks, e.g. Pytorch.
- Prior experience with Computer Vision (e.g. CS-442) or vision data recommended.

- **Project (60%)**
 - Project proposal (15%)
 - Project progress reports (15%)
 - Final project presentation (15%) and report (15%)
 - Presentation: the last week of the semester.
 - Report/presentation templates provided.
- **Homework (40%)**
 - 3-5 homework and notebook assignments.
- **Late Policy:** total of 3 late days with no penalty. 0 otherwise.

- Get inspired!
- Most of your learning will happen while doing the project and homework.
- Desired characteristics of projects:
 - Related to vision, perceptual agents, etc.
 - Cool and creative (in problem selection, solution formulation, implementation, demonstration)
 - Thought-through.
 - Well-executed.
 - Does not have to be hardware or software. Up to you.
 - You learn something from it.

- BristleBot
- Seemingly complex behavior can come out of simple pieces.
- An active agent doesn't have to look like Boston Dynamic's Atlas.



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

Existing Vision/ Robotic/FM Challenges

- E.g.
 - CVPR Embodied AI Challenges



Analysis/debunking Projects

Do ImageNet Classifiers Generalize to ImageNet?

Benjamin Recht* Rebecca Roelofs Ludwig Schmidt Vaishal Shankar
UC Berkeley UC Berkeley UC Berkeley UC Berkeley

Abstract

We build new test sets for the CIFAR-10 and ImageNet datasets. Both benchmarks have been the focus of intense research for almost a decade, raising the danger of overfitting to excessively re-used test sets. By closely following the original dataset creation processes, we test to what extent current classification models generalize to new data. We evaluate a broad range of models and find accuracy drops of 3% – 15% on CIFAR-10 and 11% – 14% on ImageNet. However, accuracy gains on the original test sets translate to larger gains on the new test sets. Our results suggest that the accuracy drops are not caused by adaptivity, but by the models' inability to generalize to slightly "harder" images than those found in the original test sets.

A Metric Learning Reality Check

Kevin Musgrave¹, Serge Belongie¹, Ser-Nam Lim²

¹Cornell Tech ²Facebook AI

Abstract. Deep metric learning papers from the past four years have consistently claimed great advances in accuracy, often more than doubling the performance of decade-old methods. In this paper, we take a closer look at the field to see if this is actually true. We find flaws in the experimental setup of these papers, and propose a new way to evaluate metric learning algorithms. Finally, we present experimental results that show that the improvements over time have been marginal at best.

From Same Photo: Cheating on Visual Kinship Challenges*

Mitchell Dawson¹[0000-0002-6719-6584], Andrew Zisserman¹[0000-0002-8945-8573],
and Christoffer Nellaker²[0000-0002-2887-2068]

¹ Visual Geometry Group (VGG), Dept. of Engineering Science, University of Oxford
{mdawson, az}@robots.ox.ac.uk

² Nuffield Dept. of Women's & Reproductive Health, Big Data Institute, IBME,
University of Oxford
christoffer.nellaker@bdi.ox.ac.uk

Abstract. With the propensity for deep learning models to learn unintended signals from data sets there is always the possibility that the network can "cheat" in order to solve a task. In the instance of data sets for visual kinship verification, one such unintended signal could be that the faces are cropped from the same photograph, since faces from the same photograph are more likely to be from the same family. In this paper we investigate the influence of this artefactual data inference in published data sets for kinship verification.





- **Hunting for Insights: Investigating Predator-Prey Dynamics through Simulated Vision and Reinforcement Learning**
- Arvind Menon, Lars C.P.M. Quaedvlieg, Somesh Mehra
- <https://arvind6599.github.io/PredatorPreyWebsite/>

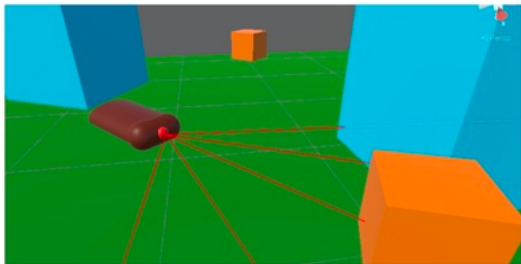


Figure 3. Visualization of the agent's vision with 5 eyes

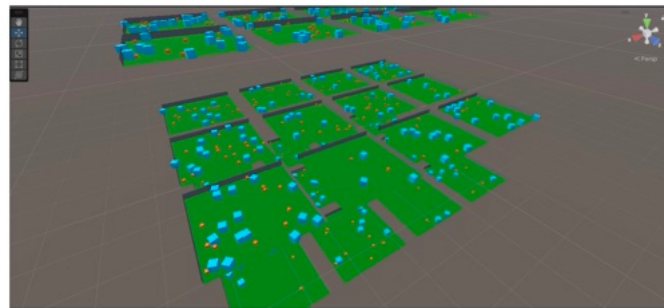
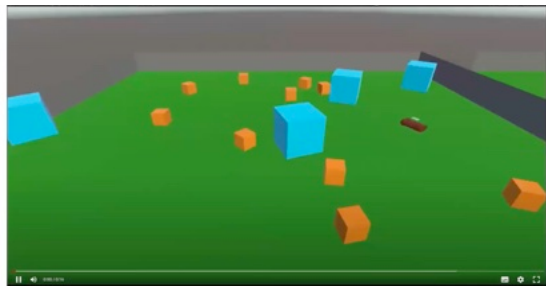


Figure 2. Multiple environments to speed up training

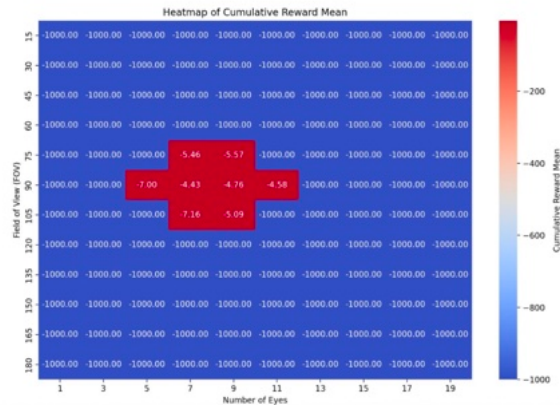
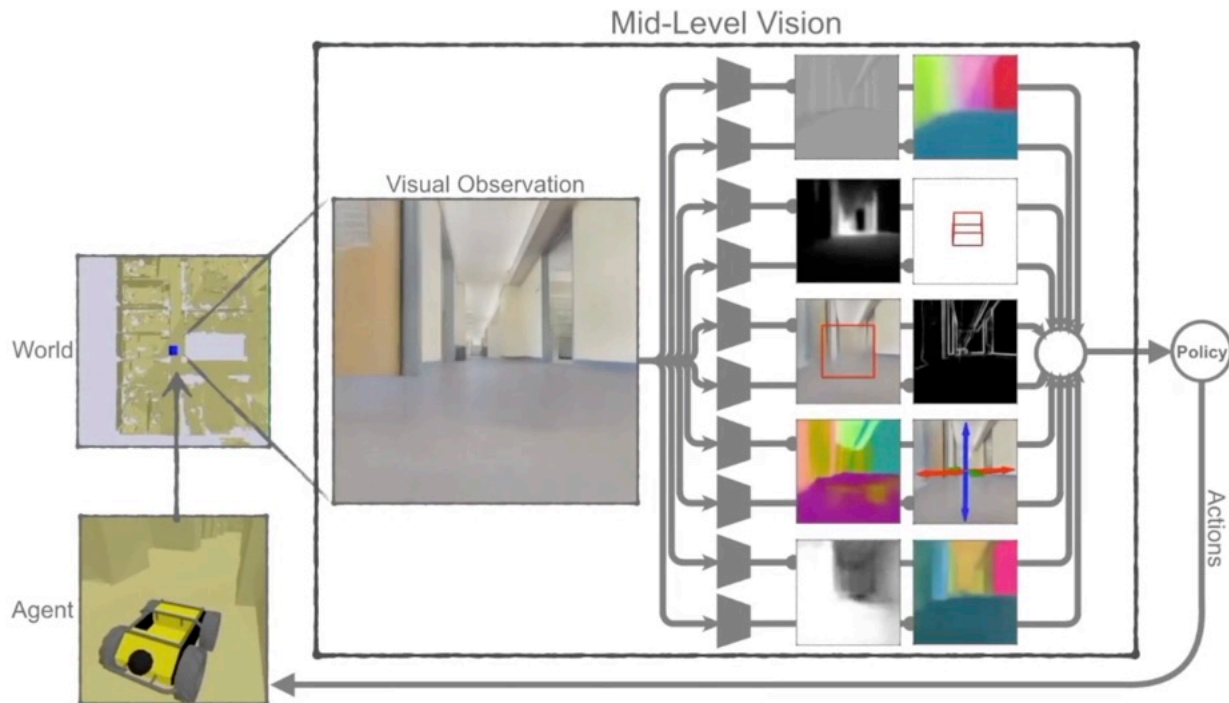


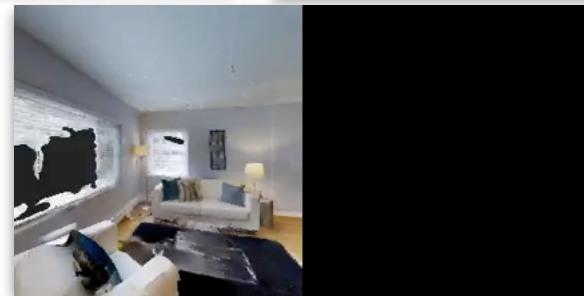
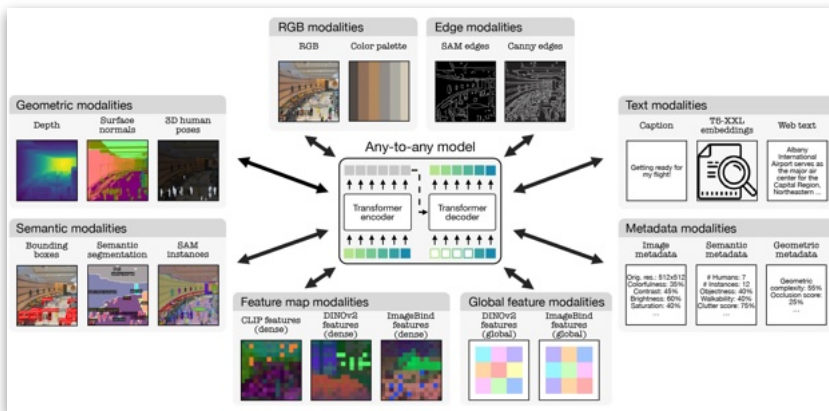
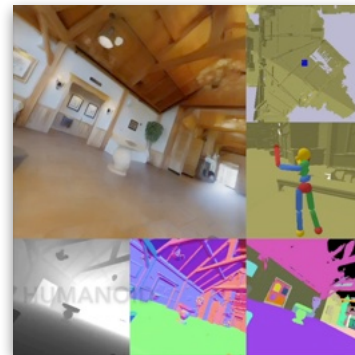
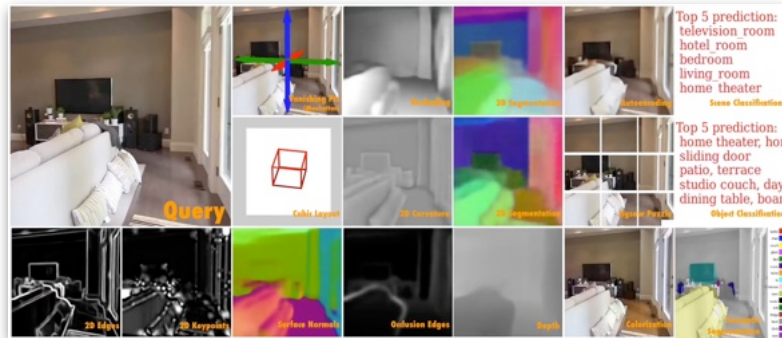
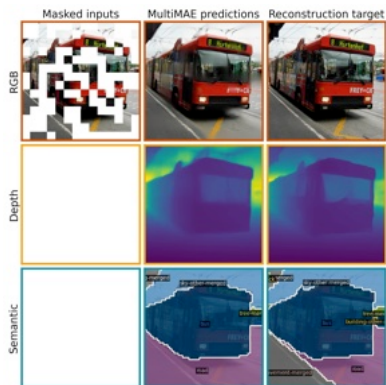
Figure 6. Exploration result starting with 9 eye and 90 fov, that converged on 7 eyes and 90 fov

- **Vision Evolution: Natural Selection shaped by environment**
- Yassine Abdennadher, Ambroise Borbely, Matthieu Andre.

EPFL Past example



- Bradley Emi
- Mid-Level Vision, CoRL 2019, CoRL 2020. Sax et al.



Projects: Think about the unsolved



- More categorical examples:
 - Make a visual agent inspired by a simple biological organism.
 - Extend the FMs you develop.
 - Simulate the sensory “umwelt” of an animal.
 - Add a camera to any visually blind (or too heavily sensor loaded) system out there.
 - Stationary cameras: extract actionable intelligence. E.g. a camera in the corner of the living room, the passive surveillance cameras, etc.
 - Apply/analyze one of the concepts discussed in the class.
 - Take a concept from one of the books and apply it computationally.
 - Any of relevant CVPR/ICCV/ECCV/NeurIPS challenges.
 - Analysis study
- We will give you feedback on your proposals. Do seek feedback from TAs extensively.

- The goal is doing something cool and learn. We're here to help you. We'll be flexible with solid efforts.
- Use the teaching team.
- The top projects are of acceptable quality to top-tier AI/vision/ML conferences.
 - (CVPR/ICCV/ECCV/NeurIPS/ICML/ICLR/CoRL/RSS/ICRA)
- The Best Project will be selected and advertised. Awarded with conference registration.
- Well-distribute your time (continuous evaluation).
- **Teaming policy:** encouraged, but $1+1 \nless 2$.
- **Merging with an existing project:**
 - In principle, okay. But 1) No work by others 2) New work done. The report should clarify.

- SCITAS
 - Anyone at EPFL can submit jobs via their scheduler (SLURM)
 - We reserved some GPUs for this course. More available if capacity unused.
- Google Cloud
- Azure
- Free tier (open to everyone)
 - Google Cloud: \$300 for 90 days
 - Azure: \$200 for 30 days

- **Front loading lectures:** About 8 weeks of lectures 2x per week, then you're free to work on your projects full time.
 - The class hours turn into support sessions for troubleshooting and discussion with TAs.
- **Books and reading material** references on Moodle and the library (use them!)
 - Additional resources and reading will be occasionally placed on Moodle
 - **Have a sharing, fun, passionate culture:** find something related interesting? An article, video, book, meme? Just post it on the forum for everyone.
- **Amir's office hours:** Tuesdays 15:00-16:00 (by prior email appointment)
- **Feedback From:**
 - Anonymous. Any thoughts and suggestions welcome.
 - <https://forms.gle/K7RKKafcEgfEJjhz7> (link on Moodle)

EQUAL OPPORTUNITY
MAKES US BRIGHTER



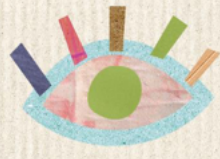
DIVERSITY IS
OUR STRENGTH



LISTEN AND DARE
TO SPEAK UP



LOOK OUT
FOR EACH OTHER



KINDNESS
BRINGS SERENITY



TOGETHER
WE GO FURTHER



Mental health

**Brilliant
ideas in a
peaceful mind**



Enjoy the Course!

<https://vilab.epfl.ch/>