

Introduction

- Instructor Mengye Ren: Assistant Professor in Data Science and Computer Science.
- Before NYU: Google and Uber ATG (self-driving)
- Agentic Learning AI Lab: Enabling future agentic AIs to learn and adapt flexibly in the real world.
- Representation learning, self-supervised learning meta-learning, continual learning, visual perception and planning



Mengye Ren
Assistant Professor



Chris Hoang
PhD Student



Jack Lu
PhD Student



Ryan Teehan
PhD Student



Ying Wang
PhD Student



Yanlai Yang
PhD Student



Wun Ting Chan
Master Student



Amelia (Hui)
Dai
Master Student



Xavier
(Xiaoyang)
Jiang
Master Student



Jinran Jin
Master Student



Dahye Kim
Master Student



Anurup Naskar
Master Student



Arjun Prasad
Master Student



Zhenbang Yang
Master Student



Yuen-Hei
Yeung
Master Student



Weizhen Zhou
Master Student



Azwar
Abdulsalam
Visiting Researcher



Steven Luo
Visiting Researcher



Frank
(Zequan) Wu
Visiting Researcher



Will Wu
Undergraduate Student

Office Hours and Communication

- Office Hour: Thursday 2:00pm – 3:00pm Room 508, 60 5th Ave
- By appointment. Scheduling link in the syllabus.
- Non-admin related question: Use CampusWire
- Admin related question (e.g. homework): Email TAs

TA Intro



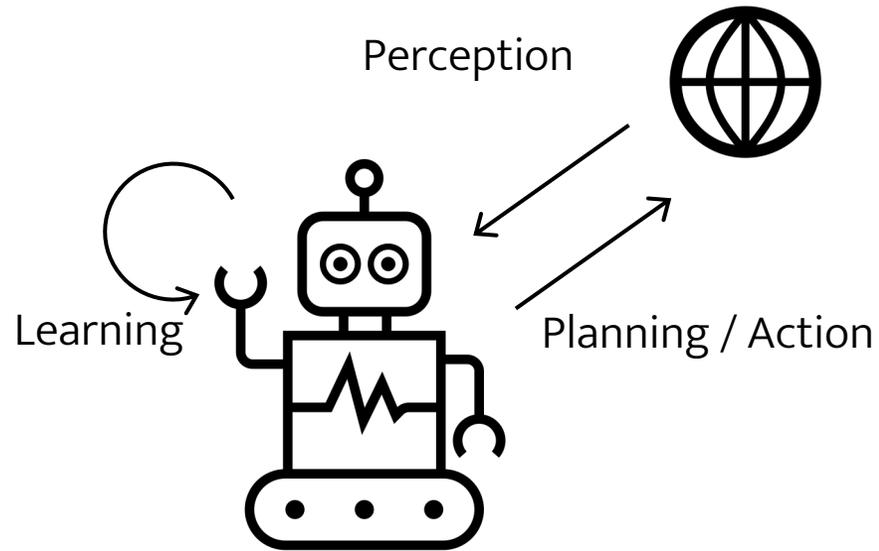
Ying Wang
Thu 1-2PM
Room 763
Latent Planning



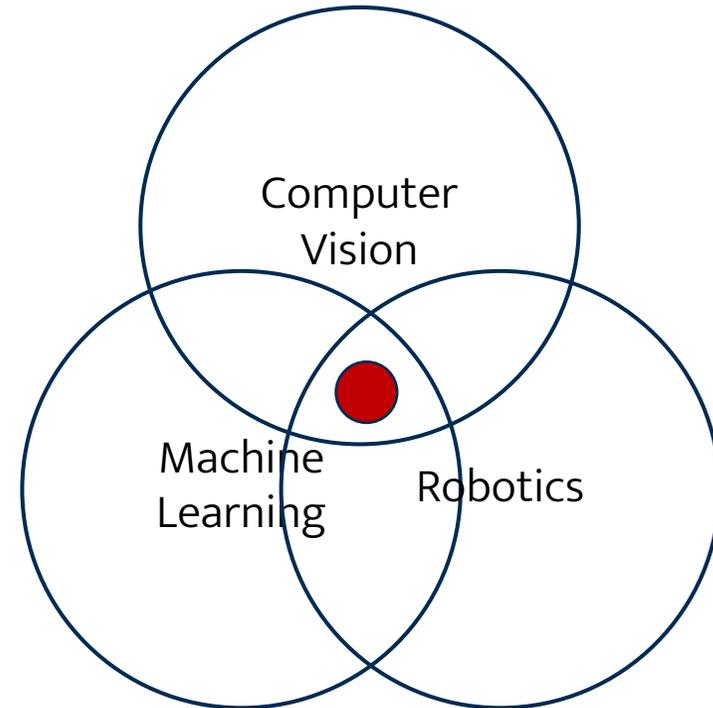
Ellis Brown
Tue 11AM-12PM
Room 402
Spatial Understanding

Introduction

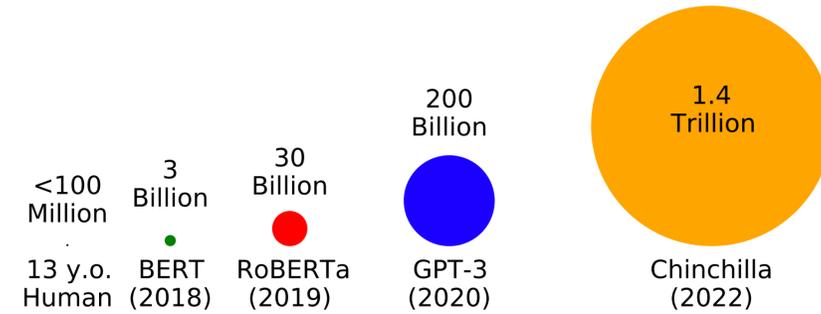
What is this course about?



A General-Purpose Learning Agent



Human vs. Machine Learning



Language Vs. Embodied Video?



Yann LeCun  
@ylecun



* Language is low bandwidth: less than 12 bytes/second. A person can read 270 words/minutes, or 4.5 words/second, which is 12 bytes/s (assuming 2 bytes per token and 0.75 words per token). A modern LLM is typically trained with 1×10^{13} two-byte tokens, which is 2×10^{13} bytes. This would take about 100,000 years for a person to read (at 12 hours a day).

* Vision is much higher bandwidth: about 20MB/s. Each of the two optical nerves has 1 million nerve fibers, each carrying about 10 bytes per second. A 4 year-old child has been awake a total 16,000 hours, which translates into 1×10^{15} bytes.

The Development of Embodied Cognition: Six Lessons from Babies

- Multimodal
- Incremental
- Physical
- Explore
- Social
- Use language

Abstract The embodiment hypothesis is the idea that intelligence emerges in the interaction of an agent with an environment and as a result of sensorimotor activity. We offer six lessons for *developing* embodied intelligent agents suggested by research in developmental psychology. We argue that starting as a baby grounded in a physical, social, and linguistic world is crucial to the development of the flexible and inventive intelligence that characterizes humankind.

Linda Smith

Psychology Department
Indiana University
Bloomington, IN 47405
smith4@Indiana.edu

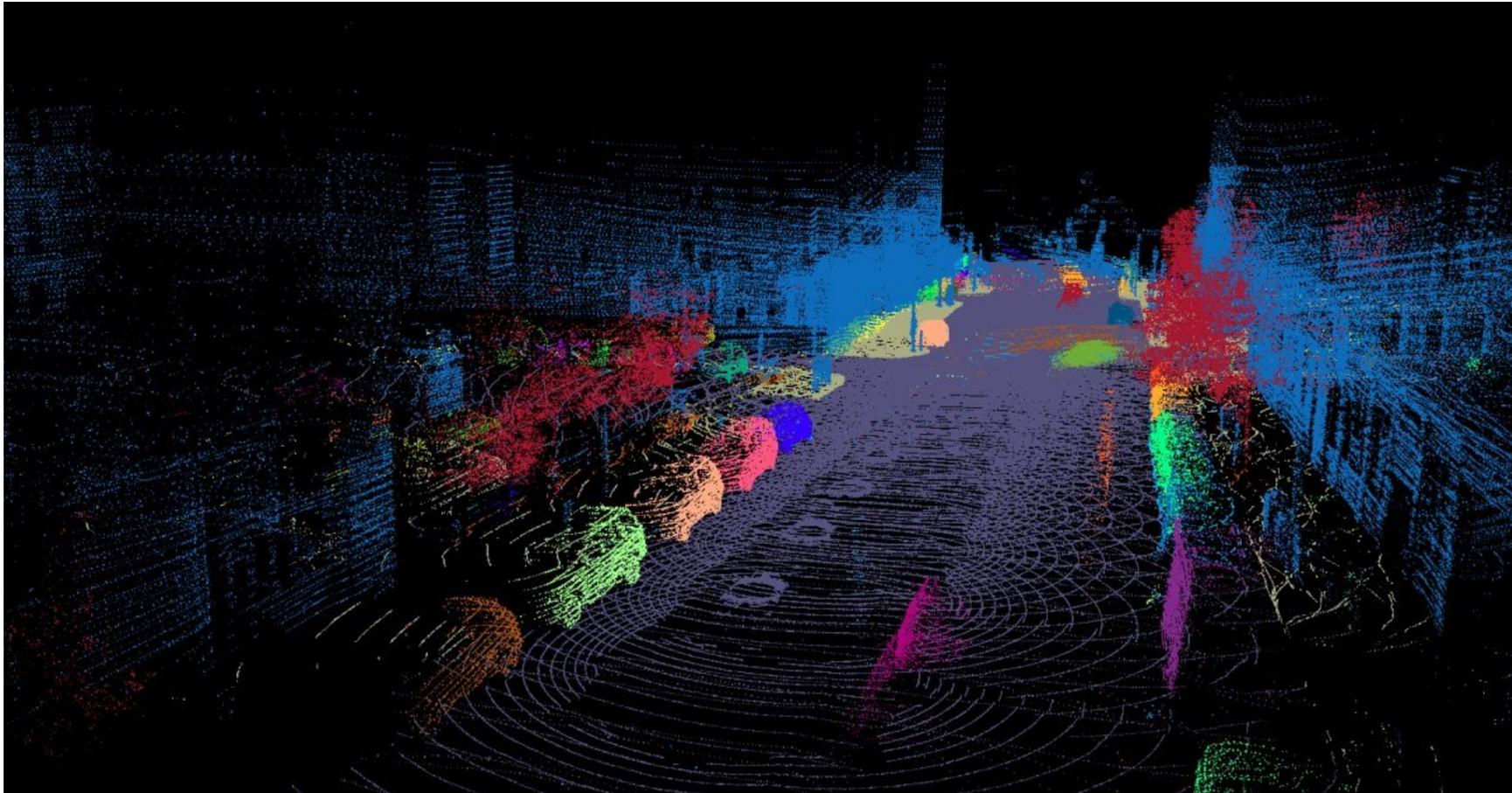
Michael Gasser

Computer Science Department
Indiana University
Bloomington, IN 47405
gasser@Indiana.edu

Keywords

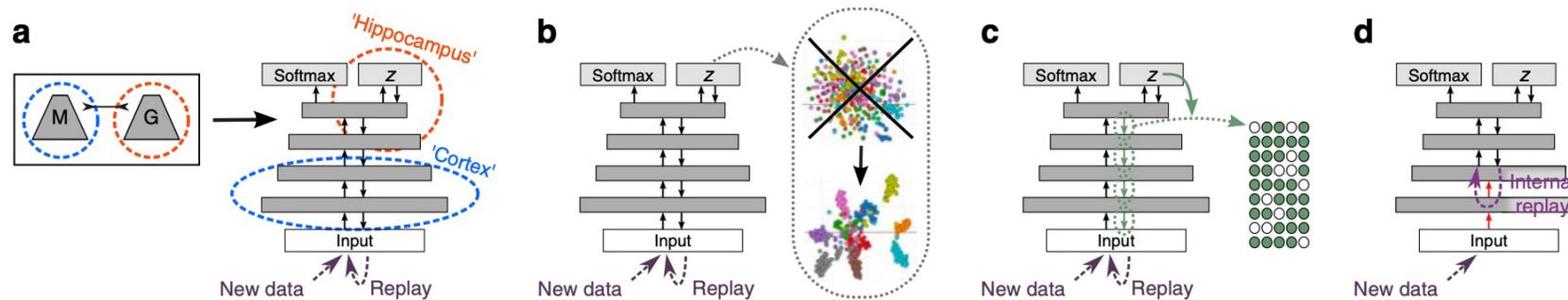
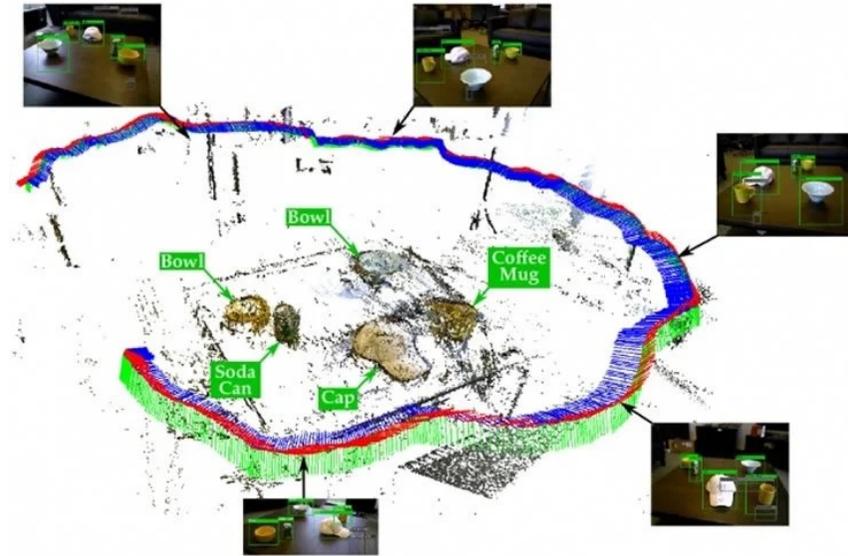
Development, cognition, language,
embodiment, motor control

Physical: Geometric and Temporal Structure



Incremental: Learning and Memory

- Spatial memory (mapping)
- Episodic memory (autobiography)
- Semantic memory (rule learning)
- Procedural memory (skill learning)
- Replay: Generative or storage, consolidation



van de Ven, 2020

Explore: Learning Efficiency



Learning Objectives

- Solving embodied problems using deep learning tools
- Leverage geometric and temporal structure from real-world and simulated data
- Understand and design learning and memory in embodied agents

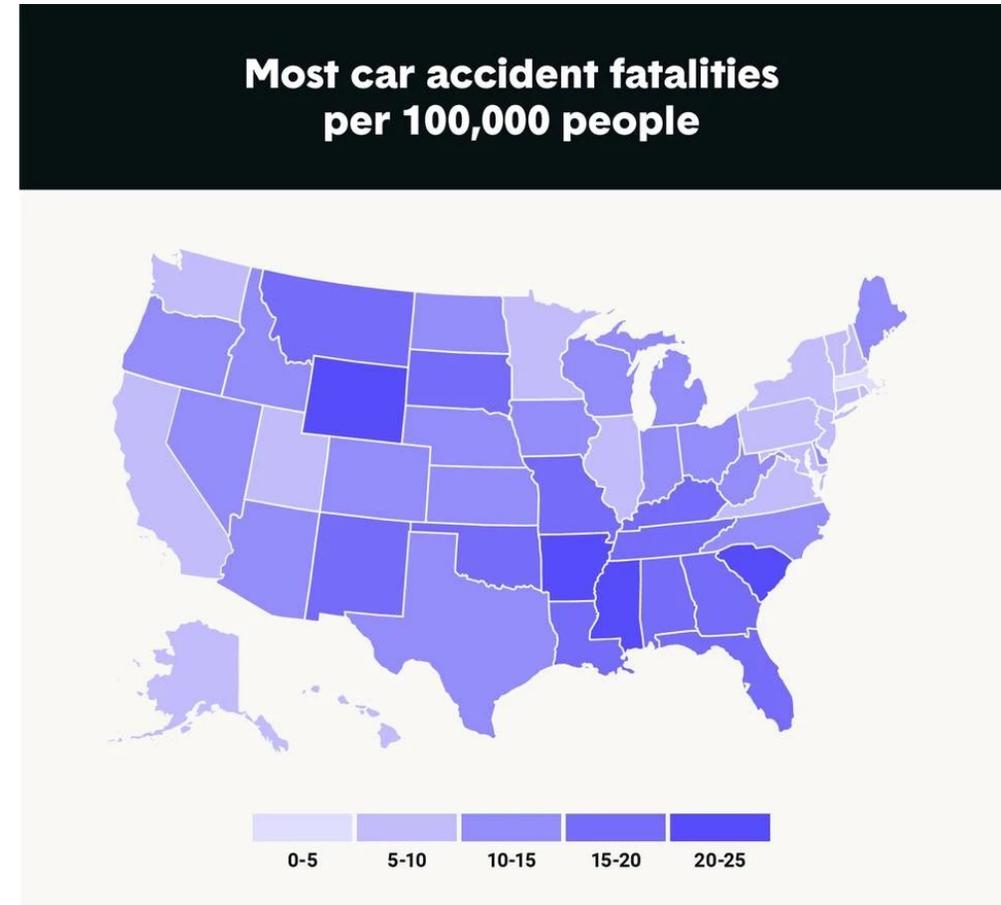
- Advanced graduate level course on a collection of topics
- Combination of lectures, seminars, and research projects
- Develop research skills and conduct cutting edge research

Applications



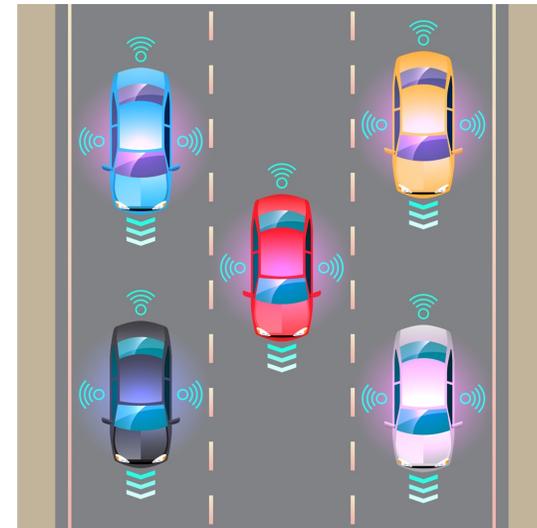
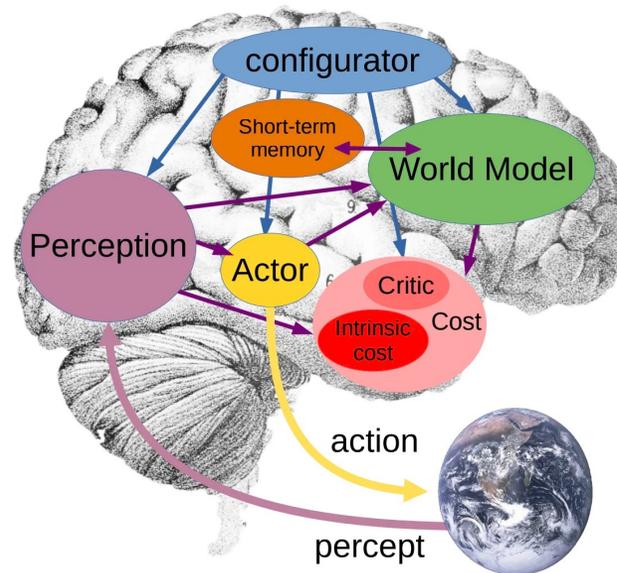
Why Self-Driving Cars

- >1 million people die every year from road
- One of the top 10 leading causes of death and injuries
- Environmental causes
 - Car/battery manufacturing

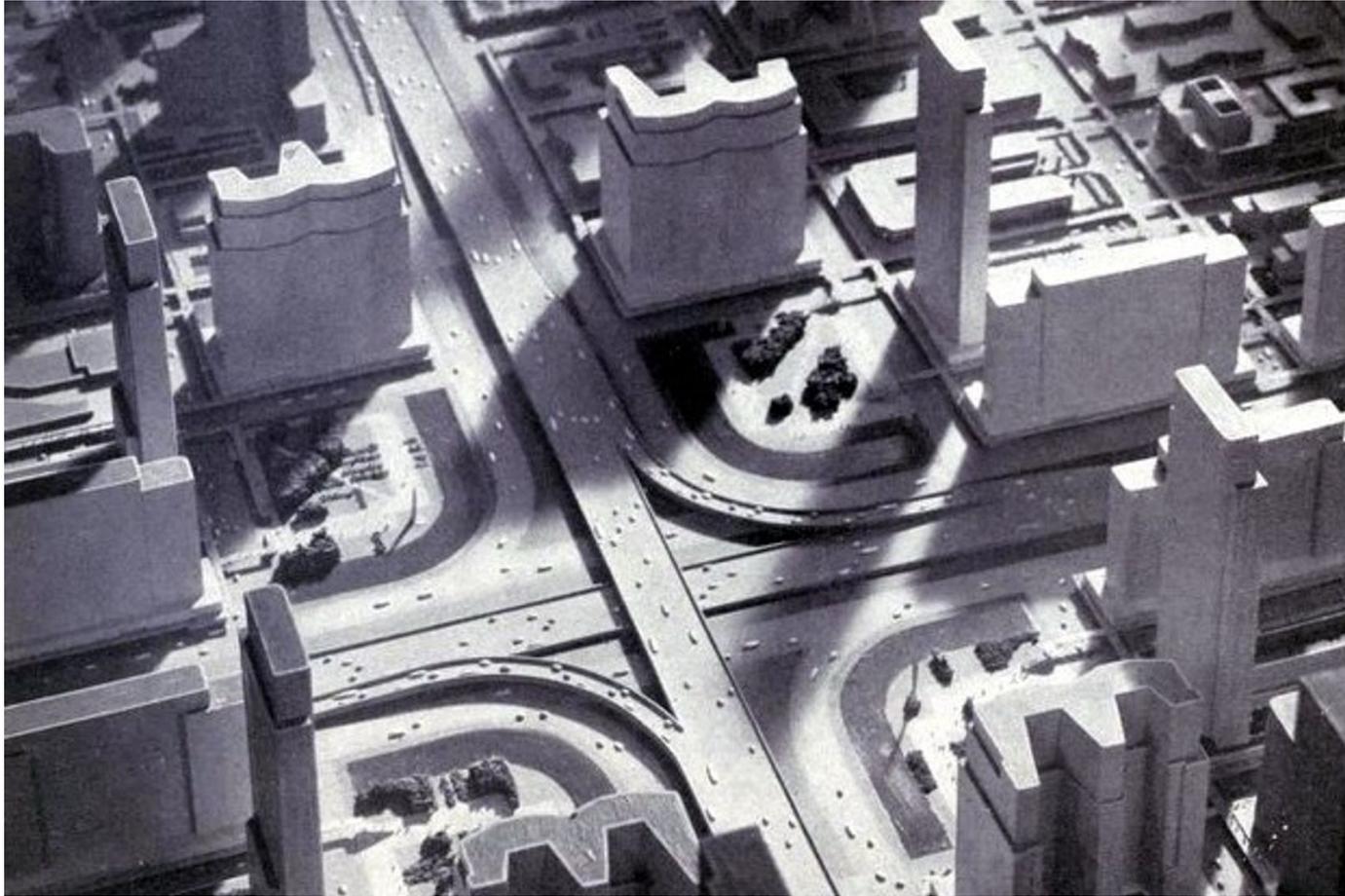


A Test-Bed for General Embodied Intelligence

- Perception, world model, mapping, planning, and control
- Long-term hierarchical planning
- Closed loop learning
- Multi-agent social cognition, intent inference, communication
- Rule-based learning



A Brief History of Self-Driving Cars



Norman Bel Geddes' "Magic Motorways" 1939

1966 LUNAR Stanford University



Also in 1966

MASSACHUSETTS INSTITUTE OF TECHNOLOGY

PROJECT MAC

Artificial 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".

1988 ALVINN

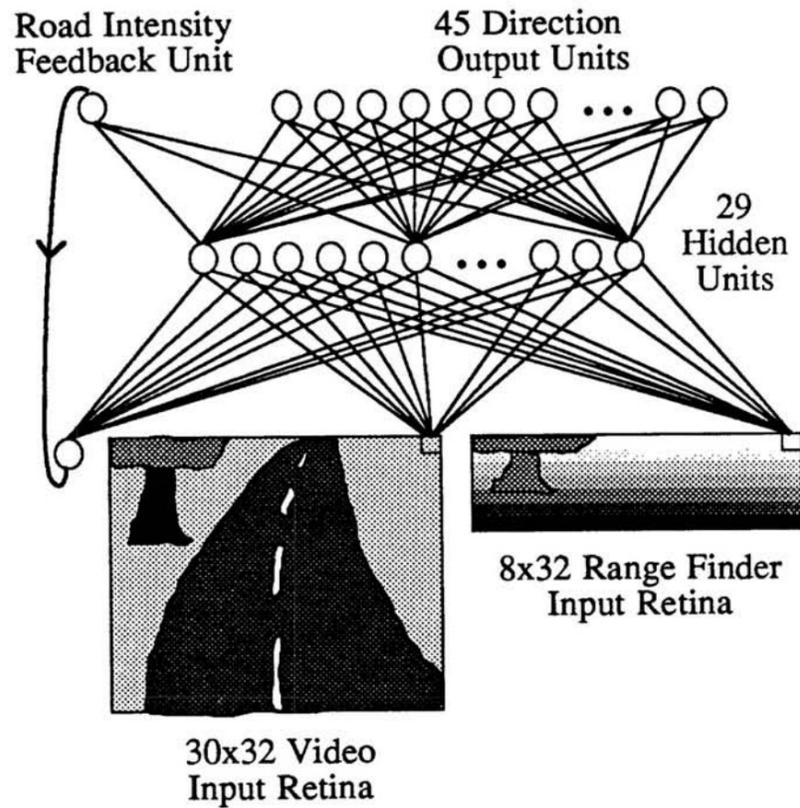


Figure 1: ALVINN Architecture



2004 DARPA Grand Challenge



2005 DARPA Grand Challenge



2007 DARPA Urban Challenge



<https://www.youtube.com/watch?v=aHYRtOvSx-M>

What Has Changed?

- Proof-of-Concept -> L4 autonomy
- Basic obstacle avoidance and planning -> Joint perception and planning, learning from massive data
- No deep learning involved -> Fully deep learning stack

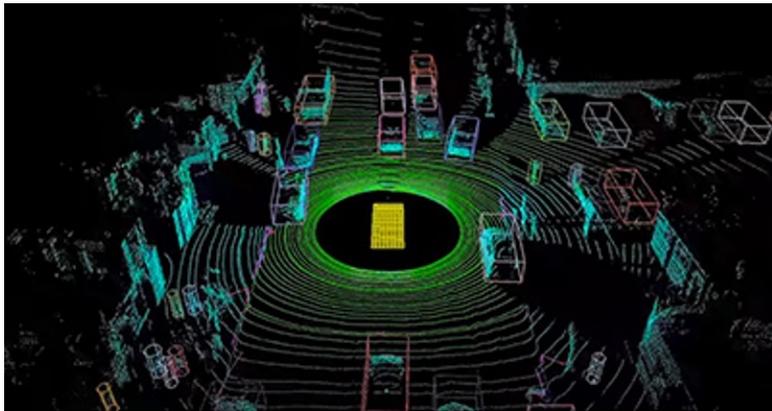


What is Missing General Embodied Intelligence?

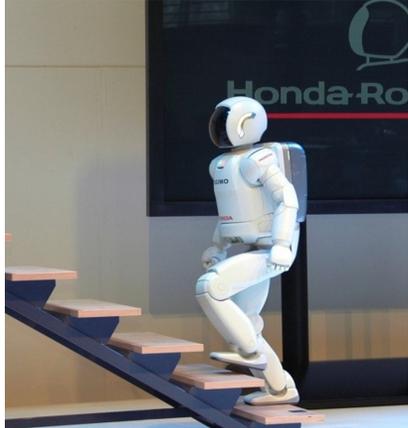
- Self-driving as an example for a more general intelligence
- It has perception, mapping, planning, and multi-agent communication.
- Where are we in general embodied intelligence?
- Perception: Exploring new environments, recognizing new signs, objects, etc.
- Learning: Learning from world modeling, future prediction, causal relations, from language instructions, etc. Learning efficiency.
- Memory and mapping: Efficient exploration of new environment without maps.
- Adaptation: Adapting to different hardware, and environments.

Better Perception

- Deep learning, semantic understanding
- Massively labeled data for training
- Sensor fusion: Camera, Multi-View, LiDAR, Radar, Motion



Better Control



ASIMO
<http://www.honda.co.jp/ASIMO/>

Honda Asimo, 2011



© Honda Motor Co., Ltd. and its subsidiaries and affiliates. All Rights Reserved



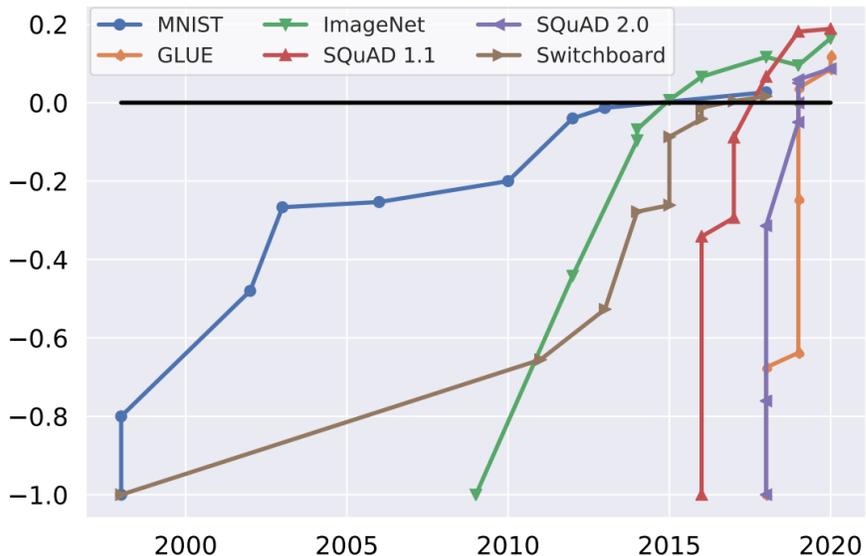
Unitree B2W, 2024



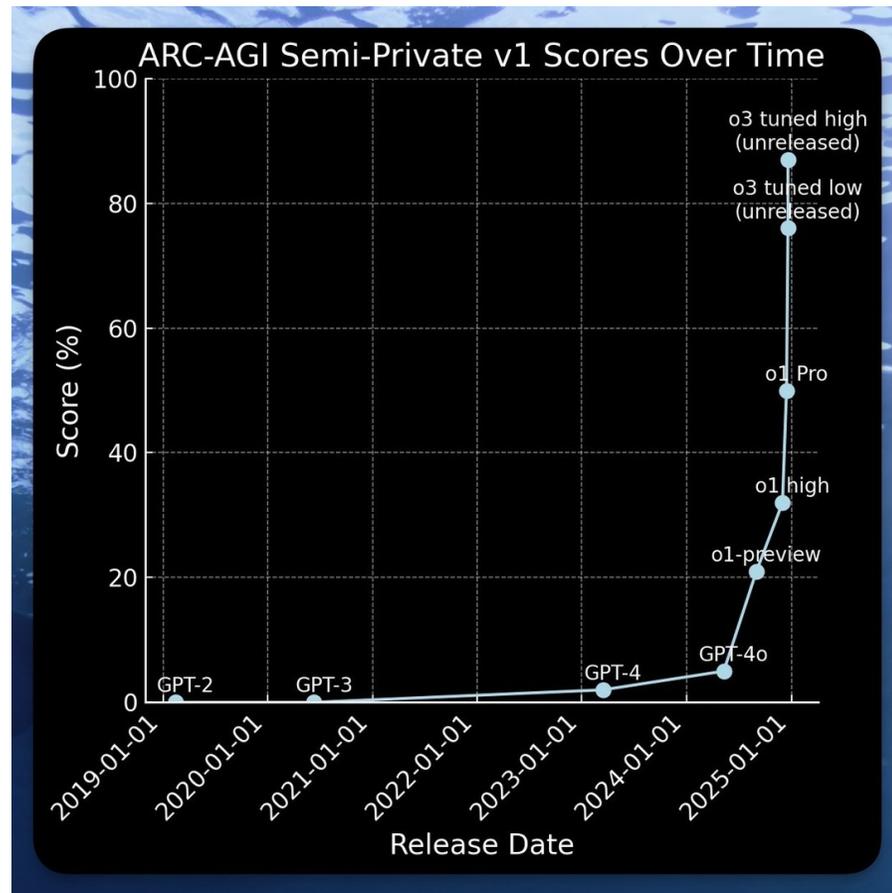
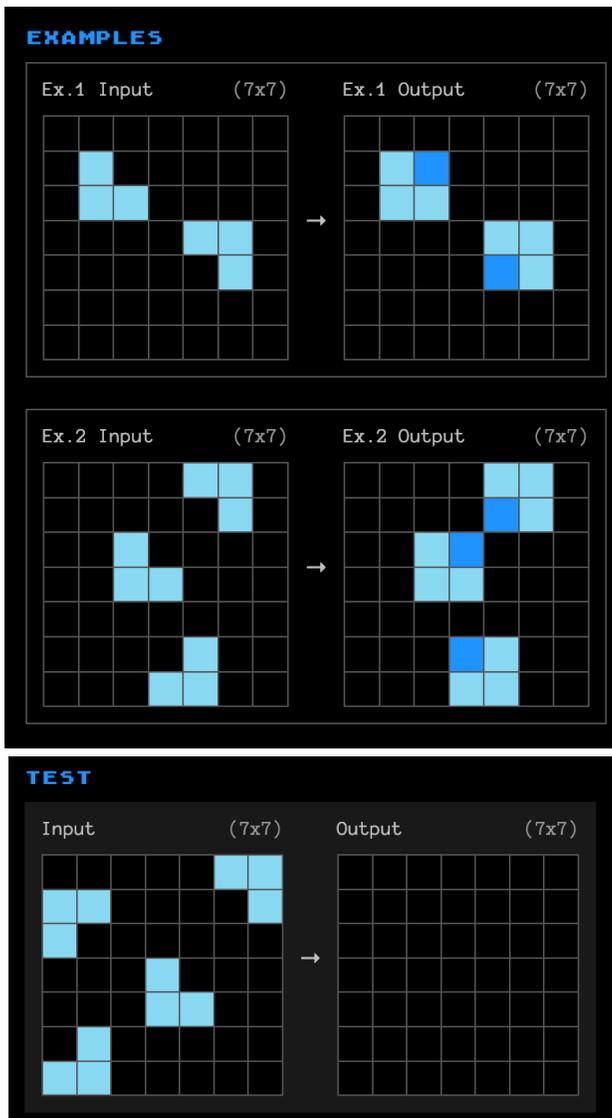
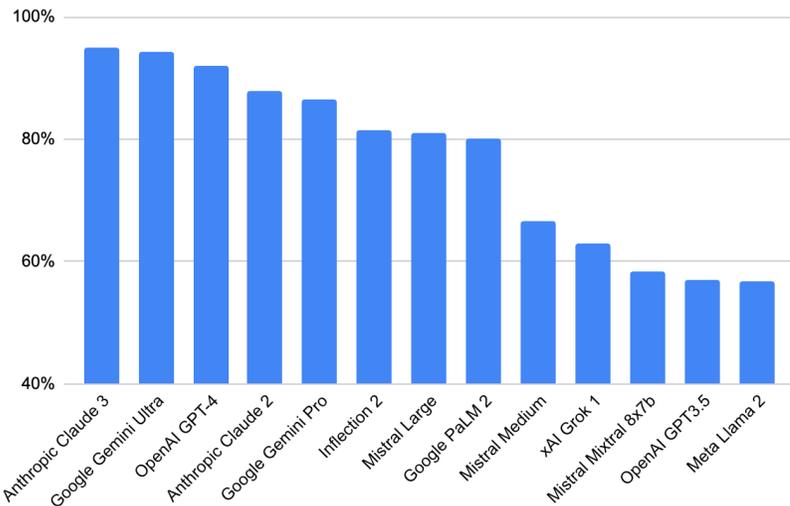
Boston Dynamics 2009-2022

Better Reasoning and Abstraction

Kiela et al. 2021



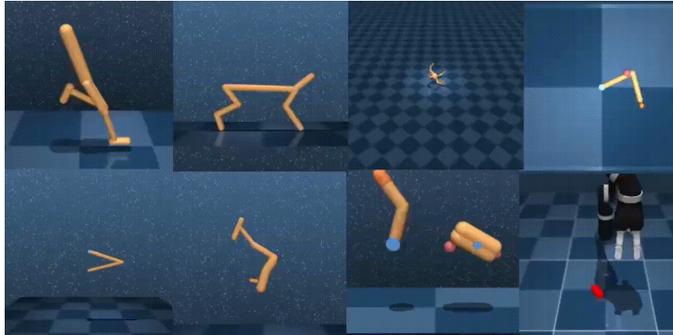
GSM 8K Accuracy



Riley Goodside, 2024



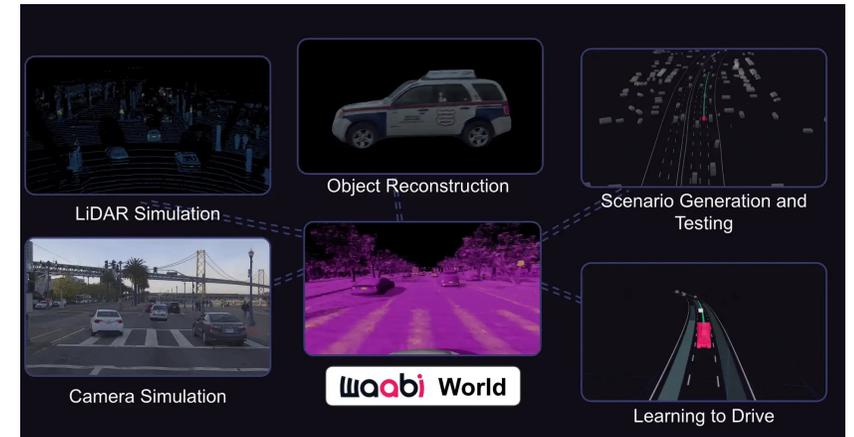
Better Simulation and Benchmarks



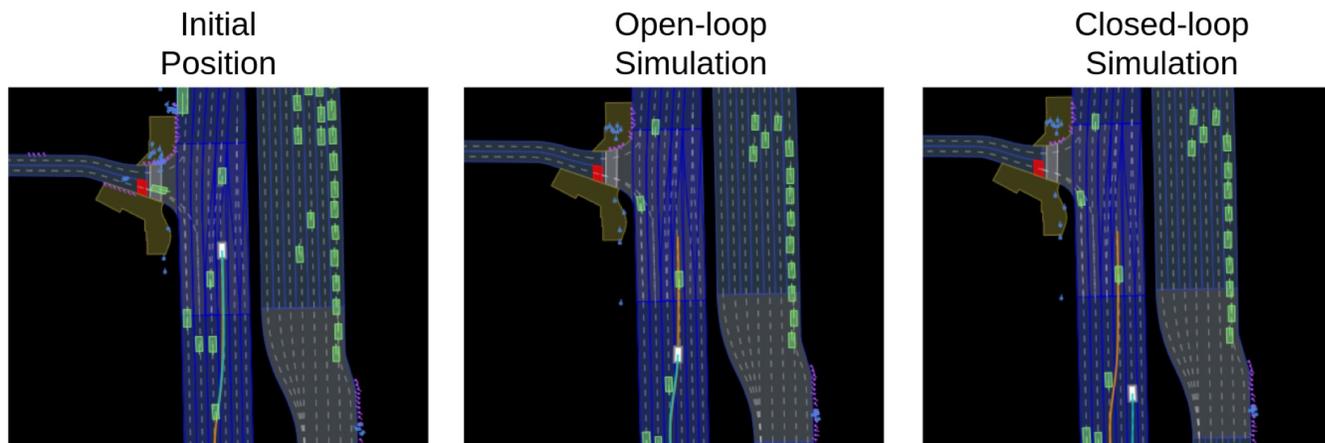
DM Control Suite



MineCraft & MineDojo



Waabi World



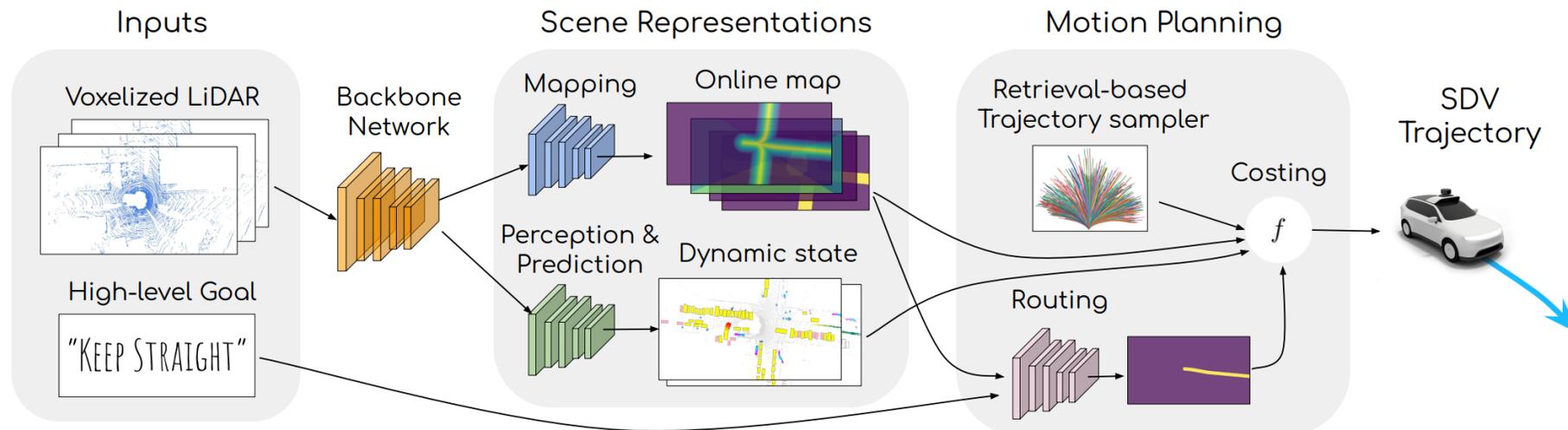
nuPlan



Nvidia Isaac Sim

End-to-End Learning

- Continuous and differentiable modules for end-to-end learning.
- We know how to optimize deep networks.
- Maintain rich information throughout decision making.
- Representations & output space modeling.



The Learning Question

- Humans can learn driving in 20 hours
- Current ML requires hundreds of millions of examples
- Ability to learn from noisy streaming data
- Ability to generalize and perform abstraction

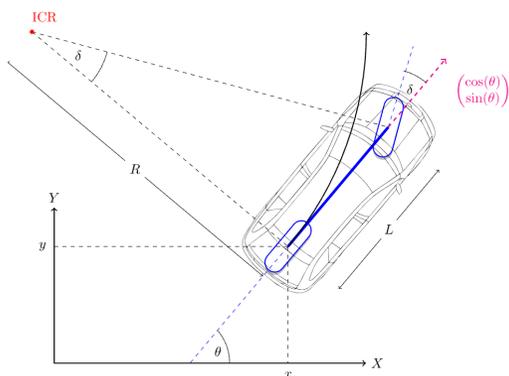
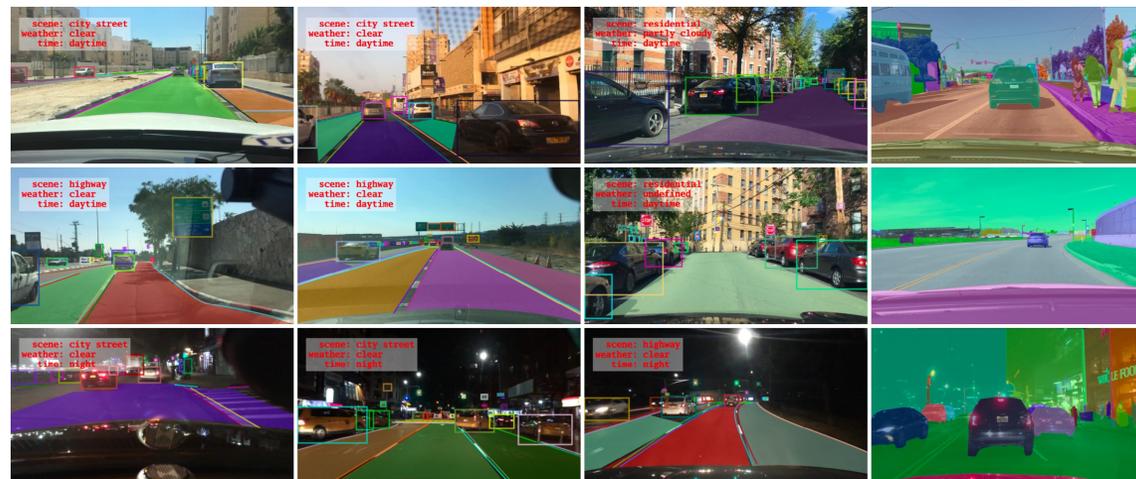
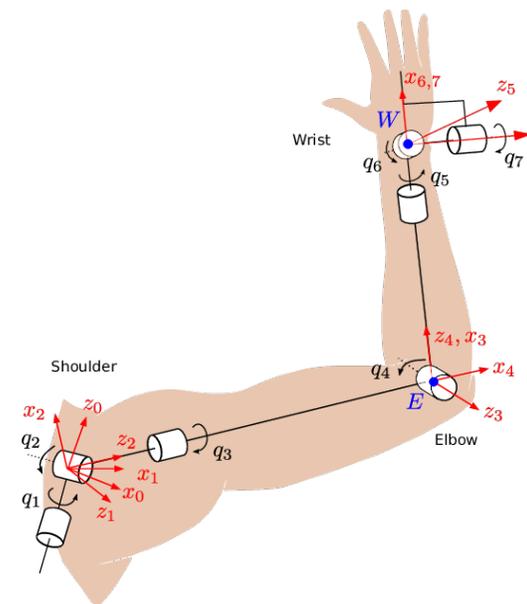


Learning Challenges?

- Learning from physical world, spatial and physical outputs
- Exploration, learning efficiency
- Imbalance, causality
- Flexibility, adaptation, continual learning
- Modular Integration
 - Sensory modality
 - Memory
 - Reasoning
 - Planning

Model Driven Vs. Data Driven

Yu et al.
2018



A collage of images and text describing the Open X-Embodiment dataset. The central text reads: "1M Episodes from 311 Scenes", "34 Research Labs across 21 Institutions", "22 Embodiments", "527 Skills", "60 Datasets", and "1,798 Attributes • 5,228 Objects • 23,486 Spatial Relations". Surrounding this text are various images of robotic arms performing tasks such as "pick anything", "pour", "sweep the green cloth to the left side of the table", "push T", "stack cups", "pick red block", "place the black bowl in the dish rack", "Cable Routing", "pick green chip bag from counter", "set the bowl to the right side of the table", and "Door Opening". The names of the datasets are listed at the bottom: "Jaco Play", "ALOHA", "Taco Play", "Bridge", and "RT-1".

Open X-Embodiment, 2024



Modular Vs. End-to-End

Modular

End-to-End

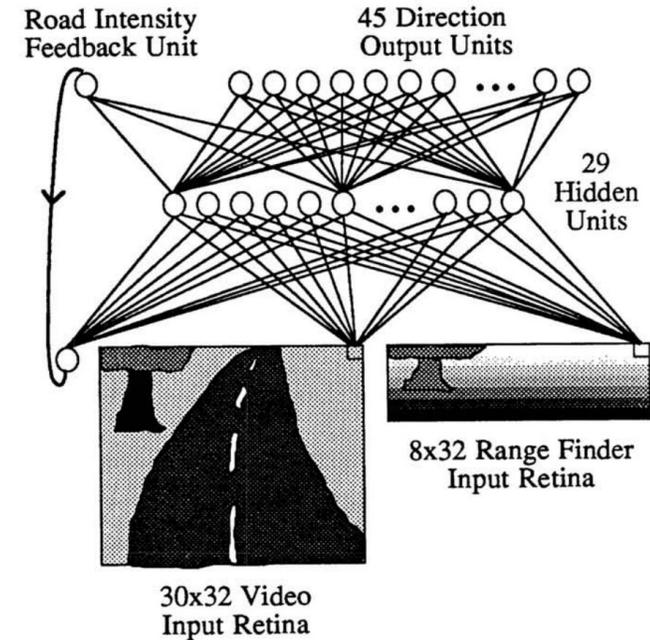
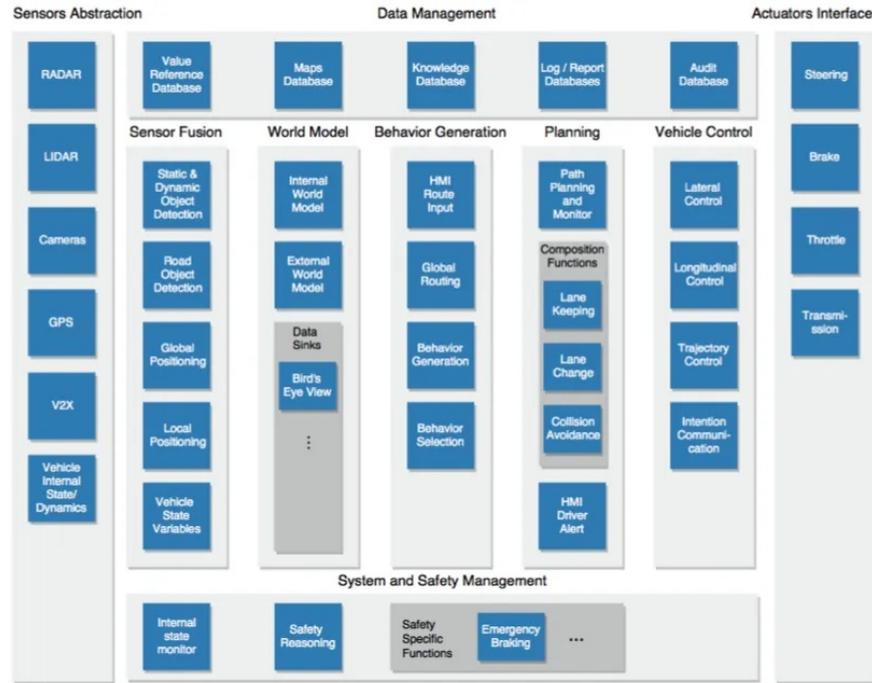
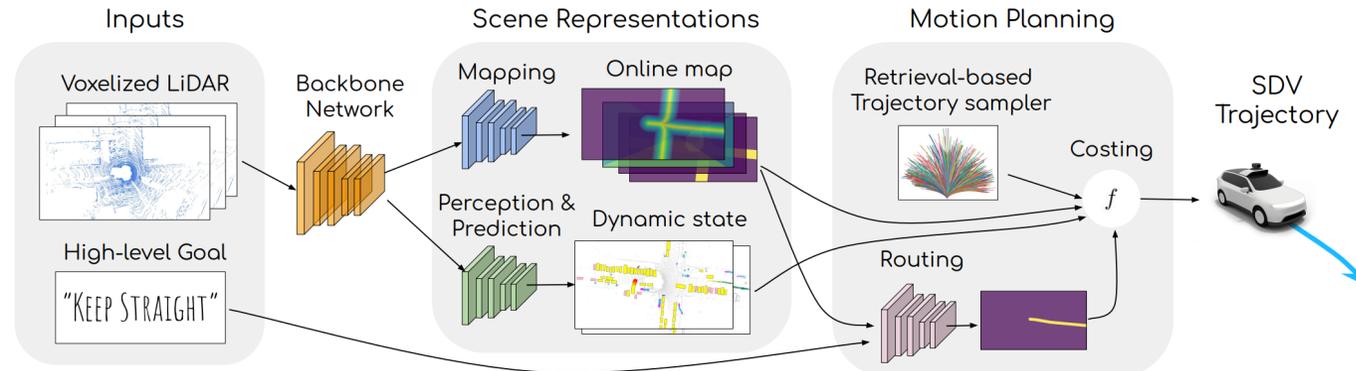
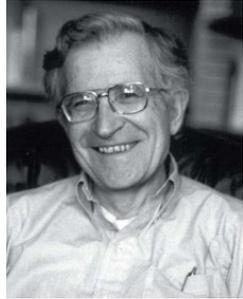


Figure 1: ALVINN Architecture

Modular
and
End-to-End
Differentiable



Nature Vs. Nurture



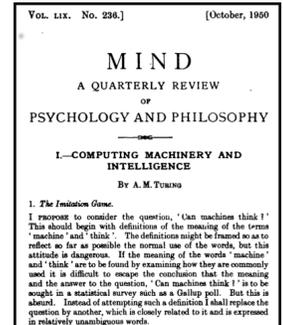
Noam Chomsky

There's an obvious answer to that: the knowledge is built in. You and I can learn English, as well as any other language, with all its richness because we are designed to learn languages based upon a common set of principles, which we may call universal grammar.



Yarek Waszul

In fact, if someone came along and said that a bird embryo is somehow “trained” to grow wings, people would just laugh.

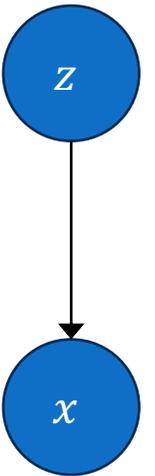


COMPUTING MACHINERY AND INTELLIGENCE (Turing, 1950)

Instead of trying to produce a programme to simulate the adult mind, why not rather try to produce one which simulates the child's? If this were then subjected to an appropriate course of education one would obtain the adult brain. Presumably the child brain is something like a notebook as one buys it from the stationer's. Rather little mechanism, and lots of blank sheets. (Mechanism and writing are from our point of view almost synonymous.) Our hope is that there is so little mechanism in the child brain that something like it can be easily programmed.

Why Do We Need Learning in Real-World Agents?

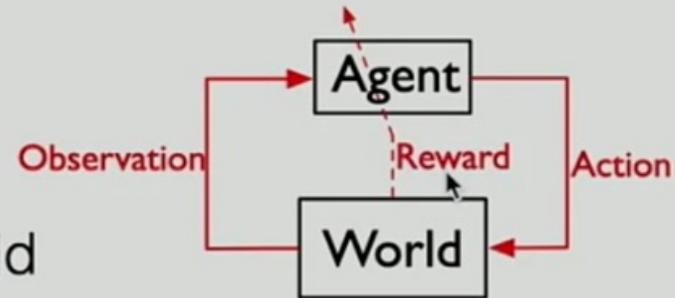
- Symbolic Argument: You can represent infinite variations with finite length description of an abstract system. The underlying system remains the same.
- Big World Hypothesis: We always need learning in exploring new environments. There will always be something unknown. There will always be room for improvement. There won't be enough capacity to store all existing knowledge.
- Given resource constraints, we should only adapt to the current moment and generalize to the future.



Rich Sutton on the Quest of Continual Learning

General (domain independent)

— contains nothing specific to any world



Experiential

— grows from **runtime experience**, not from a special training phase

Able to discover its own abstractions in state and time

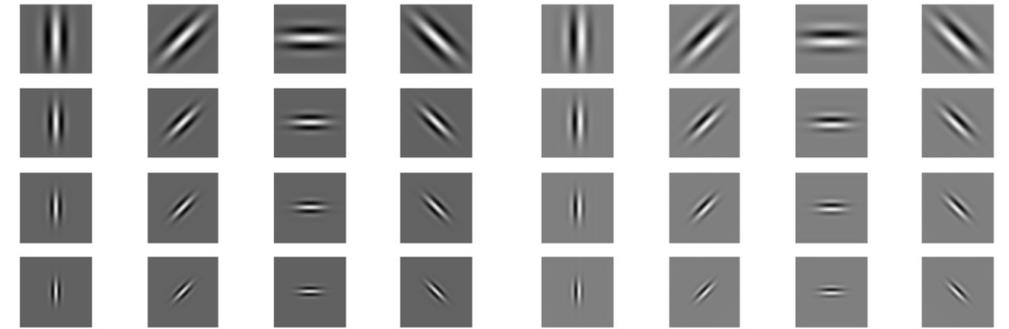
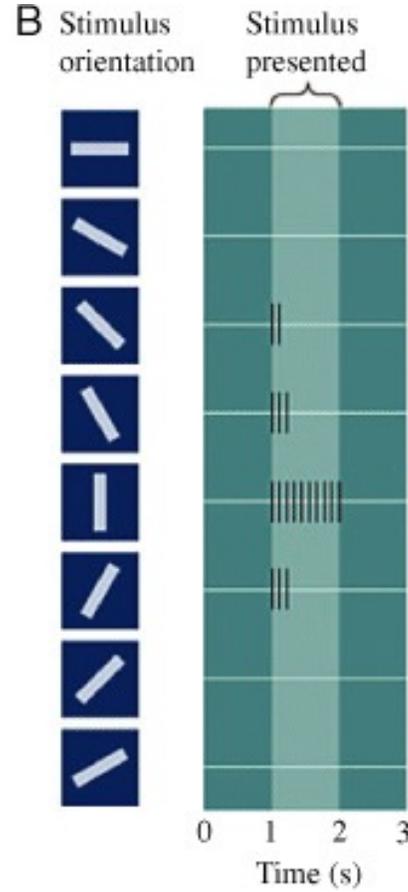
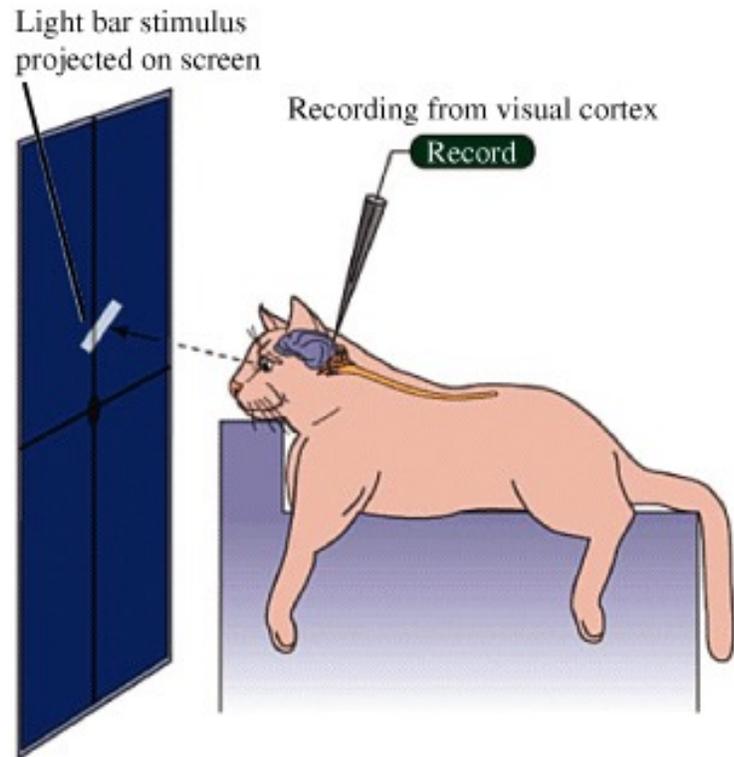
— in an open-ended way (limited only by computational resources)

Insights from Biology

- Animals and humans are examples of real-world embodied intelligence.
- There is evidence that higher cognitive functions in animals require more learning and adaptation.
- Evidence on visual acuity and perception in mammals.
- Survival needs:
 - Changing reward functions in the wild.
 - Problem solving skills.

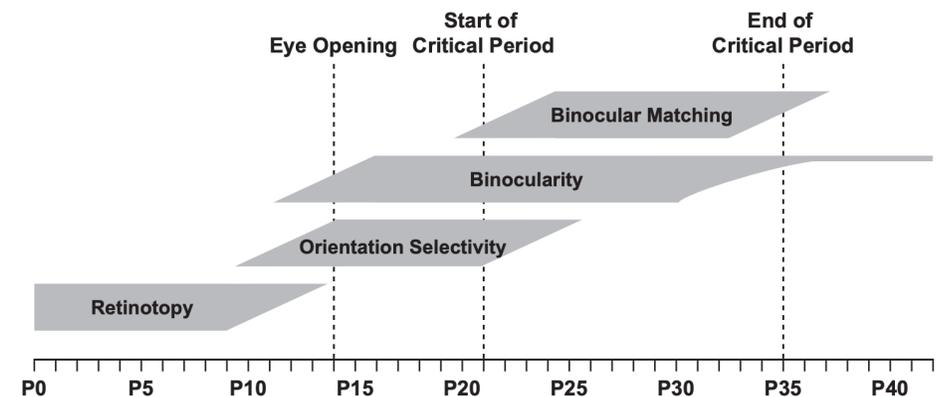
Hubel and Wiesel's Experiments

A Experimental setup



Simple Cells, Gabor Filters

Wiesel, T.N., and Hubel, D.H. (1963)



Espinosa and Stryker (2012)

Human Developmental Periods

Sensorimotor learning

Simple Reflexes (birth-1 month)

Infants use reflexes such as rooting, sucking, following moving objects with the eyes, and grasping objects. (For example: Infant closes their hand when a toy touches their palm.)

Primary Circular Reactions (1-4 months)

A primary circular reaction is when an infant tries to reproduce an event that happened by accident because they find it to be pleasurable. (For example: Intentionally mouthing a toy bunny.)

Secondary Circular Reactions (4-8 months)

Child becomes more focused on the world and begins to intentionally repeat an action in order to trigger an environmental response. (For example: purposefully picking up a pacifier to put it in their mouth.)

Coordination Of Secondary Circular Reactions (8-12 months)

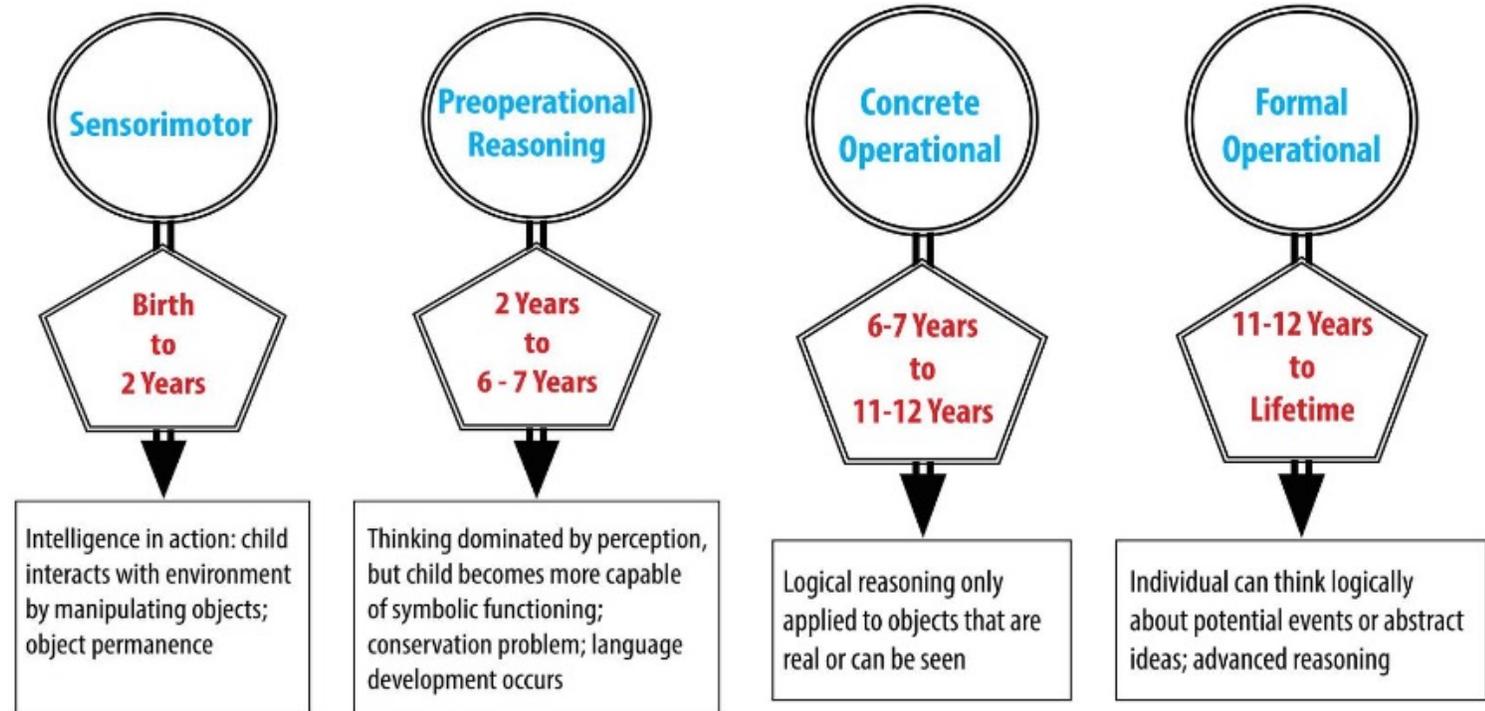
Child acts intentionally and follows steps to achieve goals. Child begin to do things intentionally and understands object permanence. (For example: Child will push one toy aside to get to a second toy partially concealed underneath.)

Tertiary Circular Reactions (12-18 months)

Child discovers new means to meet goals and begins to modify earlier behaviors to meet existing needs. Piaget described children in this stage as "young scientists". (For example: Child repeatedly drops/throws a set of plastic keys and observes how they move through space.)

Internalization of schemas (18-24 months)

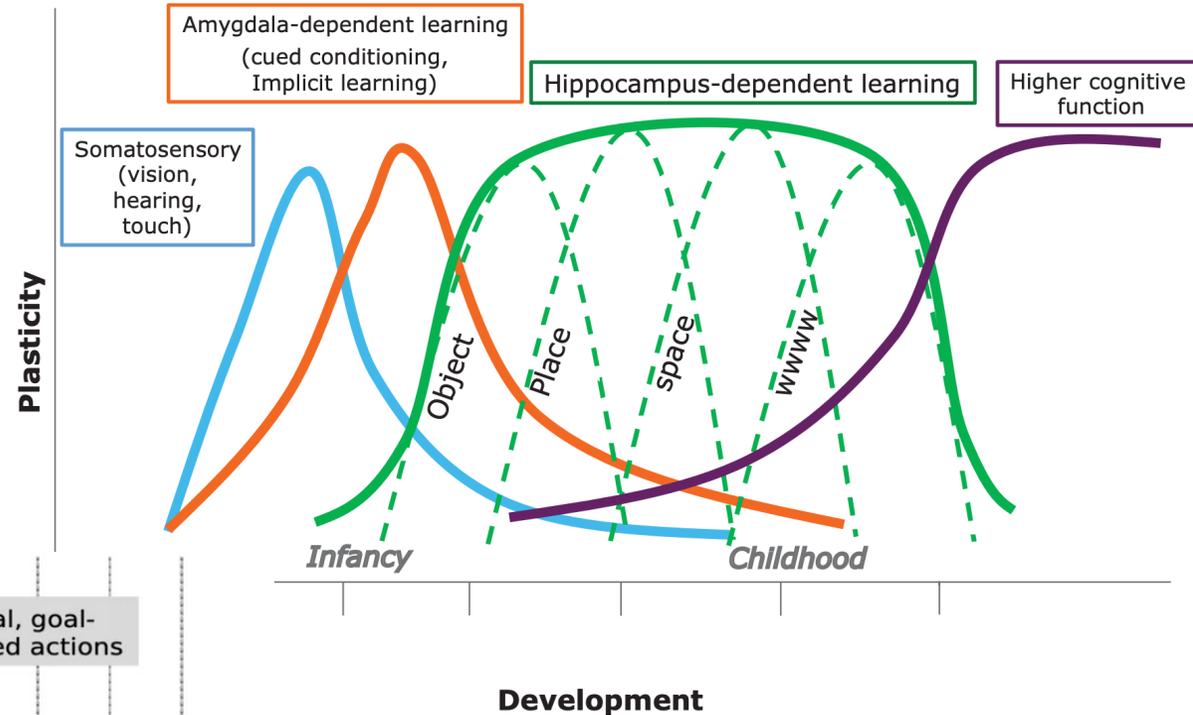
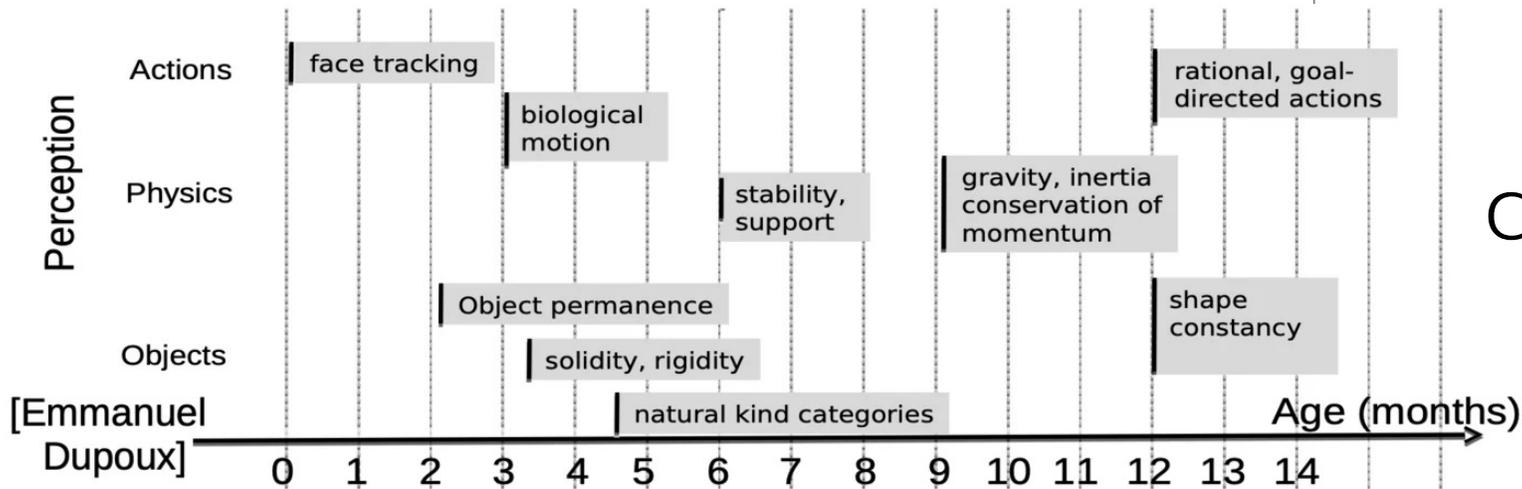
Child begins to use symbols and form mental representations. The beginnings of insight and creativity are associated with this stage. (For example: Child pushes a chair across the kitchen and climbs up on it to reach a cookie on the counter.)



Piaget's Theory of Cognitive Development

Insights from the Brain

- Perception and motion
- Low-level to high-level representation
- Hippocampus and memory
- Abstraction



Critical Periods of Development

Takeaways

- General embodied and agentic intelligence requires learning in embodied and agentic environments.
- A learning system requires both model inductive biases and data.
- There is a wide space in between pure modular vs. pure model-free (end-to-end).
- Still a huge gap in terms of learning efficiency, robustness, and flexibility.

Logistics

Grading

- In-Class Participation (10%)
- Paper Review (15%)
- Paper Presentation (30%)
- Project (45%):
 - Project Proposal (10%)
 - Project Meetings (10%)
 - Report (15%)
 - Presentation (10%)



[Course Syllabus](#)

CampusWire



3052

<https://campuswire.com/p/GEFF369DC>

People who sign up using this link will get 'student' access to the class.

Gradescope

- For submitting weekly paper reviews
- Entry code: ZJPZ3J

Introduce Yourself

- Share an introductory post on Campuswire.
 - Please update your profile picture.
 - Say what you are interested in.
 - This also helps you know each other and find a team partner on the course project.

In-Class Participation

- You get marks for asking good questions in Q&A periods of in-class student presentations, guest lectures, etc.
- Announce your name before asking a question.
- 10% worth of marks

Paper Reviews

- 15% of the total mark
- Select a paper from the suggested reading list (recent only), or find a recent paper of your interest (with approval)
- One topic each week
- Due before the lectures: W3 – W10

Paper Reviews

- Please follow the template in the Syllabus.
- Components: Summary, Strengths, Weaknesses,
- New this year: Paper idea illustration.

Topic Presentation

- Sign Up: Students have to sign up for a slot by Week 3.
- Calendar: Week 7 – 13 (Tentatively).
- Approximately 4 students present on each topic.
- Each student will conduct a 15-20 min presentation on 1-2 designated recent papers including necessary backgrounds.

Course Project

- 45% of the total mark
 - Project Proposal (10%)
 - Project Meetings (10%)
 - Report (15%)
 - Presentation (10%)
- Project Meetings: Pair with one of the TAs and meet biweekly.
- Week 14 + Week 15 Project Presentation. Week 14: 2% Bonus.
- Target: Top-tier conference workshops and above.
- Key dates in syllabus.

Project Report Template

- Page limit: 8
- Including Table and Figures
- Does not include Appendix.
- What is enough?
- When you need to squeeze white space and cut down content.

DS-GA 3001 Course Project Report Template and Instructions

Student Name 1
Affiliation
Address
email

Student Name 2
Affiliation
Address
email

Team [Team ID]

Abstract

This document outlines the instructions for the course project of DS-GA 3001. Your project proposal and final report should use this space for abstract.

1 Course Project

The course project constitutes 45% of your overall grade. The goal of the final project is to let students develop hands-on skills of implementing embodied learning systems for concrete real-world tasks such as toy games, long-form egocentric video understanding, self-driving simulation, indoor navigation, and robotic manipulation.

Direction 1: End-to-end representation learning for perception and planning

- Existing end-to-end learning-based planning frameworks are mostly focused on supervised learning of labeled objects and human demonstrations.
- Demonstrations and rewards are forms of labels.
- How do we achieve label-efficient learning and exploration through self-supervision?
- Is planning and action necessary for a label-efficient algorithm for perception?
- Explore the full spectrum from end-to-end learning to modular designs.

Direction 2: Foundation models for spatial and physical intelligence

- Foundation models are trained with discrete tokens and are less familiar with the 3D world to perform exact perception and planning.
- Augment pretrained foundation models with the ability to perceive and plan under physical precision in embodied environments.
- How do we enhance robustness in real-world environments?
- Can be synthetic/realistic, 2D/3D environments.
- Can models with geometric designs beat generic foundation models in terms of learning efficiency?

Direction 3: Continual learning for embodied intelligence

- How do we apply continual learning algorithms to embodied tasks?
- Skill learning, open world learning
- Memory design, retrieval augmentation, continuous finetuning
- Incremental learning with experience/action abstraction
- Defining subgoals and subproblems
- Replay with spatial and physical constraints
- Actively choosing learning objectives

Other Directions?

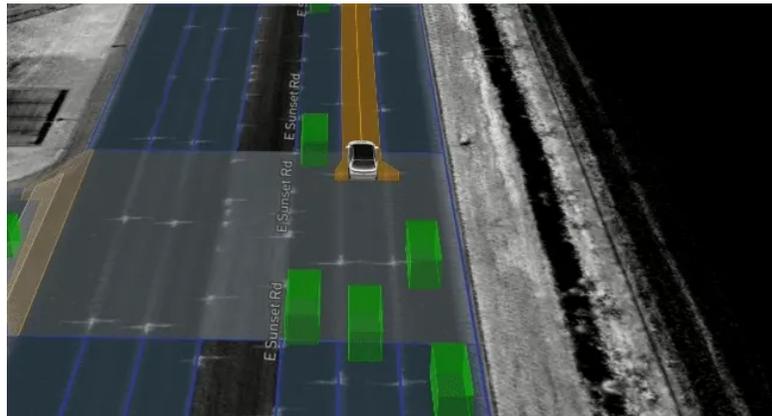
- You are allowed to form your own research ideas.
- Talk to me early in the semester and get approval.

Embodied Environments

- You **must** demonstrate your project in an embodied environment.
- You can focus on one aspect of the algorithm. No need for a full stack.
- Your TAs will showcase demos on some exemplar environments.



Habitat indoor home



NuPlan self-driving



Ego-Exo4D Egocentric Videos

GenAI Policy

- AI may not be used in weekly paper reviews and paper presentations (except AI illustrations).
- AI may be used towards coding assistance and report writing assistance in the course project.
- The use of AI can still impact the grade if the report contains poor writings and non-factual statements.

Content

- Introduction, Brief History
- Deep Learning and Structured Outputs
- 3D Vision and Mapping
- Latent Representation Learning and Object Discovery
- World Models and Forecasting
- End-to-End Planning
- Continual Learning
- Few-Shot Learning
- Foundation Model Agents

What's Next

- Next Week:
 - Deep Learning with Structured Outputs
 - Tutorial on HPC (Cloud Burst) by Ying