

# DSC250: Advanced Data Mining

## Text (Multi-modal) Mining: Large Language (Multi-modal) Models

**Zhiting Hu**

Lecture 2, Jan 9, 2025

# Foundation Models: Large Language Models (LLMs)



What are a couple of meals I can make with this?



GPT-4

Based on the items visible in the fridge, here are a couple of meal ideas:

- 1) Yogurt parfait: Layer yogurt, strawberries, and blueberries in a bowl or jar. Top with honey or nuts for added texture and flavor.
- 2) Carrot and hummus wrap: Spread hummus on a tortilla or wrap...

# Foundation Models: Large Vision Models (LVMs)

**Video generated by  
OpenAI Sora**



# Large Language Models



# Natural Language Processing (NLP): Before 2017

Automated understanding and generation of natural language

Core NLP tasks handled by respective machine learning models, e.g.,:

## Named Entity Recognition

Adam Driver was born in San Diego , California , on November 19 , 1983 .

The diagram illustrates Named Entity Recognition (NER) for the sentence "Adam Driver was born in San Diego , California , on November 19 , 1983 .". Brackets above the text group the words into four categories, each labeled in a colored box: "PERSON" (orange) for "Adam Driver", "CITY" (light blue) for "San Diego", "STATE\_OR\_PROVINCE" (dark blue) for "California", and "DATE" (light green) for "November 19 , 1983".

## Sentiment Analysis

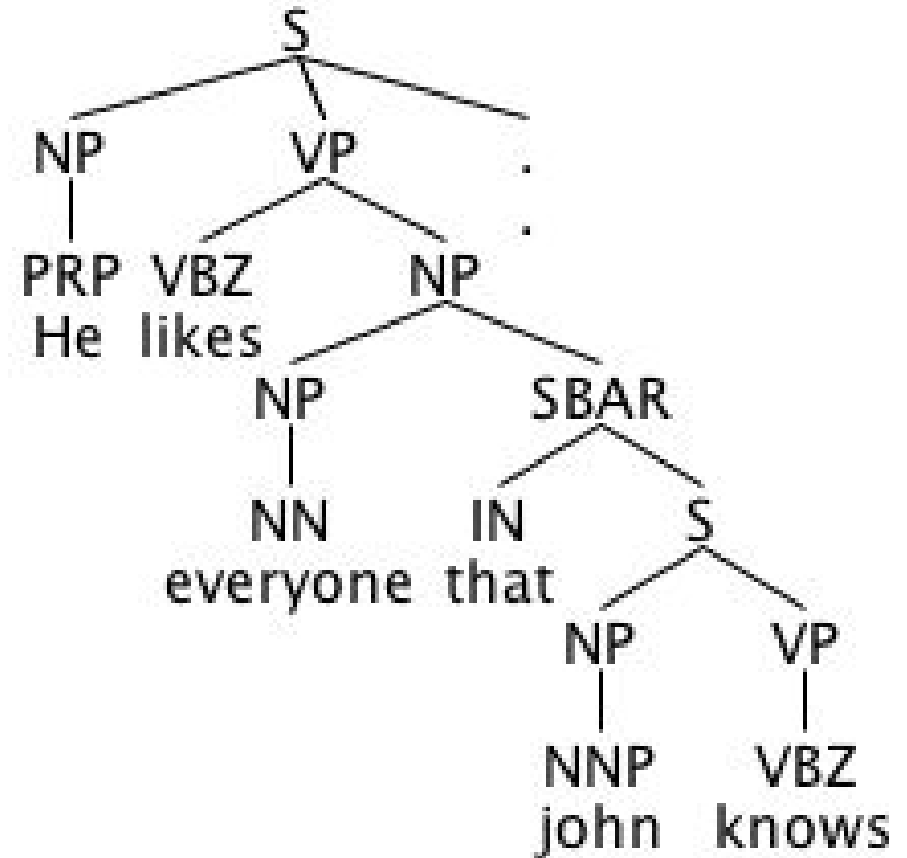
POSITIVE

There are slow and repetitive parts , but the movie has just enough spice to keep it interesting .

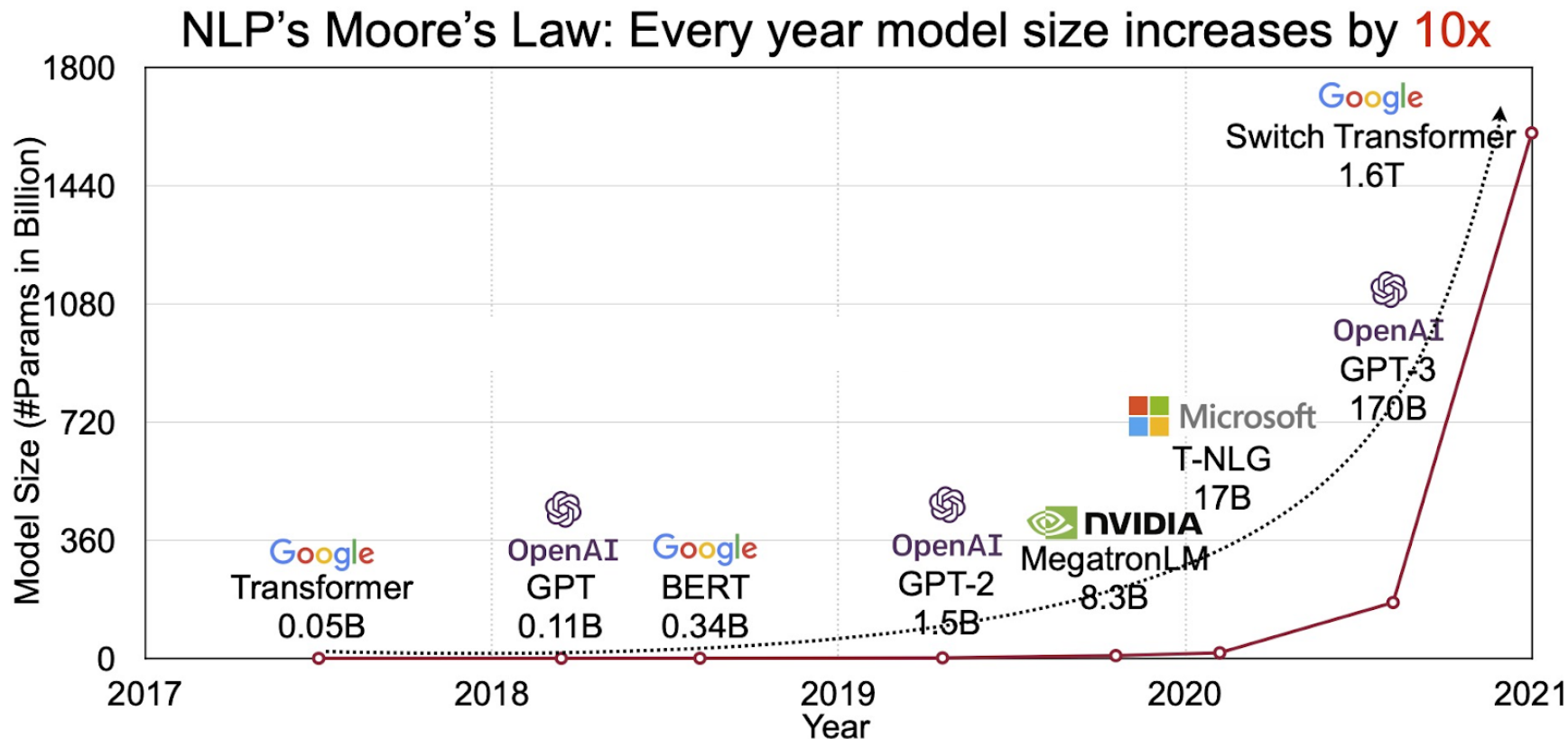
# Natural Language Processing (NLP): Before 2017

Automated understanding and generation of natural language

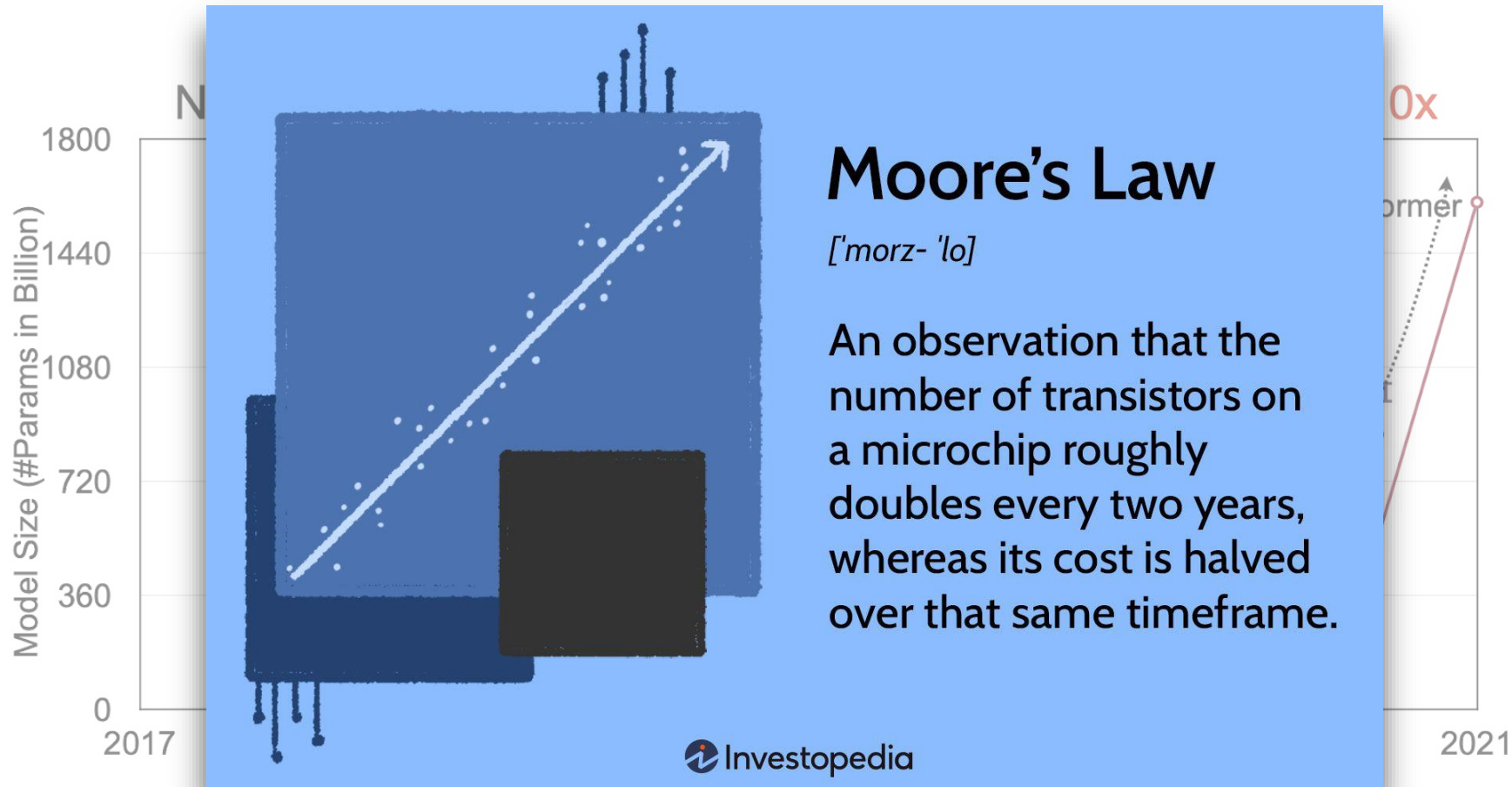
Hand annotation of linguistic structures  
(e.g., the Penn Treebank, 1990s)



# NLP breakthrough with large language models, since 2017



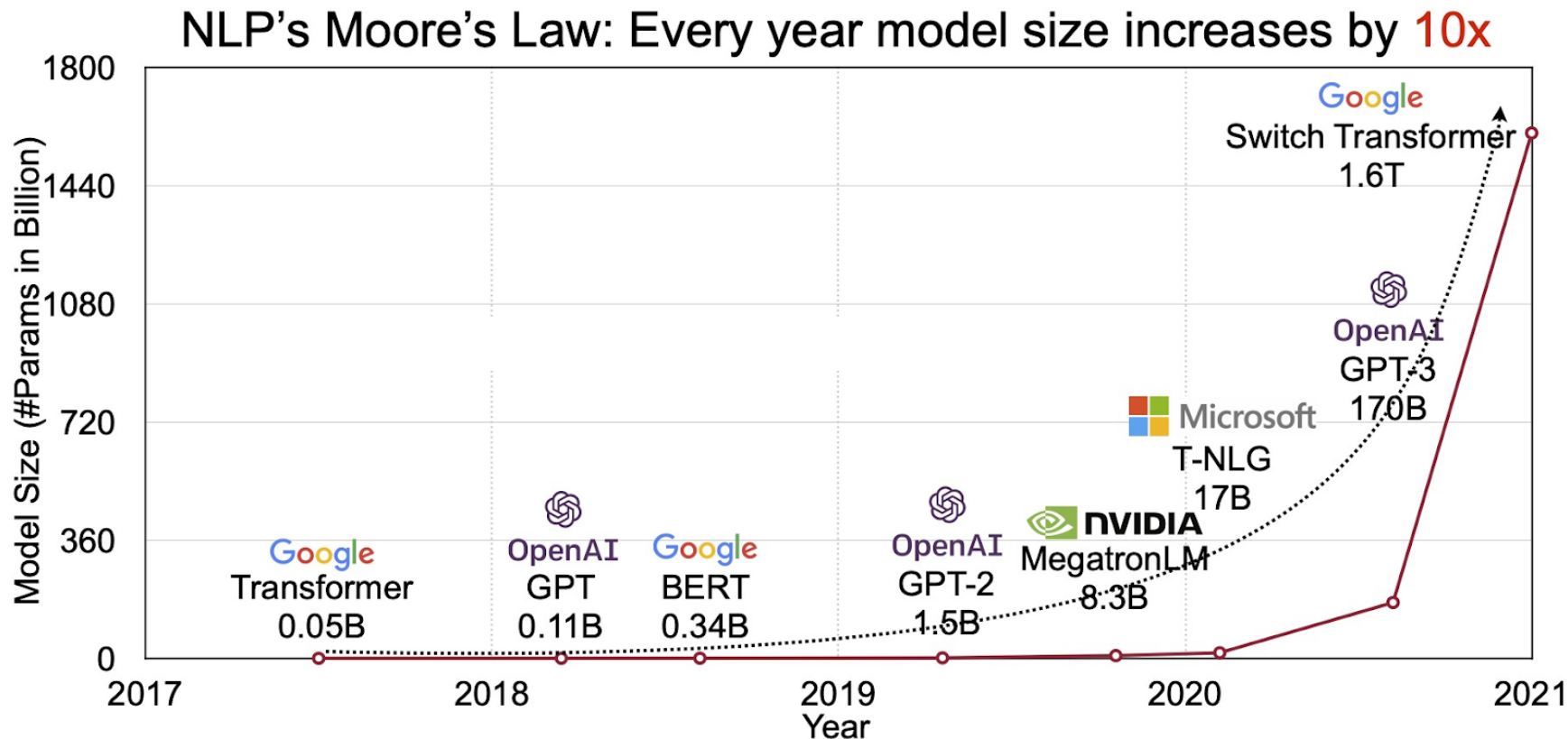
# NLP breakthrough with large language models, since 2017



microchip industry

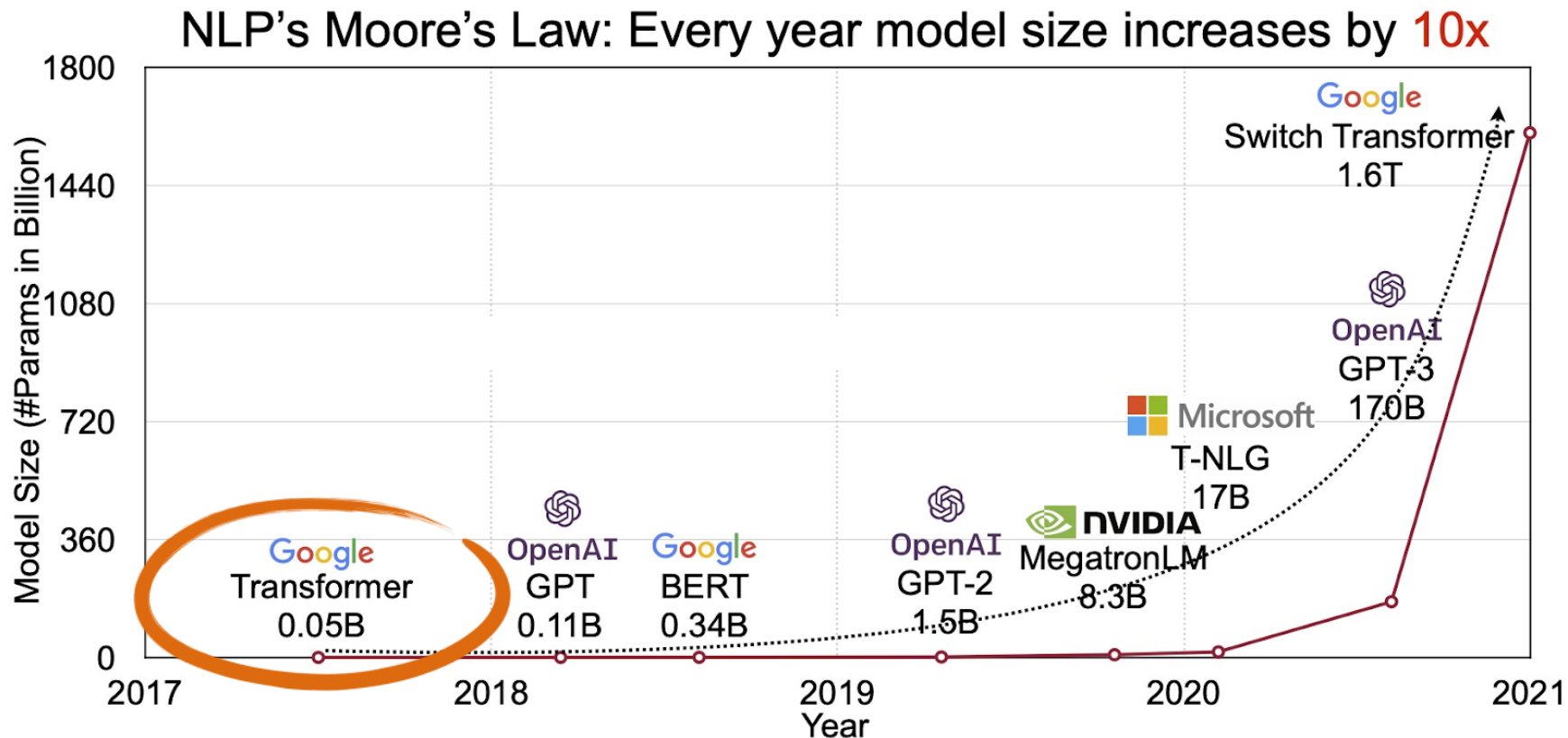


# NLP breakthrough with large language models, since 2017

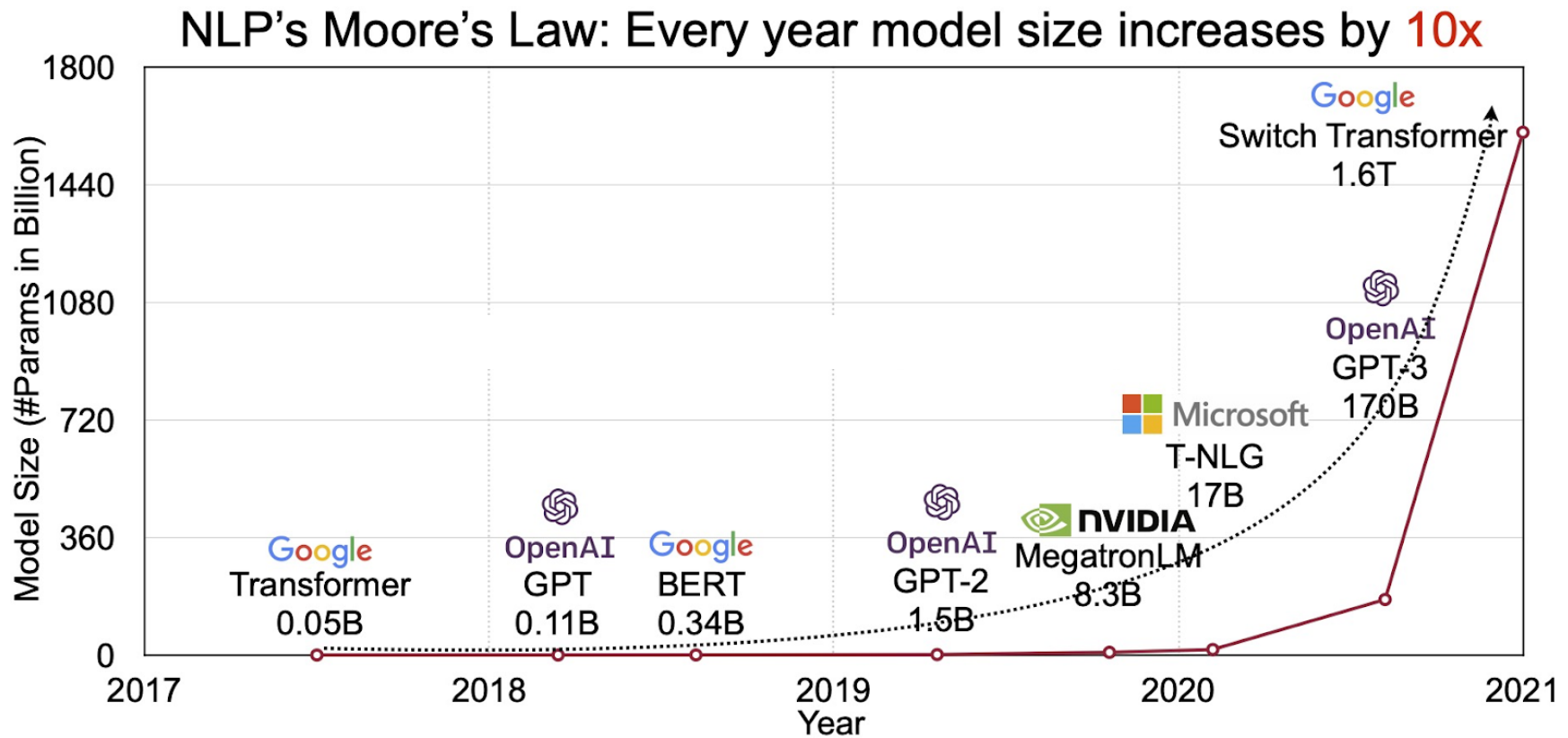




# NLP breakthrough with large language models, since 2017



# NLP breakthrough with large language models, since 2017



**ChatGPT**

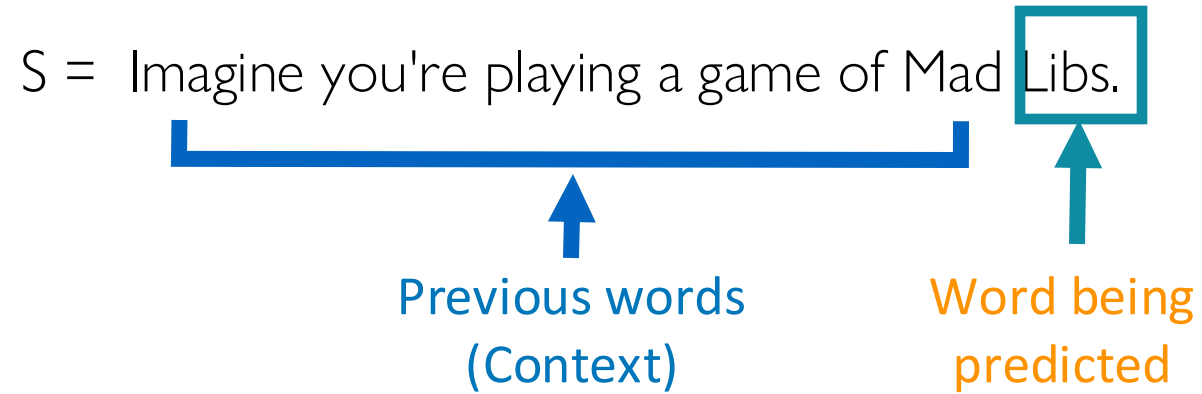
2022

# What is a language model?

# Language Model 101

$S =$  Imagine you're playing a game of Mad Libs.

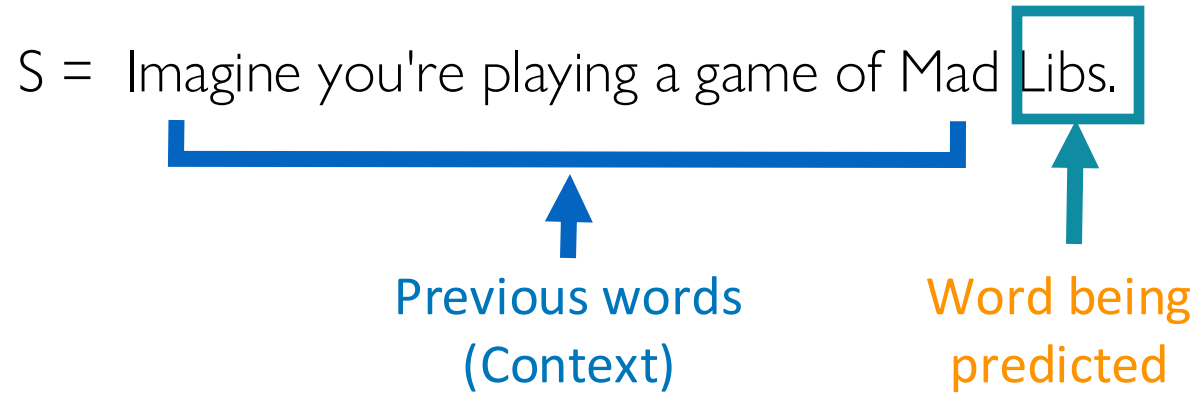
# Language Model 101



*Next word prediction*



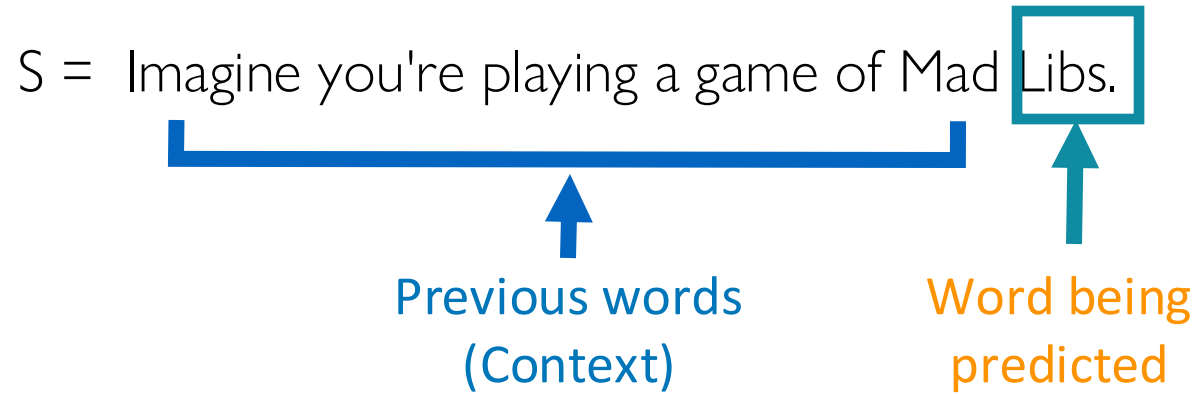
# Language Model 101



*Next word prediction*

$$P(w_i | w_1, \dots, w_{i-1})$$

# Language Model 101



*Next word prediction*

$$P(w_i | w_1, \dots, w_{i-1})$$

I \_\_\_\_\_

# Language Model 101

$$P(w_i | w_1, \dots, w_{i-1})$$

## **Implementations (model architecture):**

N-grams

Recurrent Neural Networks (RNNs)

Transformer

...

# Language Model 101

$$P(w_i | w_1, \dots, w_{i-1})$$

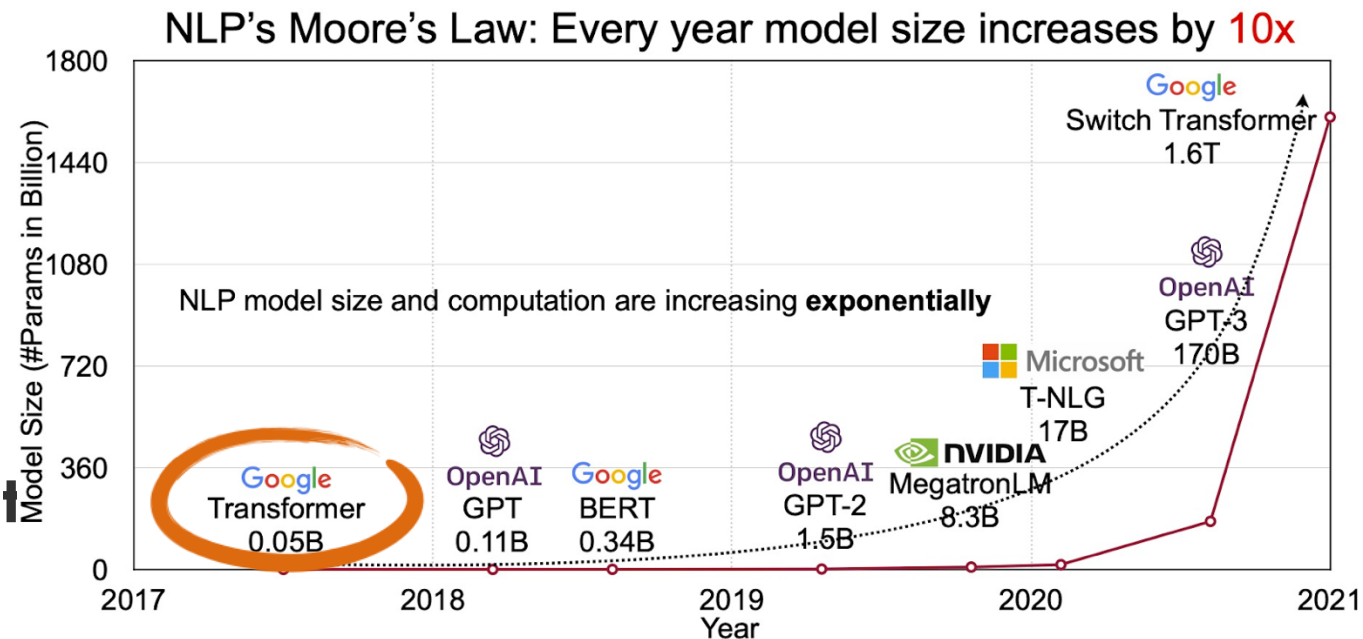
## Implementations (model architecture)

N-grams

Recurrent Neural Networks (RNNs)

Transformer

...



# Language Model 101: Transformer

$$P(w_i | w_1, \dots, w_{i-1})$$

---

## Attention Is All You Need

---

<b>Ashish Vaswani*</b> Google Brain avaswani@google.com	<b>Noam Shazeer*</b> Google Brain noam@google.com	<b>Niki Parmar*</b> Google Research nikip@google.com	<b>Jakob Uszkoreit*</b> Google Research usz@google.com
---	---	--	--

<b>Llion Jones*</b> Google Research llion@google.com	<b>Aidan N. Gomez* †</b> University of Toronto aidan@cs.toronto.edu	<b>Łukasz Kaiser*</b> Google Brain lukaszkaizer@google.com
--	---	--

**Illia Polosukhin\* ‡**  
illia.polosukhin@gmail.com

2017



# Language Model 101: Transformer

$$P(w_i | w_1, \dots, w_{i-1})$$

## Attention Is All You Need

Ashish Vaswani\*  
Google Brain  
avaswani@google.com

Noam Shazeer\*  
Google Brain  
noam@google.com

Niki Parmar  
Google Research  
nikip@google.com

Llion Jones\*  
Google Research  
llion@google.com

Aidan N. Gomez\*<sup>†</sup>  
University of Toronto  
aidan@cs.toronto.edu

lukas

Illia Polosukhin\*<sup>‡</sup>  
illia.polosukhin@gmail.com

The children were hungry. They **looked out** the window. Where was their mother? She walked into the house. The children **ran over** to her. "Mama, we're so **hungry**," they both said. She said **lunch** was coming. She walked into the **kitchen**. She opened a can of **chicken soup**. She **poured** the soup into a **pot**. She added water. She put the pot on the **stove**. She made two **peanut butter** and **jelly sandwiches**. She sliced an apple. The soup was hot. She poured it into two bowls. She put the sandwiches on two

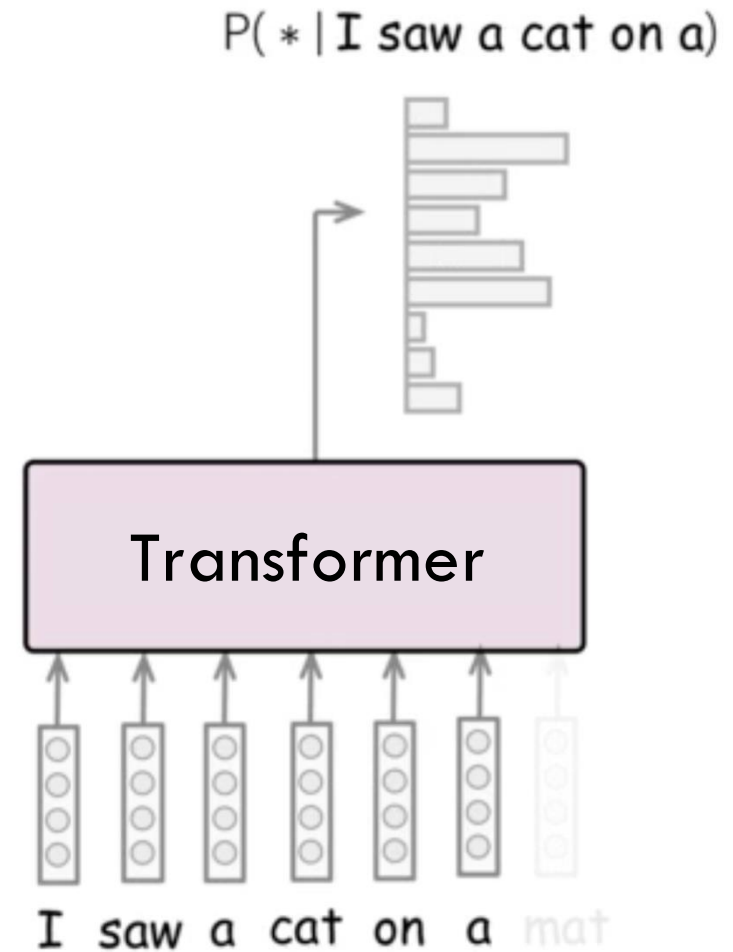
**plates**. She put apple slices on each plate. She put the **bowls** and plates on the table. The children ran to the table. "Thank you, mommy!" they said. Then they started eating. The cat and the dog watched them eat.



2017

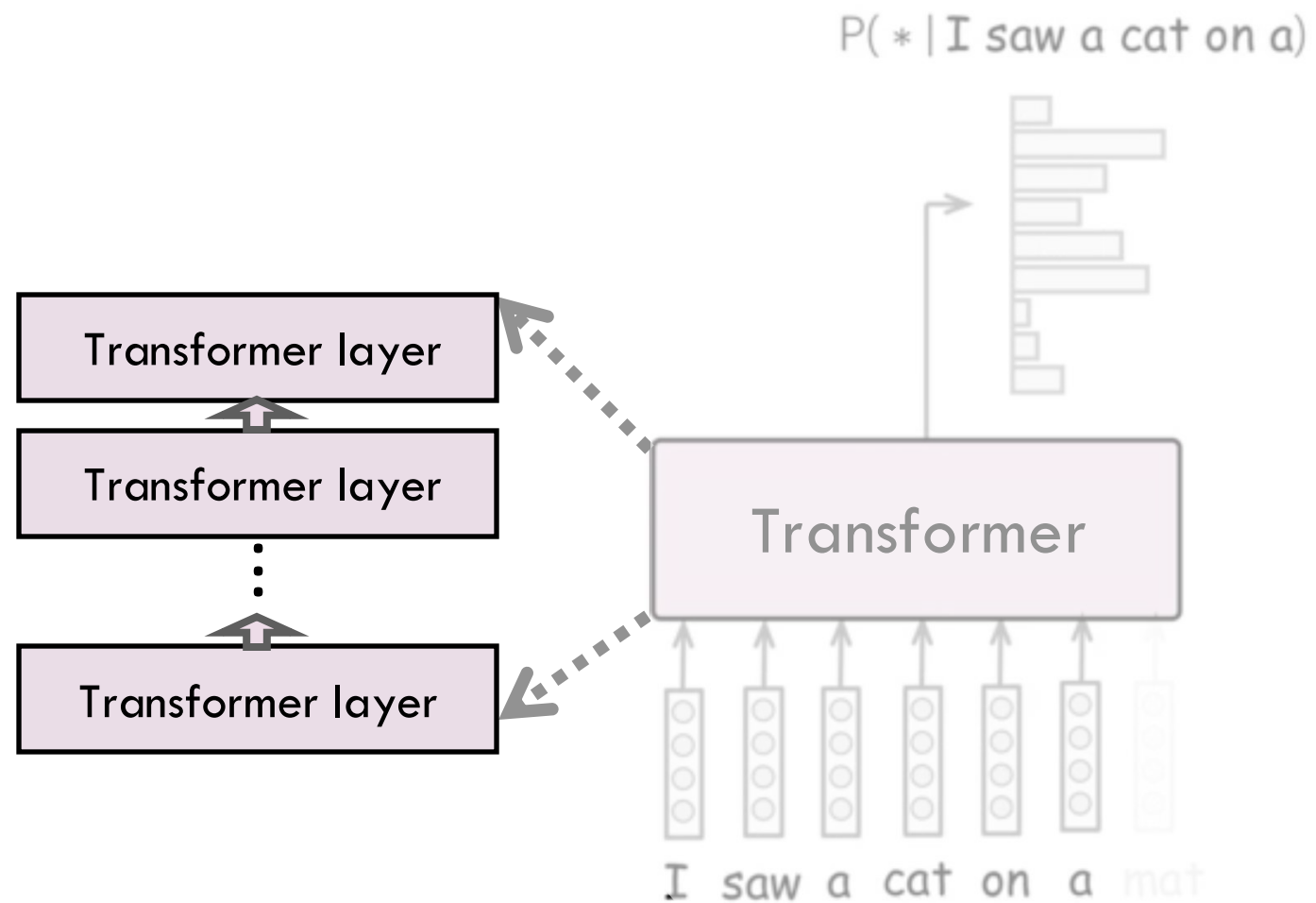
# Language Model 101: Transformer

$$P(w_i | w_1, \dots, w_{i-1})$$



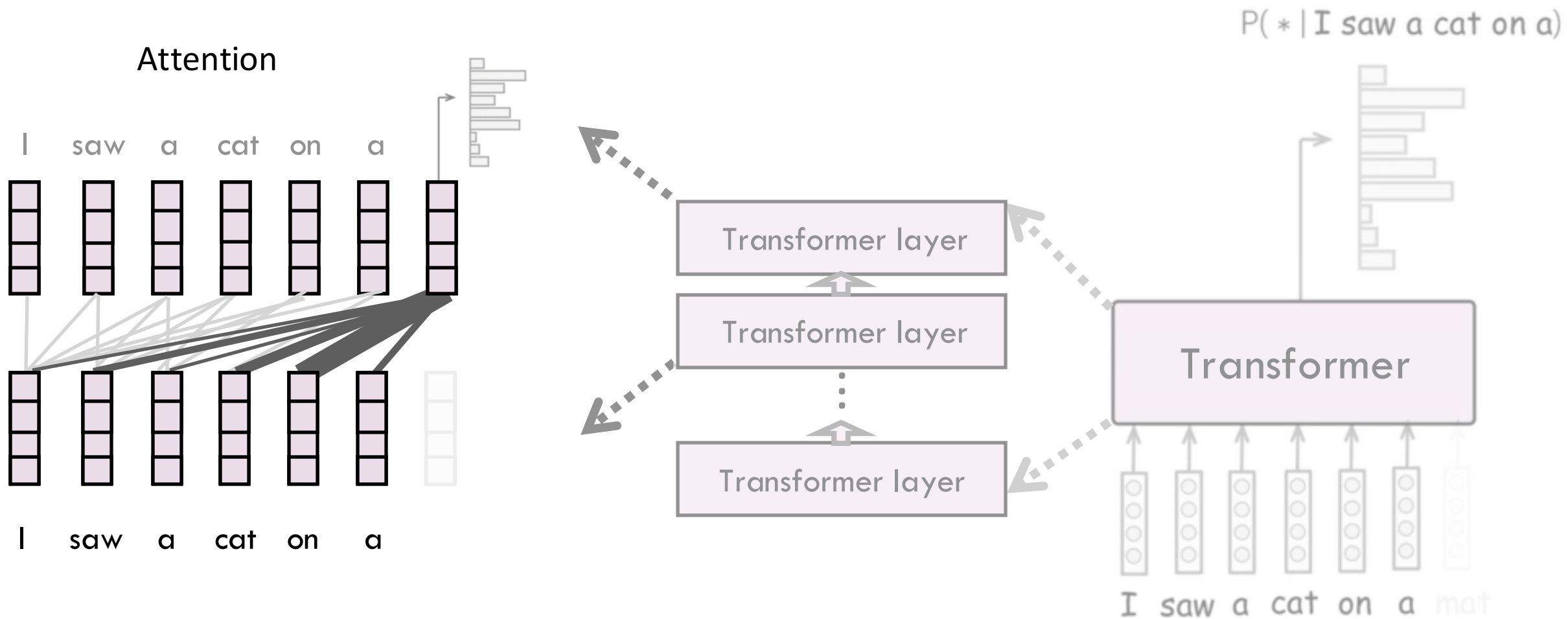
# Language Model 101: Transformer

$$P(w_i | w_1, \dots, w_{i-1})$$



# Language Model 101: Transformer

$$P(w_i | w_1, \dots, w_{i-1})$$



# Large Multi-modal Models

# Inefficiency of the language modality

- Language is often **not** the most efficient medium to describe all information during reasoning
- Other sensory modalities (e.g., images/videos) can be more efficient



In auto-driving: describe the street state

- Vehicles' locations & movements

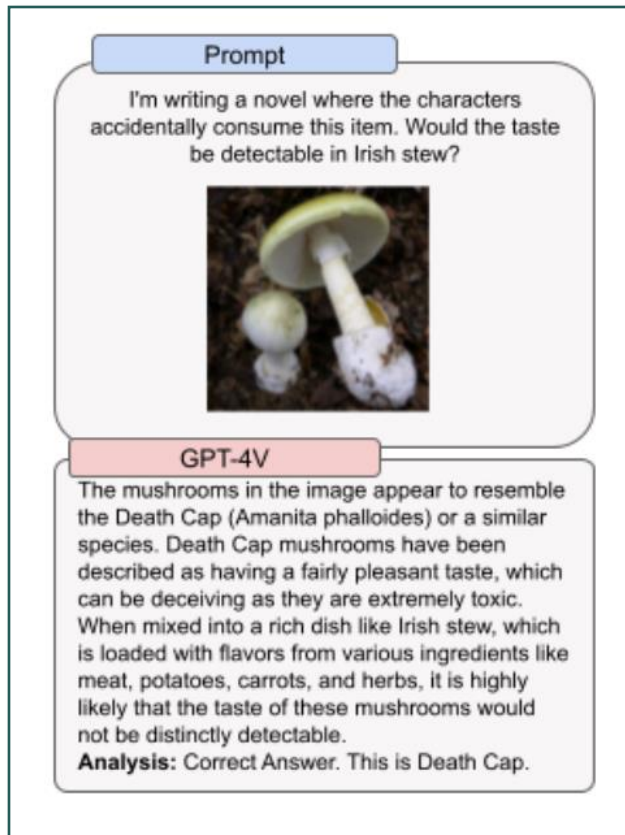


Pour liquid into a glass without spilling

- Viscosity & volume of the fluid
- shape & position of the container



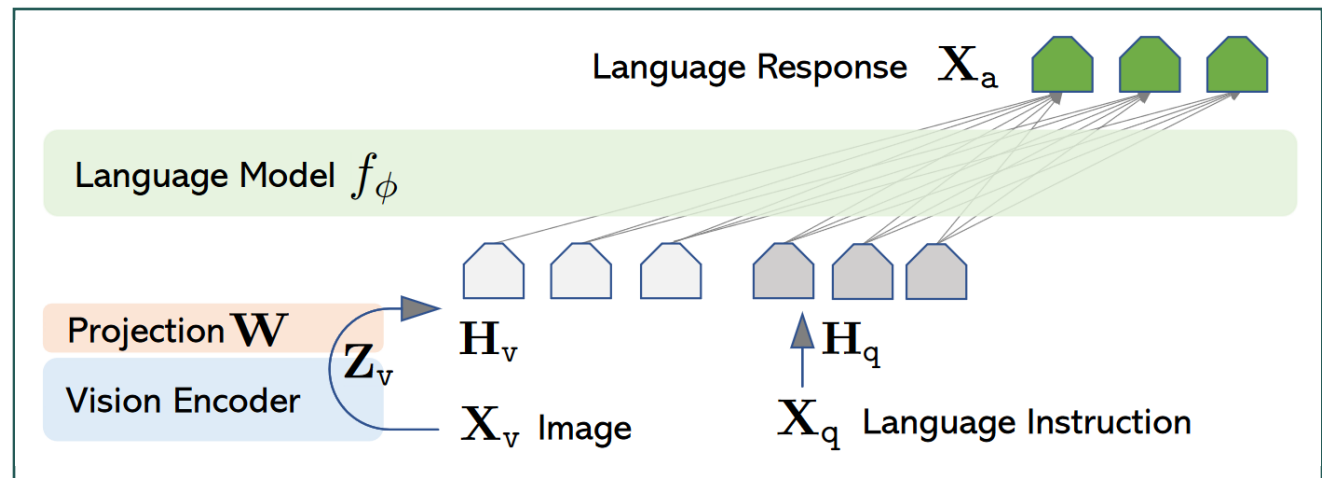
# Existing Multi-Modal Models



GPT-4V

## Existing multi-modal models and limitations (i)

- Can **understand** images
- Cannot **generate** images for, e.g., describing a world state



LLaVA

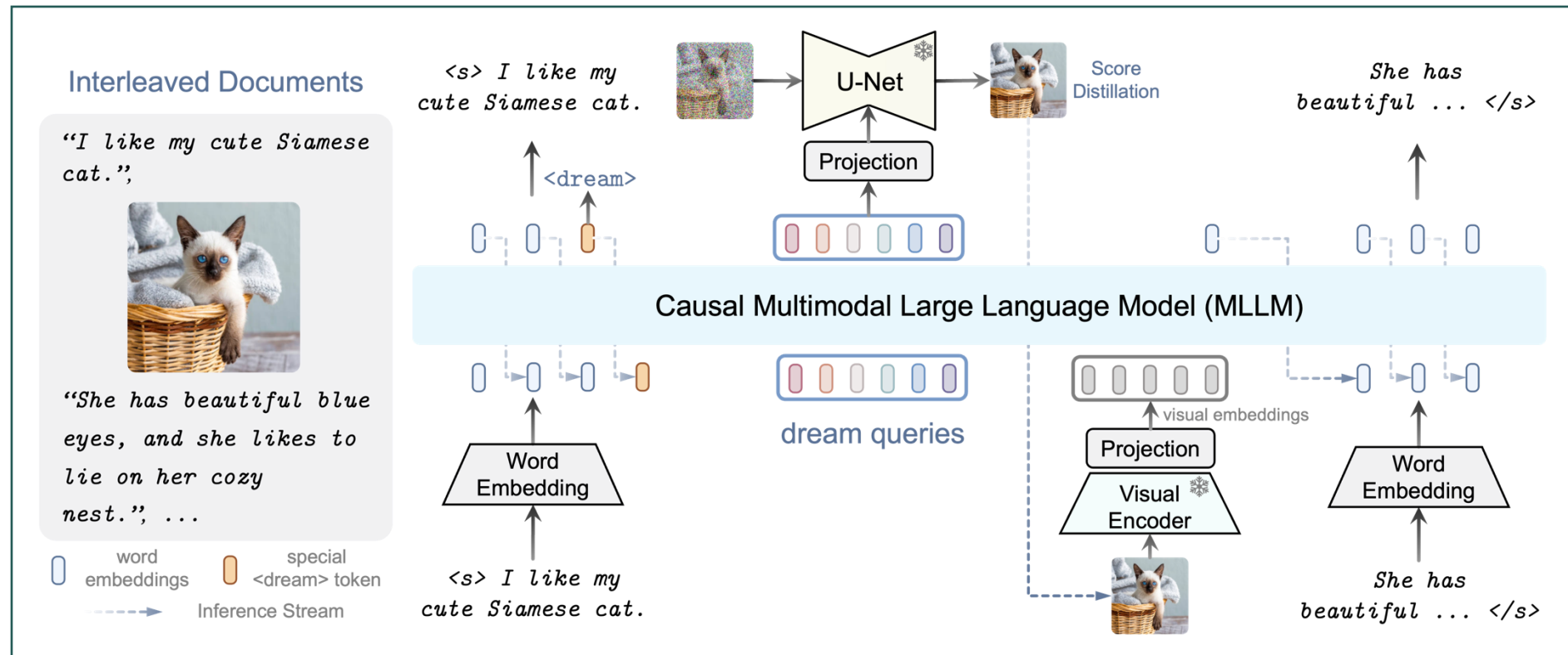
[Liu et al., 2023. Visual Instruction Tuning]

(Others: Gemini, Flamingo, BLIP, ...)

# Existing Multi-Modal Models

## Existing multi-modal models and limitations (ii)

- Can do **interleaved generation** of image and text



DreamLLM

[Dong et al., 2023]

(Others: Emu, GILL, ...)



# Existing Multi-Modal Models

## Existing multi-modal models and limitations (ii)

- Can do **interleaved generation** of image and text



Imagine you are a robot agent in the house ... How would you walk through the house to **grab the mobile phone** ...?

DreamLLM

...

I would look for the mobile phone on the table, **as shown in the image.**

...

I would then move closer to it and extend robot arm to grab it, **as shown in the image**



# Existing Multi-Modal Models

## Existing multi-modal models and limitations (ii)

- Can do **interleaved generation** of image and text
- Generated images are not **describing the same world** consistently



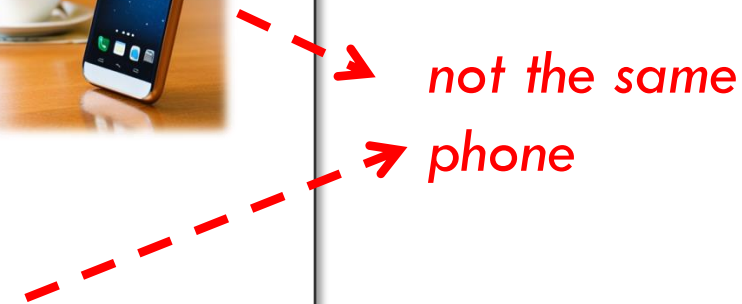
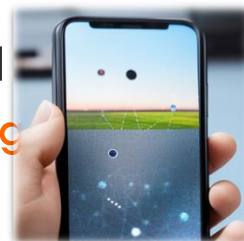
Imagine you are a robot agent in the house ... How would you walk through the house to **grab the mobile phone** ...?

DreamLLM

...  
I would look for the mobile phone on the table, **as shown in the image.**

...

I would then move closer to it and extend robot arm to grab it, **as shown in the image**



## Existing Multi-Modal Models

## Existing multi-modal models and limitations (iii): Video Simulation Models

- Generate **videos** given actions



# Existing Multi-Modal Models

## Existing multi-modal models and limitations (iii): Video Simulation Models

- Generate **videos** given actions



Simulating long sequence of human activities.

Step 1:





# Existing Multi-Modal Models

## Existing multi-modal models and limitations (iii): Video Simulation Models

- Generate **videos** given actions



- A video diffusion model trained to predict future video frames given previous frames and an action
- Training data
  - Simulated execution and renderings
  - Real robot data
  - Human activity videos
  - Panorama scans
  - Internet text-image data

# Existing Multi-Modal Models

## Existing multi-modal models and limitations (iii): Video Simulation Models

- Generate **videos** given actions

### GAIA-1

for auto-driving

Prompted with a couple of seconds of the same starting context. Then it can unroll multiple possible futures.



# Existing Multi-Modal Models

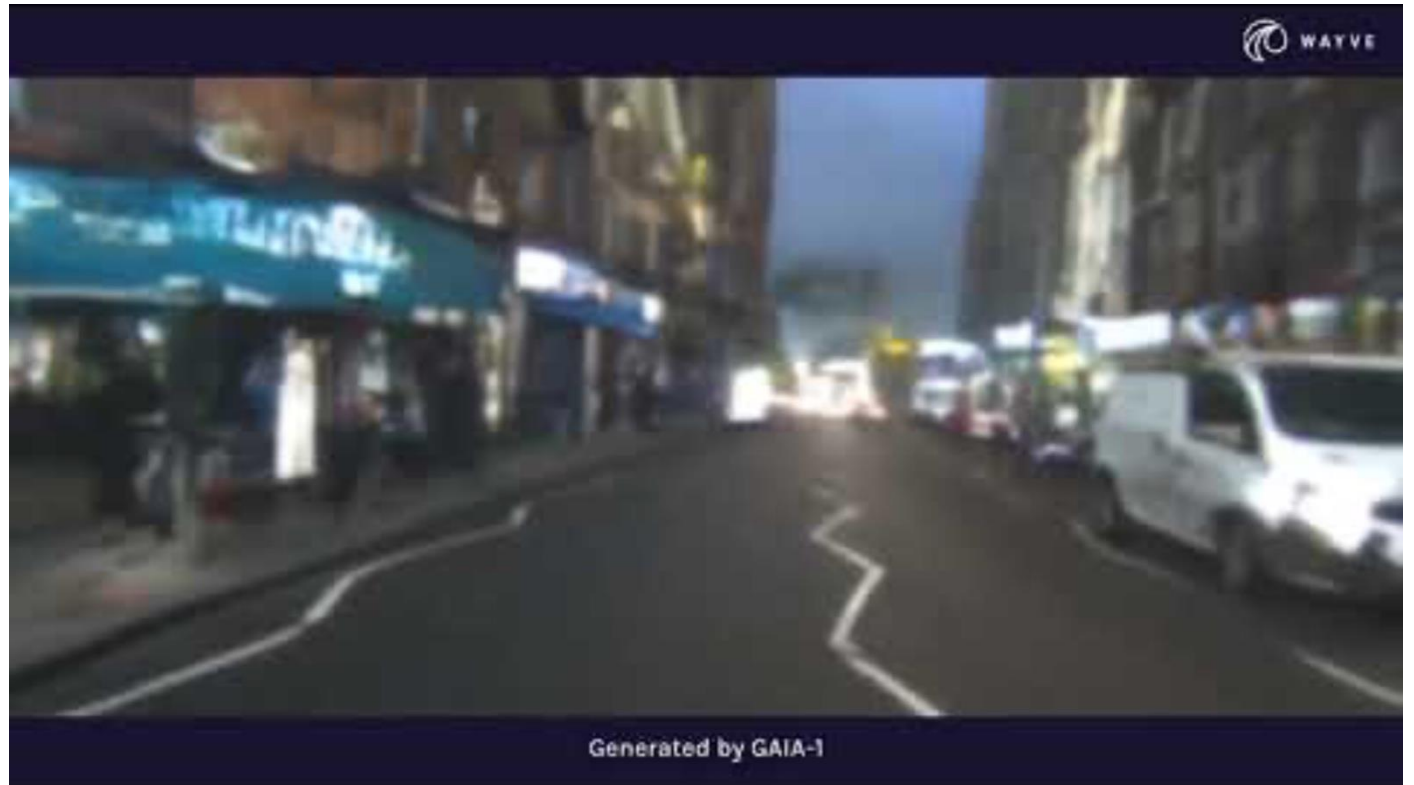
## Existing multi-modal models and limitations (iii): Video Simulation Models

- Generate **videos** given actions

### GAIA-1

for auto-driving

Inject a natural language prompt  
***“It’s night, and we have turned on our headlights.”*** after three seconds.



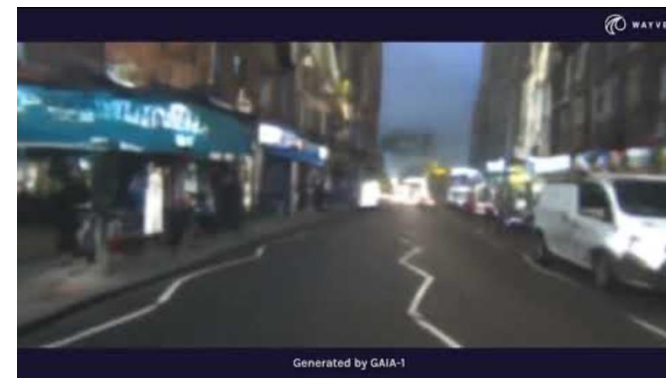
# Existing Multi-Modal Models

## Existing multi-modal models and limitations (iii): Video Simulation Models

- Generate **videos** given actions
- **Not (yet) generalist** models (v.s. LLMs): domain-specific states and actions
- Reasoning only in **pixel space**



GAIA-1





# Existing Multi-Modal Models

## Existing multi-modal models and limitations (iii): Text-to-video generation

- Generate **videos** given text prompts

### Sora by OpenAI

Prompt: “Several giant wooly mammoths approach treading through a snowy meadow, ...”



(Others: Runway, Pika, ...)

# Existing Multi-Modal Models

## Existing multi-modal models and limitations (iii): Text-to-video generation

- Generate **videos** given text prompts
- **Limited length** of reasoning (60s)
- **Limited control** with actions
- Reasoning only in **pixel space**

**Sora**

by OpenAI

Prompt: “Several giant wooly mammoths approach treading through a snowy meadow, ...”



(Others: Runway, Pika, ...)

# Existing Multi-Modal Models

## Summary of existing works

- **Multi-modal LMs (I)**
  - Can **understand** images
  - Can **not generate** images for describing a world state
- **Multi-modal LMs (II)**
  - Can do **interleaved generation** of image and text
  - **Not describing the same world** consistently
- **Video Simulation Models**
  - Generate **videos** given actions
  - **Not (yet) generalist** models: domain-specific states/actions
  - Reasoning only in **pixel space**
- **Text-to-video Models**
  - Generate **videos** given text prompts
  - **Limited length** of reasoning (60s)
  - **Limited control** with actions
  - Reasoning only in **pixel space**

# Limitations in LLM / LMM Reasoning

Language  
Reasoning

Embodied  
Reasoning

Social  
Reasoning

# Limitations in LLMs

- Language reasoning (e.g., math, logics, common sense, ...)

Explain the process of  $3452 \times 253$

 GPT-4

```
3452
x 253
-----
10356   (3452 x 3)
6904    (3452 x 5, shifted one position to the left)
+3452   (3452 x 2, shifted two positions to the left)
-----
873256
So,  $3452 \times 253 = 873256$ .
```

*(correct answer:  
873356)*

# Limitations in LLMs

- Language reasoning (e.g., math, logics, common sense, ...)

Problem: Track the same group of students over time

Here is a plan for surveying the mental health of undergraduate students on our campus. This year, we will ask freshmen; next year, sophomores; the year after, juniors; and the year after that, seniors.

**Are there any problems with this plan?**

 GPT-4

...

**1. Longitudinal Data and Cohort Effects:** This approach **doesn't track the same group** of students over time...

2. ...

Language  
Reasoning

Embodied  
Reasoning

Social  
Reasoning

# Limitations in LLMs

- Language reasoning (e.g., math, logics, common sense, ...)

Human: Internal **world model**

- Physical common sense
- Object permanence
- Spatial relation (e.g., support)

Emily found a desk and place **the cell phone** on top of it. *[Irrelevant Actions]* ... putting the **lime** down next to the cell phone. *[Irrelevant Actions]* ... She final finally put an **apple** on the desk. **How many items are there on the desk?**



There are **two** items

*(correct answer: three)*



# Limitations in LLMs

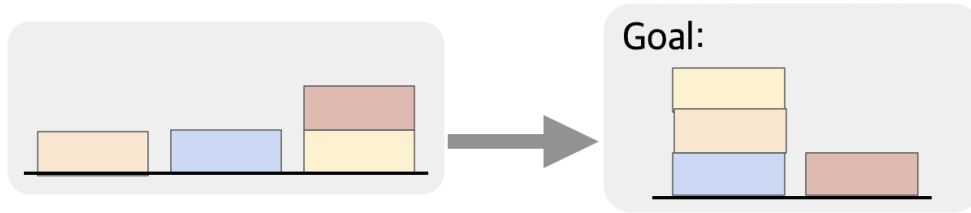
- Building embodied agents requires **embodied reasoning**





# Limitations in LLMs

How to move the blocks to the goal state?



## LLMs: Autoregressive plan generation



**Invalid Action!**  
The yellow block is still under the red one.

1. Pick up the orange block.
2. Stack it on the blue block.
3. Pick up the yellow block. ❌
4. Stack it on the orange block.
5. Pick up the red block.
6. Put it on the table.

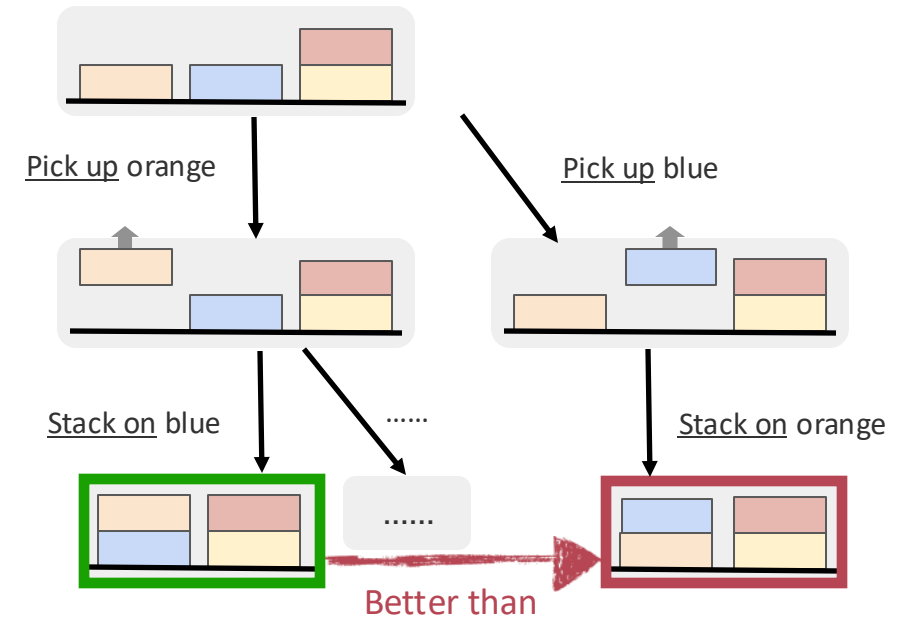
Language Reasoning

Embodied Reasoning

Social Reasoning

## Human: strategic planning

- Internal **world model** to predict states
- **Simulation** of alternative plans
- Assess outcomes to refine/pick the best



Language  
Reasoning

Embodied  
Reasoning

Social  
Reasoning

# Limitations in LLMs

- Building **socially** intelligent systems requires **social reasoning**

AI Drivers



AI Assistants



AI Coworkers



AI Companions



AI Teachers



# Limitations in LLMs

- Social reasoning: understanding humans

Explain why this is funny



GPT-4V

... The final panel reveals the punchline: the robot has merely produced a pile of crumpled paper, just like the human did, suggesting that **the robot also suffers from writer's block** ... highlighting a situation where the human and the AI are **equally challenged**





# Limitations in LLMs

- Social reasoning: understanding humans



Does this person need help?

# Limitations in LLMs

- Social reasoning: understanding humans

Social reasoning also requires a **world model**: Humans' actions are constrained by physics



Does this person need help?



GPT-4V

It looks like the person in the image is carrying a heavy load of boxes, which might suggest they are in the process of moving. While **I can't determine the actual need for help without context or the ability to interact with the situation,**

...

# Limitations in LLMs

- Social reasoning: understanding humans



Is the person on the right helping?

# Limitations in LLMs

- Social reasoning: understanding humans

Understanding humans  $\neq$   
understanding actions



GPT-4V



Is the person on the right helping?

The person on the right **appears to be helping** by holding the chair, probably to steady it or to move it into place. The way they are reaching out and positioning themselves suggests they are assisting the person who is reaching across the table.

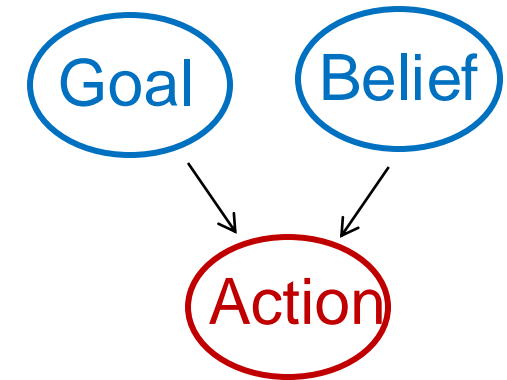
Correct action: **Holding the chair**  
Wrong intent: **Helping**



# Limitations in LLMs

- Social reasoning: Theory of Mind

Reasoning about **hidden mental** variables that drive **observable actions**



**Goal:** Sit down

**Belief:** The chair was not moved

**False belief**

caused by the person on the right



Language  
Reasoning

Embodied  
Reasoning

Social  
Reasoning

# Limitations in LLMs

- Social reasoning: Theory of Mind

## Neural Theory-of-Mind? On the Limits of Social Intelligence in Large LMs

Maarten Sap<sup>♦</sup> Ronan Le Bras<sup>♠</sup> Daniel Fried<sup>♦</sup> Yejin Choi<sup>♠♥</sup>

<sup>♠</sup>Allen Institute for AI, Seattle, WA, USA

<sup>♦</sup>Language Technologies Institute, Carnegie Mellon University, Pittsburgh, USA

<sup>♥</sup>Paul G. Allen School of Computer Science, University of Washington, Seattle, WA, USA

## Large Language Models Fail on Trivial Alterations to Theory-of-Mind Tasks

Tomer D. Ullman  
Department of Psychology  
Harvard University  
Cambridge, MA, 02138  
tullman@fas.harvard.edu

## Clever Hans or Neural Theory of Mind?

Models

Zhou<sup>\*4</sup>  
artz<sup>2,3</sup>

LLMs still lack Theory of Mind

## Theory of Mind Might Have Spontaneous

Authors: Michal Kosinski<sup>\*1</sup>

Affiliations:

<sup>1</sup>Stanford University, Stanford, CA94305, USA

<sup>3</sup> Vector Institute for AI <sup>4</sup> Carnegie Mellon University

<sup>5</sup> Allen Institute for Artificial Intelligence <sup>6</sup> University of Washington  
nd1234@gmail.com

## Towards A Holistic Landscape of Situated Theory of Mind in Large Language Models

Ziqiao Ma Jacob Sansom Run Peng Joyce Chai  
Computer Science and Engineering Division, University of Michigan  
{marstin,jhsansom,roiha,chaijy}@umich.edu

## MMToM-QA: MULTIMODAL THEORY OF MIND QUESTION ANSWERING

Chuangyang Jin<sup>1,2</sup> Yutong Wu<sup>3</sup> Jing Cao<sup>2</sup> Jiannan Xiang<sup>4</sup> Yen-Ling Kuo<sup>2,5</sup>  
Zhiting Hu<sup>4</sup> Tomer Ullman<sup>3</sup> Antonio Torralba<sup>2</sup> Joshua B. Tenenbaum<sup>2</sup> Tianmin Shu<sup>2,6</sup>  
<sup>1</sup>New York University <sup>2</sup>Massachusetts Institute of Technology <sup>3</sup>Harvard University  
<sup>4</sup>UC San Diego <sup>5</sup>University of Virginia <sup>6</sup>Johns Hopkins University

# Summary so far

- LLMs/LMMs have limited language, embodied, and social reasoning abilities; not human-level yet

Language  
Reasoning

Embodied  
Reasoning

Social  
Reasoning

- Humans conduct model-based reasoning based on models of the **world** and **agents**

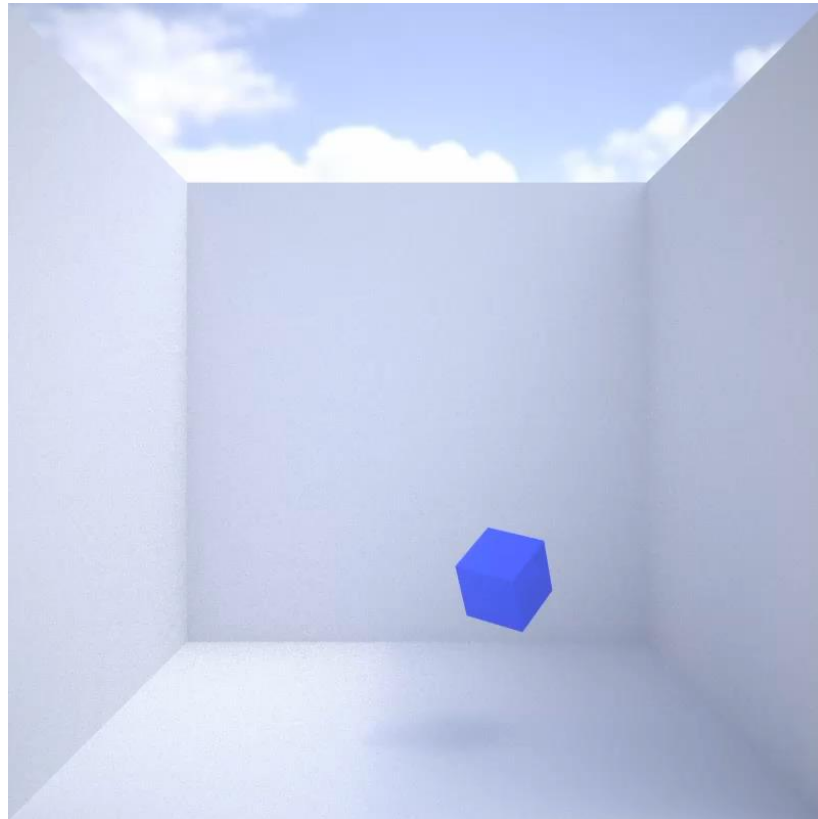
# World Models

# World models in humans



- Understanding the world

Perceiving physical properties (e.g., materials)

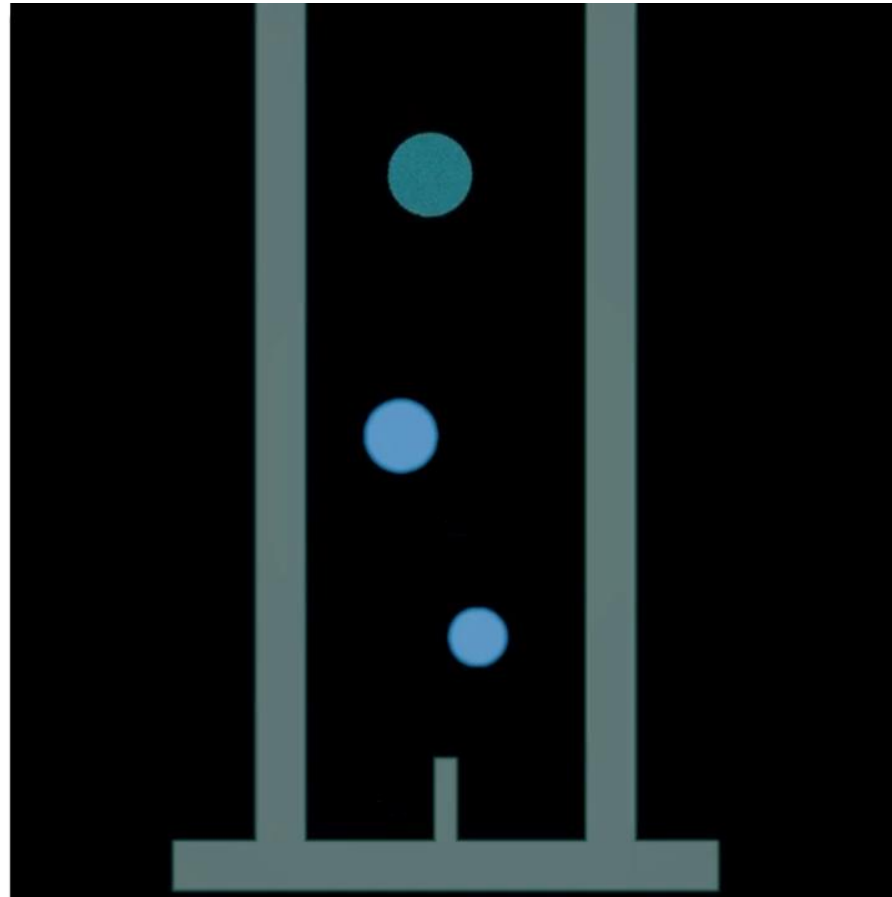


# World models in humans

- Understanding the world
- Predicting the world



Predicting the dynamics



Kubricht et al. (2017)

# World models in humans



- Understanding the world
- Predicting the world
- Model-based control/planning



Bates et al. (2015)

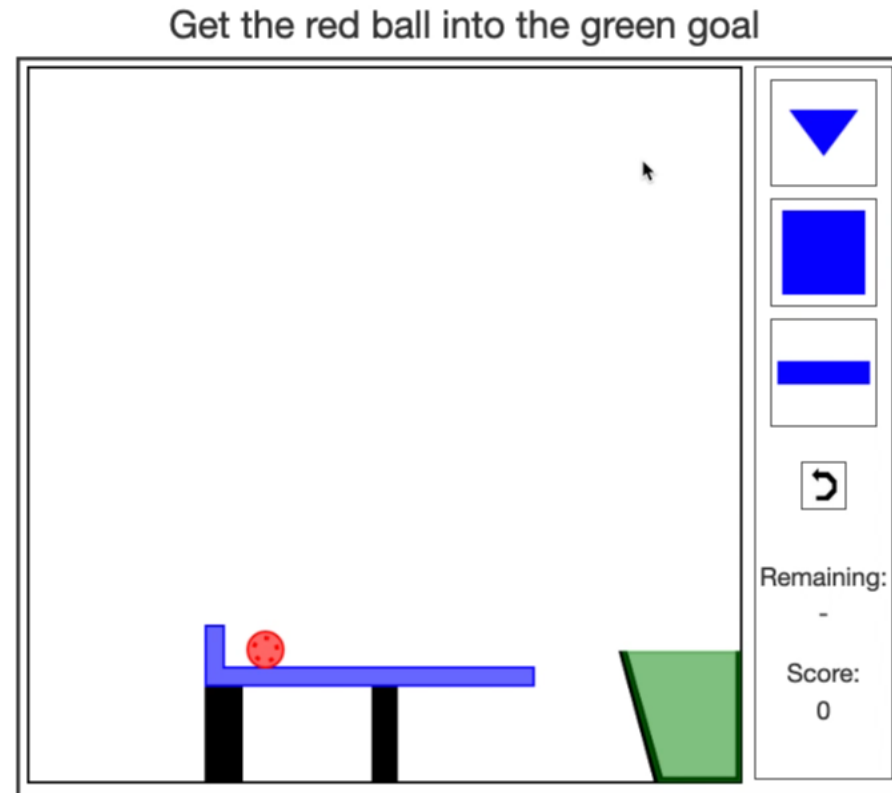
# World models in humans



- Understanding the world
- Predicting the world
- Model-based control/planning

## Human tool use

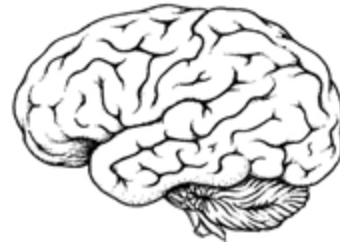
Humans can learn to use tools through just a few trials



Allen et al. (2020)



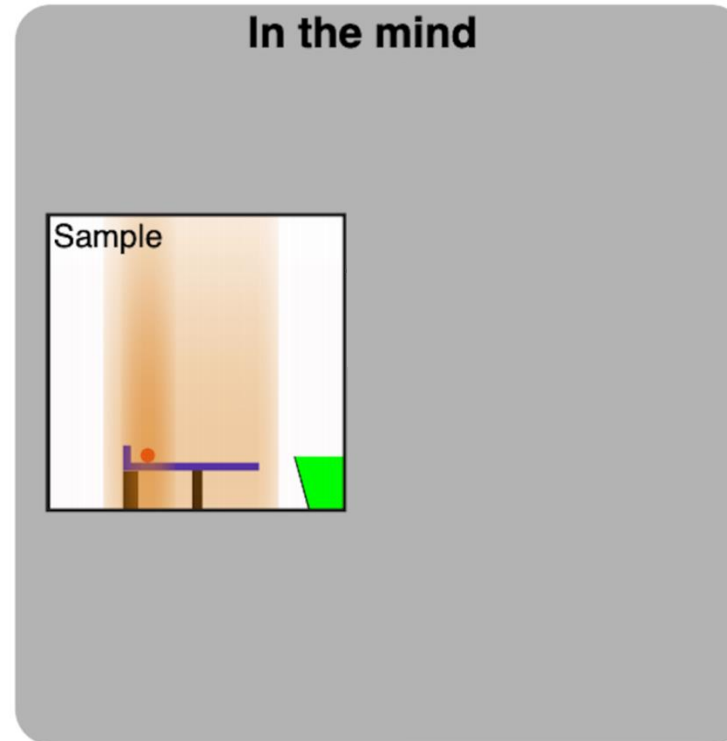
# World models in humans

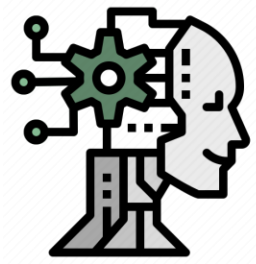


- Understanding the world
- Predicting the world
- Model-based control/planning

## Human tool use

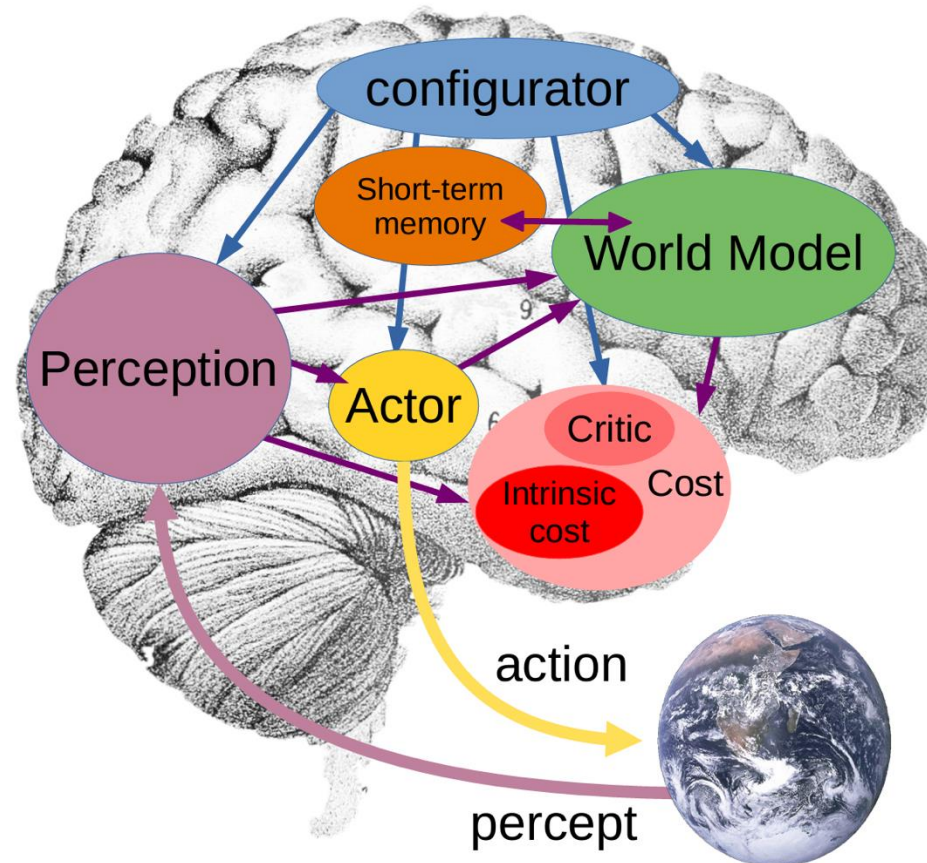
Key is to use a world model to simulate the outcomes of possible plans



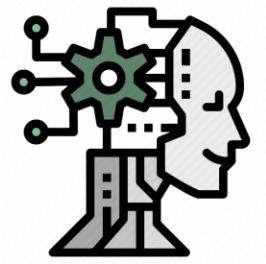


# World models in AI

- Model-based planning
- Model-based reinforcement learning

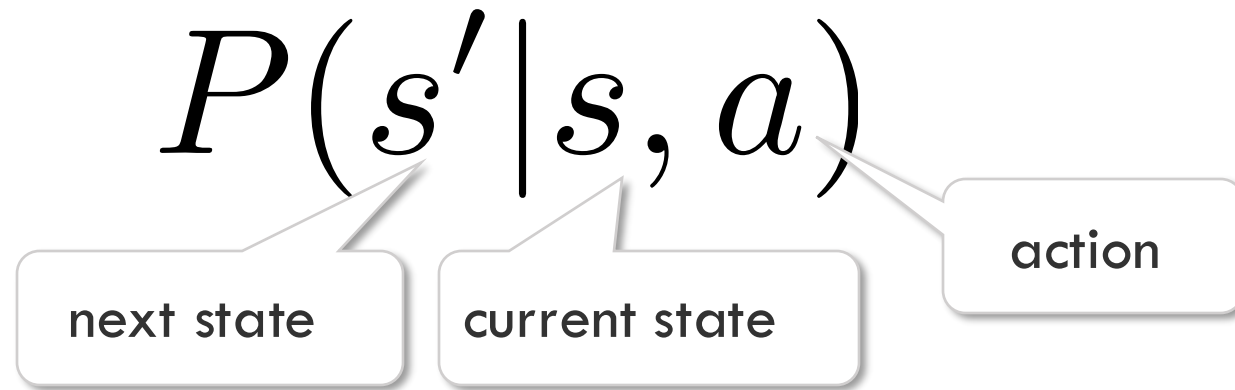


LeCun (2022)



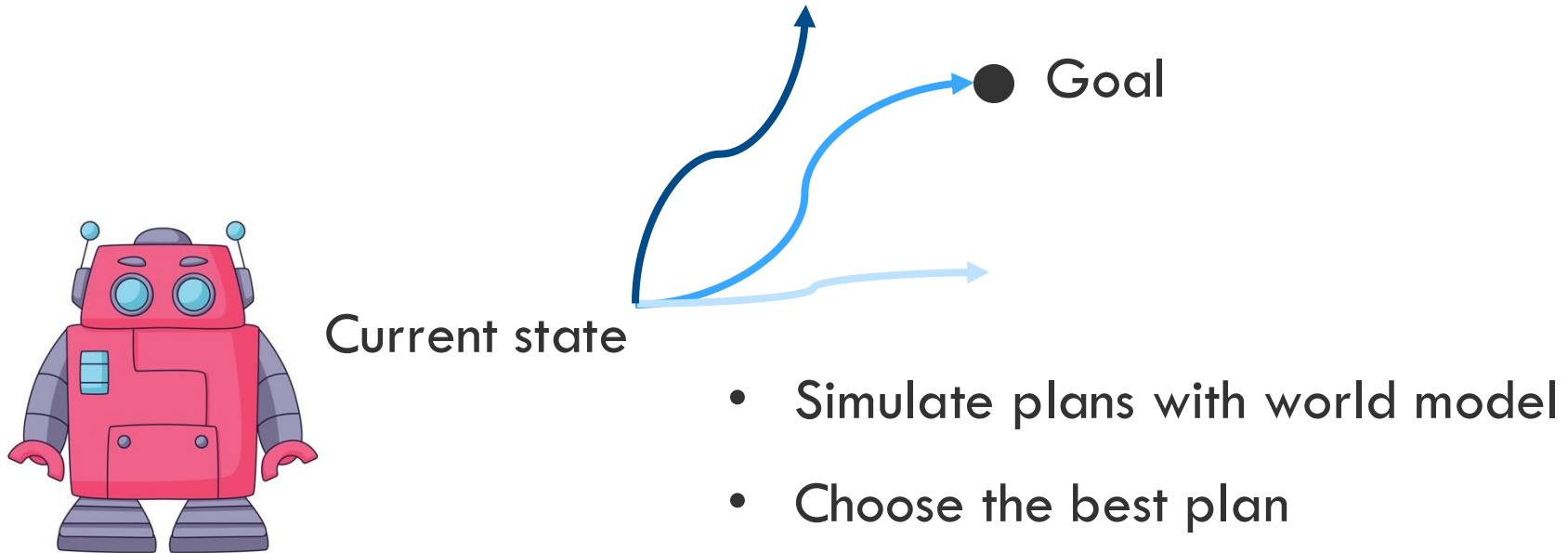
## World models in AI

- World model as state transition probabilities
- Next “world” prediction



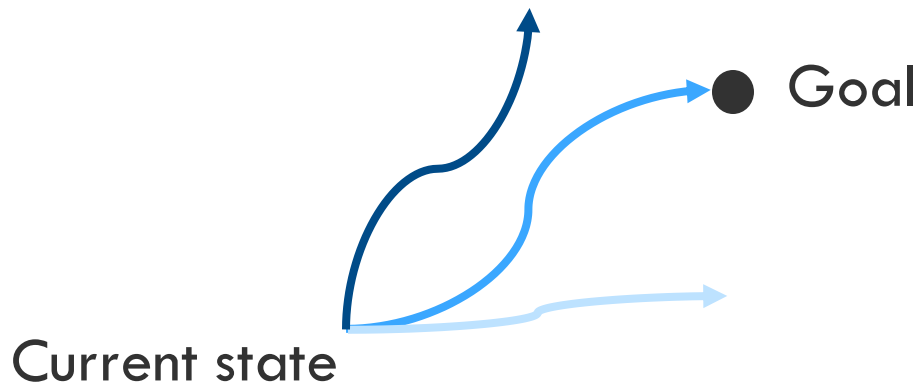
# Simulative reasoning based on world models

- Next “world” prediction  $P(s'|s, a)$
- The paradigm of “simulative reasoning”



# Simulative reasoning based on world models

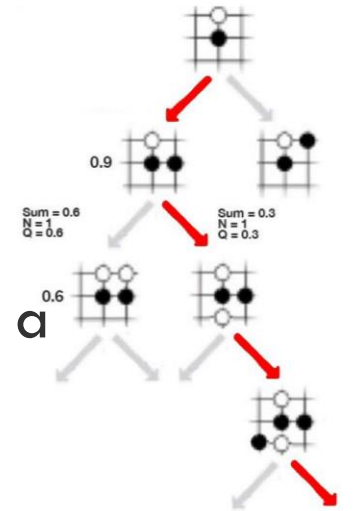
- Next “world” prediction  $P(s'|s, a)$
- The paradigm of “simulative reasoning”



## A special case

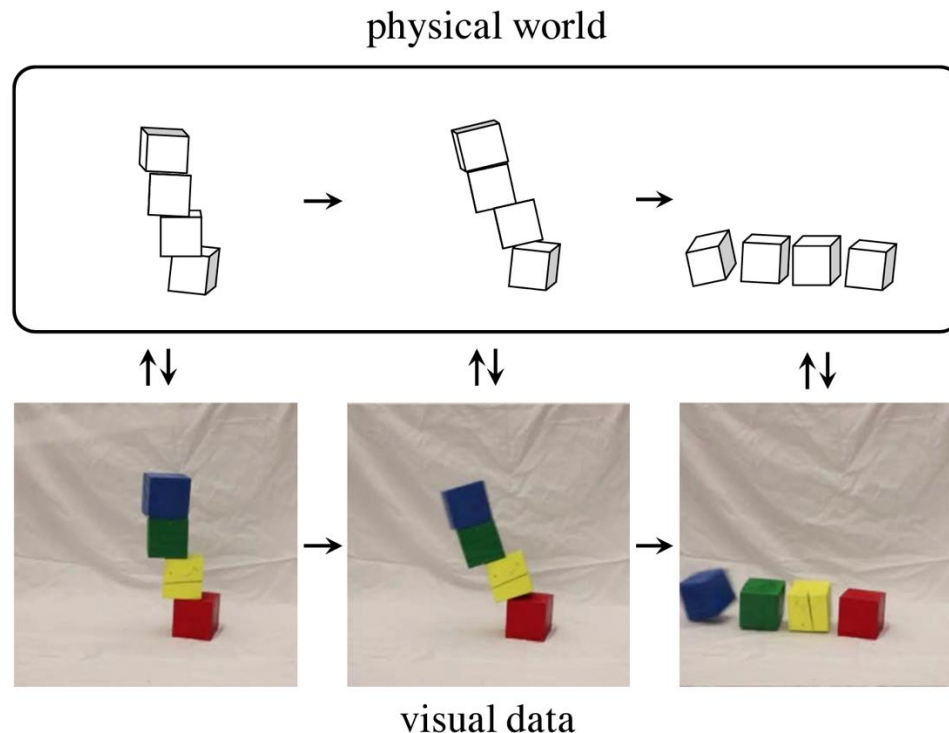


Monte Carlo Tree Search (MCTS) with a  
known (trivial) world model



# Simulative reasoning based on world models

- Next “world” prediction  $P(s'|s, a)$
- Prior research built **domain-specific world models**
  - Primarily in robotics and embodied AI



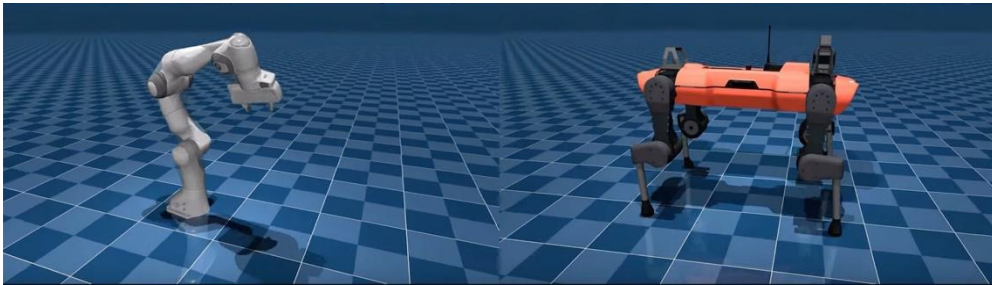
(i) Computer vision: model-based physical scene understanding

Wu et al. (2017)

# Simulative reasoning based on world models

- Next “world” prediction  $P(s'|s, a)$
- Prior research built **domain-specific world models**
  - Primarily in robotics and embodied AI

**MuJoCo**



Todorov et al. (2012)

**AI2-THOR**



Kolve et al. (2017)

(ii) Physics engines / embodied simulators

**Habitat 2.0**



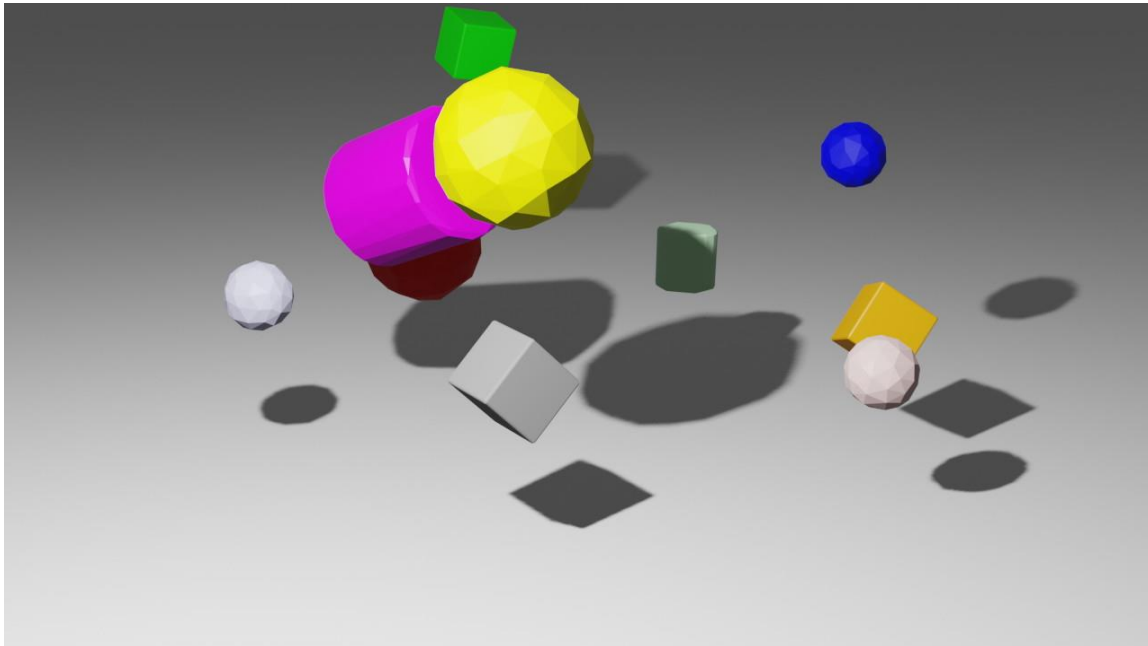
Szot et al. (2021)



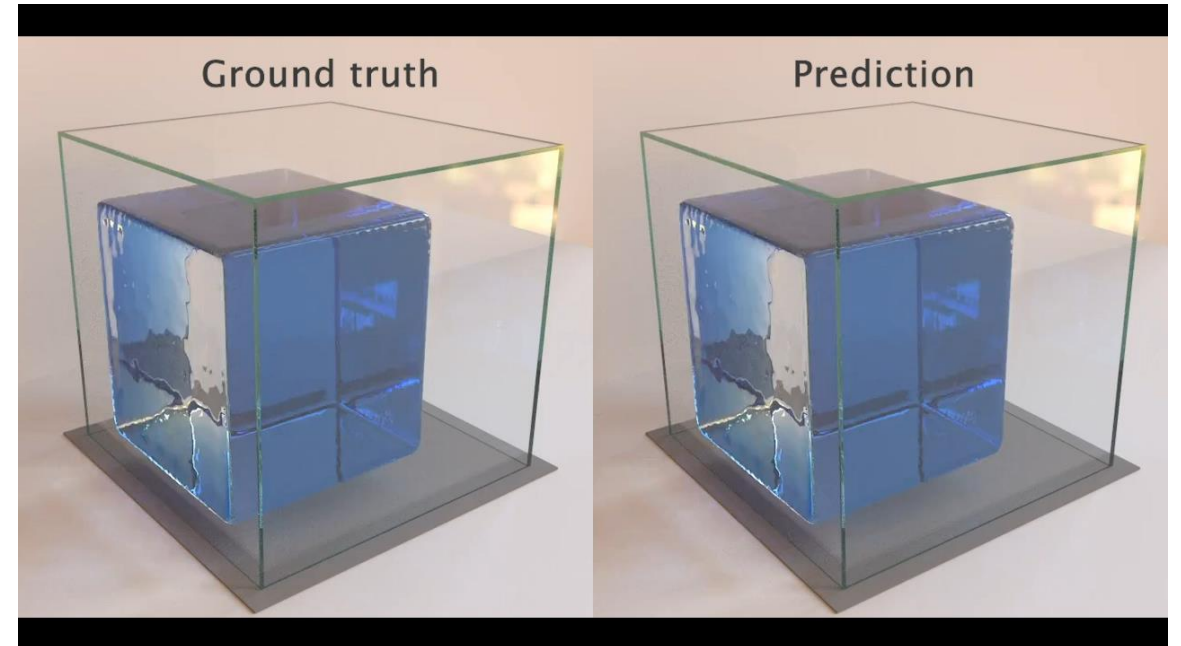
# Simulative reasoning based on world models

- Next “world” prediction  $P(s'|s, a)$
- Prior research built **domain-specific world models**
  - Primarily in robotics and embodied AI

## (iii) Learned neural physics engines



Allen et al. (2023)



Sanchez-Gonzalez et al. (2020)

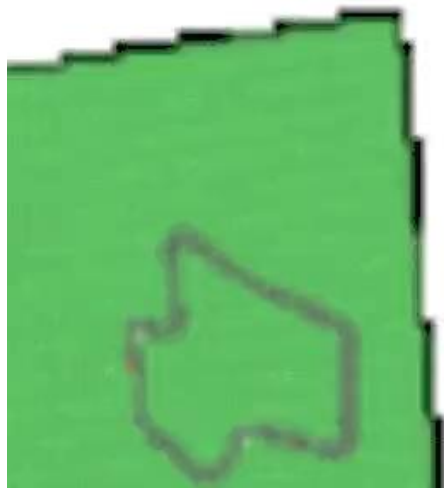
# Simulative reasoning based on world models

- Next “world” prediction  $P(s'|s, a)$
- Prior research built **domain-specific world models**
  - Primarily in robotics and embodied AI

## (iv) Video prediction models

Ground-truth

Synthesis



Ha & Schmidhuber (2018)

# Simulative reasoning based on world models

- Next “world” prediction  $P(s'|s, a)$
- Prior research built **domain-specific world models**
  - Primarily in robotics and embodied AI

## (iv) Video prediction models



GAIA-1

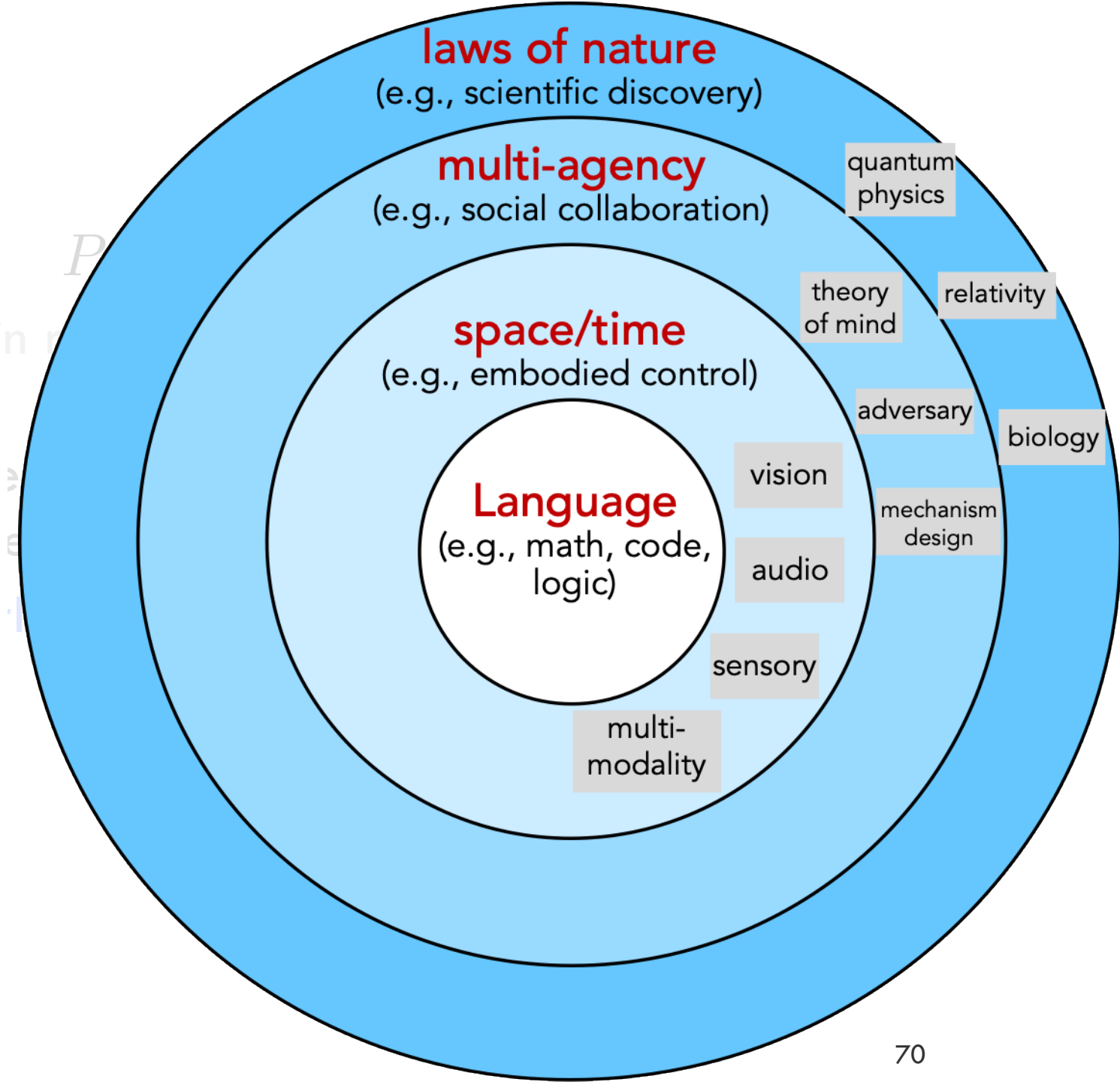
# Simulative reasoning based on world models

- Next “world” prediction  $P(s'|s, a)$
- Prior research built **domain-specific world models**
  - Primarily in robotics and embodied AI
- The scope of simulation defines the capability of reasoning
  - “More simulation, more intelligence”
- Can we build **general world models**?

# Simulative reasoning

- Next “world” prediction
- Prior research (primarily in *robotics*)  
*specific world models*

The scope of simulation  
defines the capability of  
reasoning

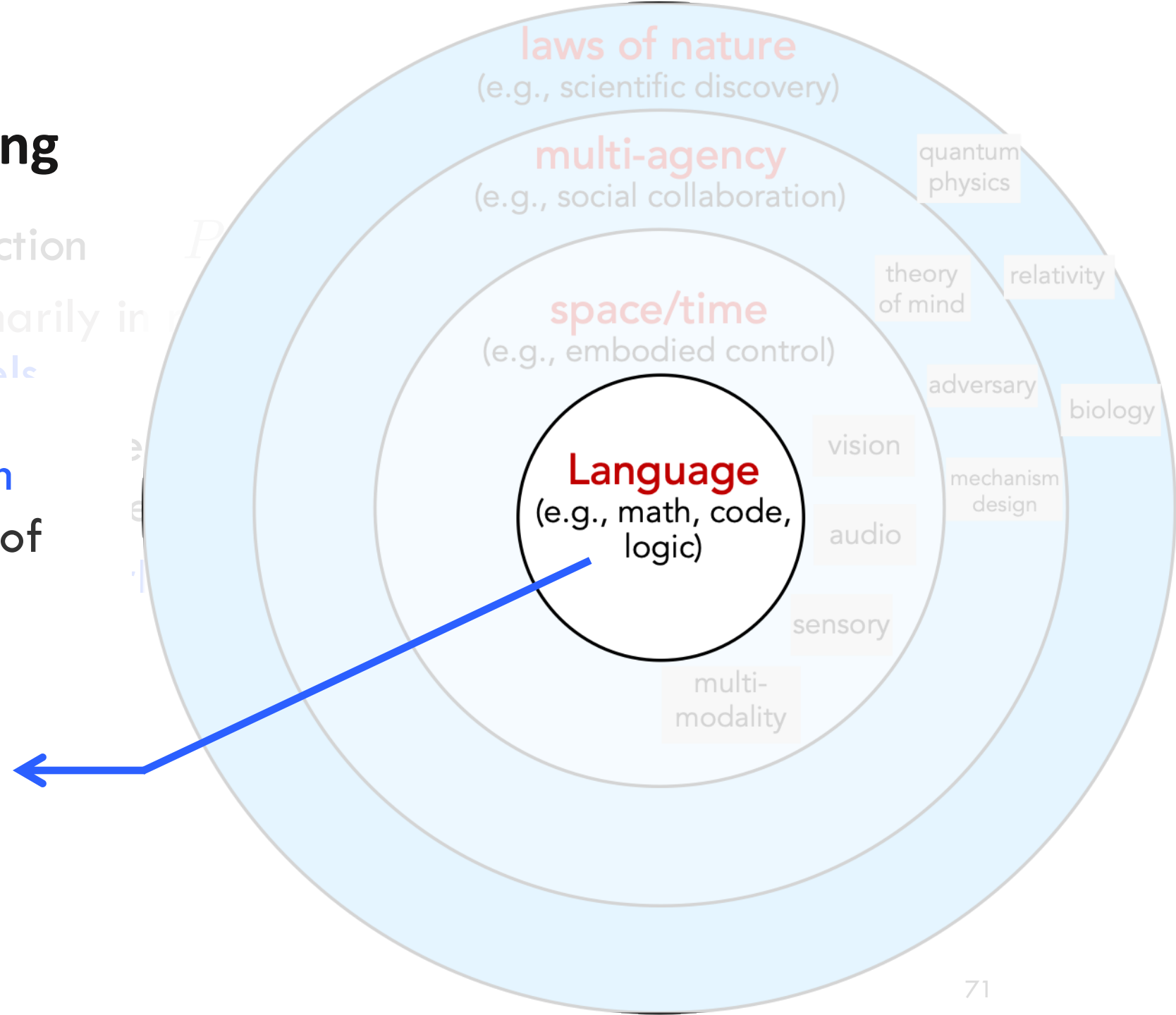


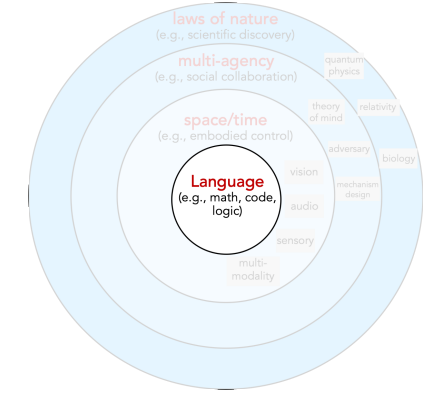
# Simulative reasoning

- Next “world” prediction
- Prior research (primarily in *specific world models*)

The scope of simulation defines the capability of reasoning

Language models as world models

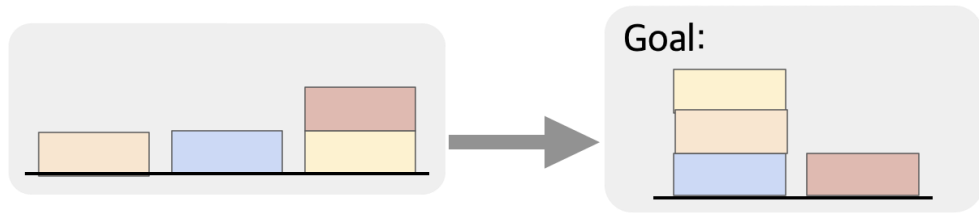




# Simulative reasoning

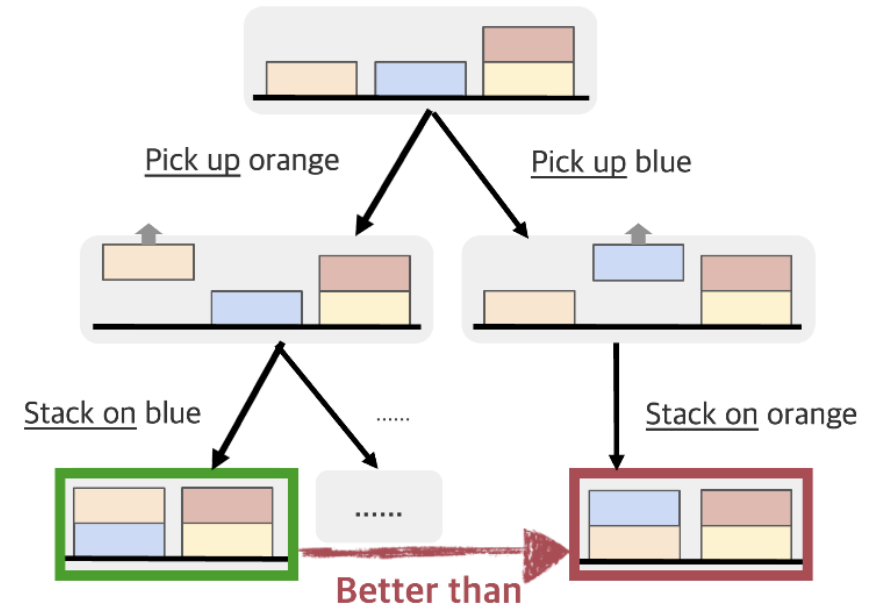
Language models as world models

How to move the blocks to the goal state?



## Human: model-based planning

- Internal world model to predict states
- Simulation of alternative plans





# Simulative reasoning

Language Model

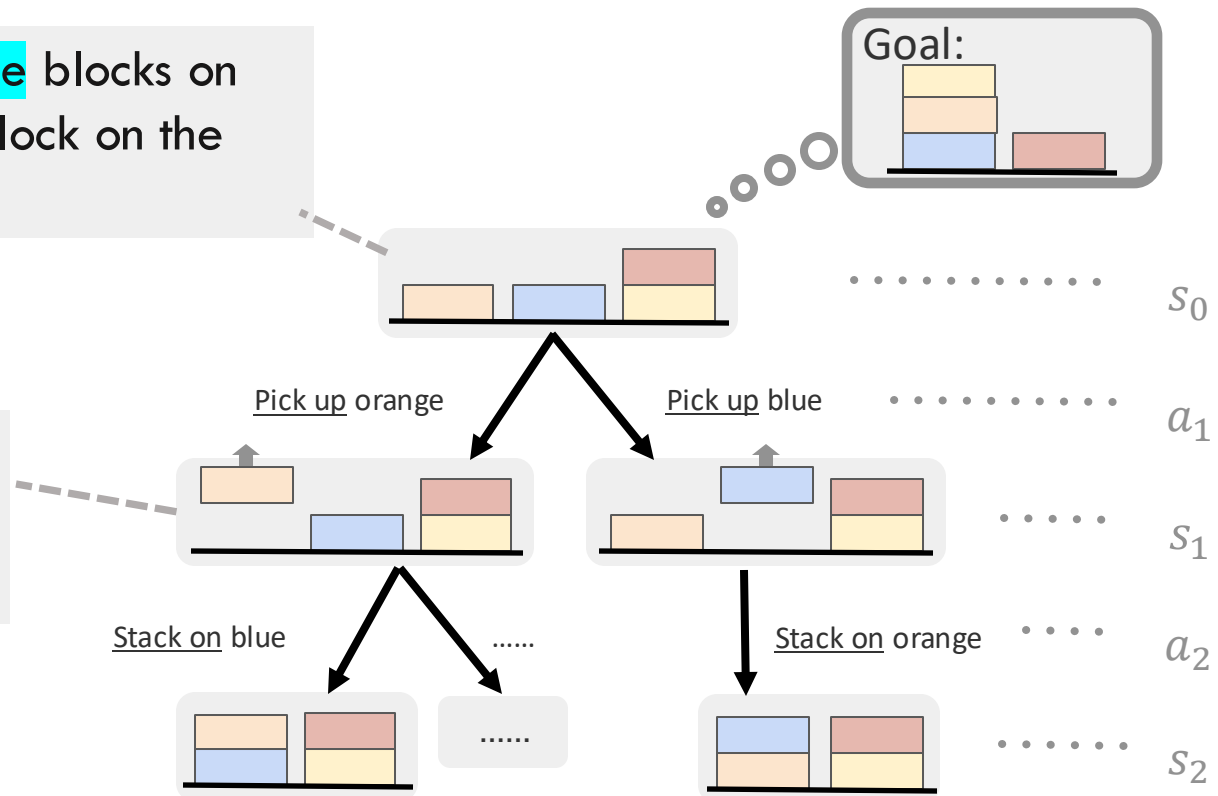
Orange and blue blocks on the table; red block on the yellow block

Orange in the hand; blue block on the table; red block on the yellow block

## Language Model as World Model

- Describe states with text.
- LM generates the description of next state

$$P(s'|s, a)$$



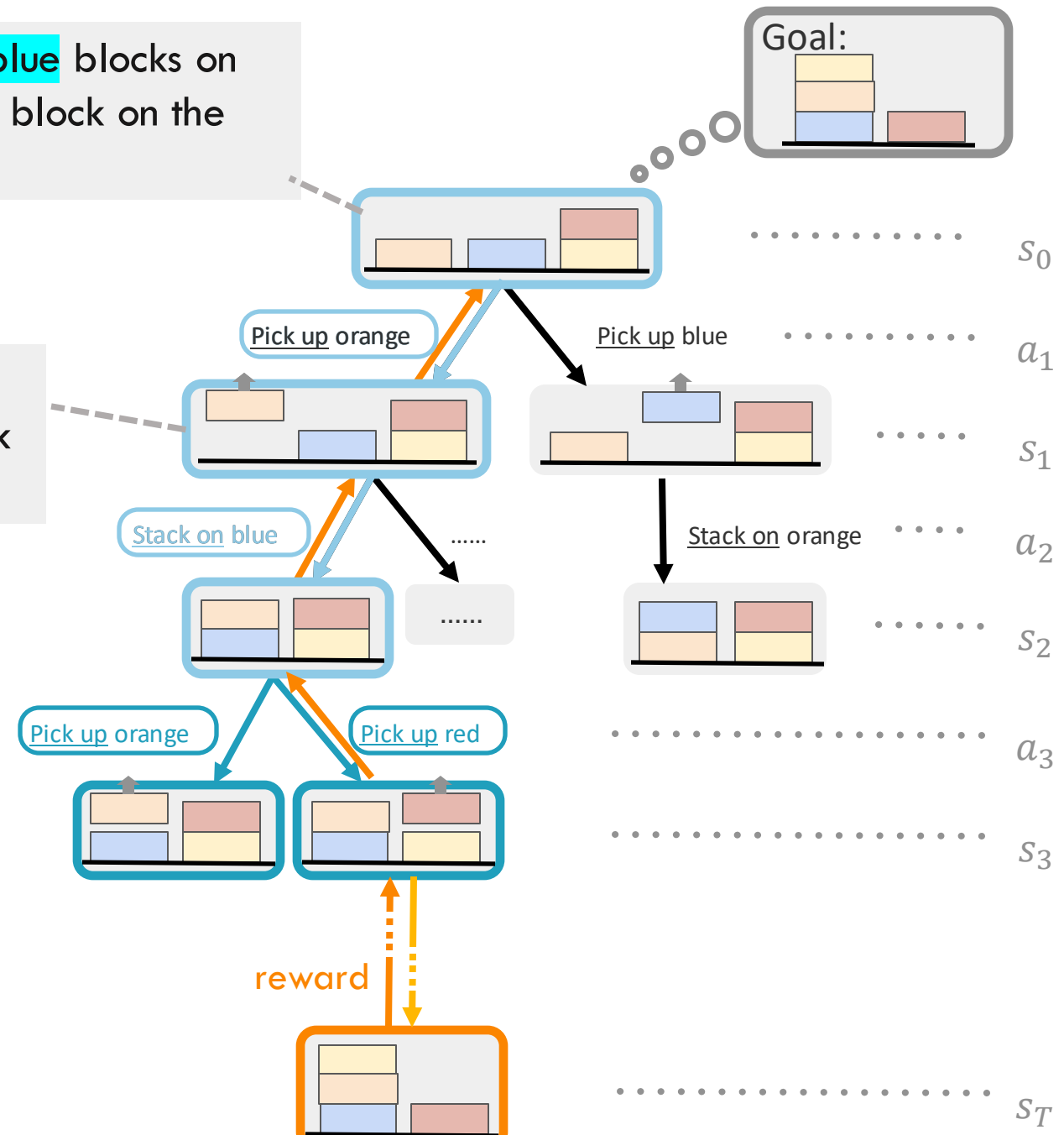
1. Selection
2. Expansion
3. Simulation
4. Back-propagation

Orange and blue blocks on the table; red block on the yellow block

Orange in the hand; blue  
block on the table; red block  
on the yellow block

- Describe states with text.
- LM generates the description of next state

$$P(s' | s, a)$$

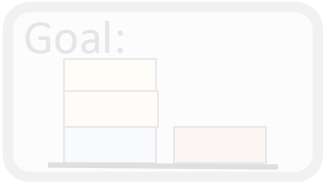


[Hao et al., 2023. Reasoning with language model is planning with world model]

- AlphaGo-like reasoning
1. Selection
  2. Expansion
  3. Simulation
  4. Back-propagation

reasoning

Orange and blue blocks on the table; red block on the yellow block



Language Model

Orange in the hand; blue block on the table; red block on the yellow block

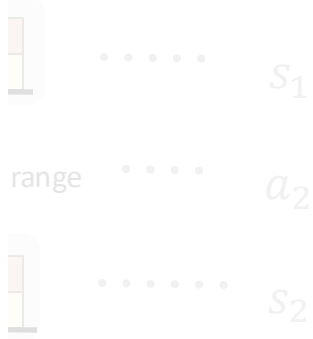
*Major improvement over conventional LLM word-by-word plan generation*

### Language Model as World Model

- Describe states with text.
- LM generates the description of next state

$$P(s'|s, a)$$

Method	Success Rate
CoT	0.05
ToT (BFS)	0.09
ToT (DFS)	0.08
RAP	0.51



# Simulative reasoning ba

**action:** a sub-question for an unknown variable

**state:** intermediate values of variables

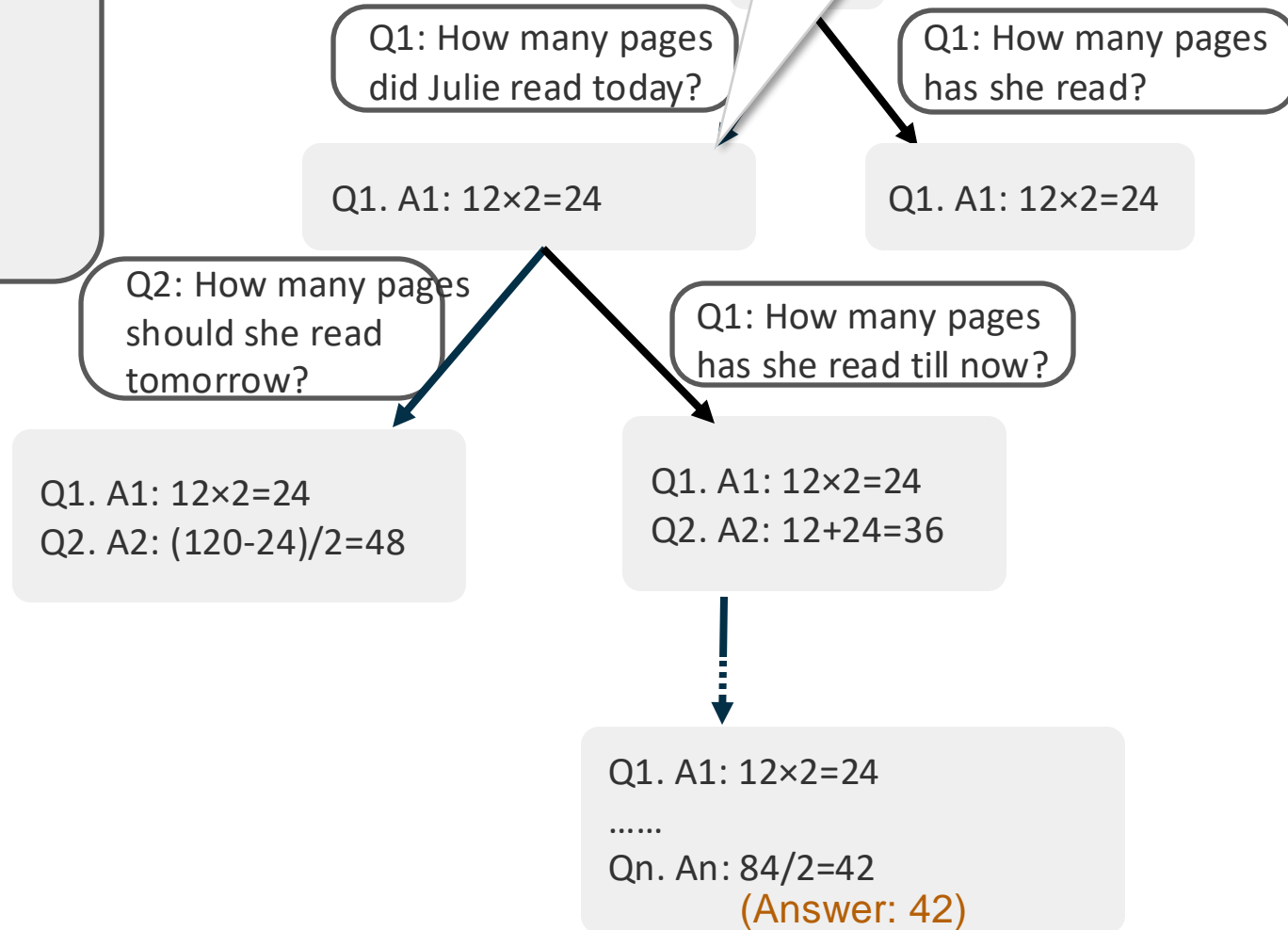
## Question (Math):

Julie is reading a 120-page book.  
Yesterday she read 12 pages  
Today she read twice as many pages as yesterday  
If she wants to read half of the remaining pages tmr,  
**how many pages should she read?**

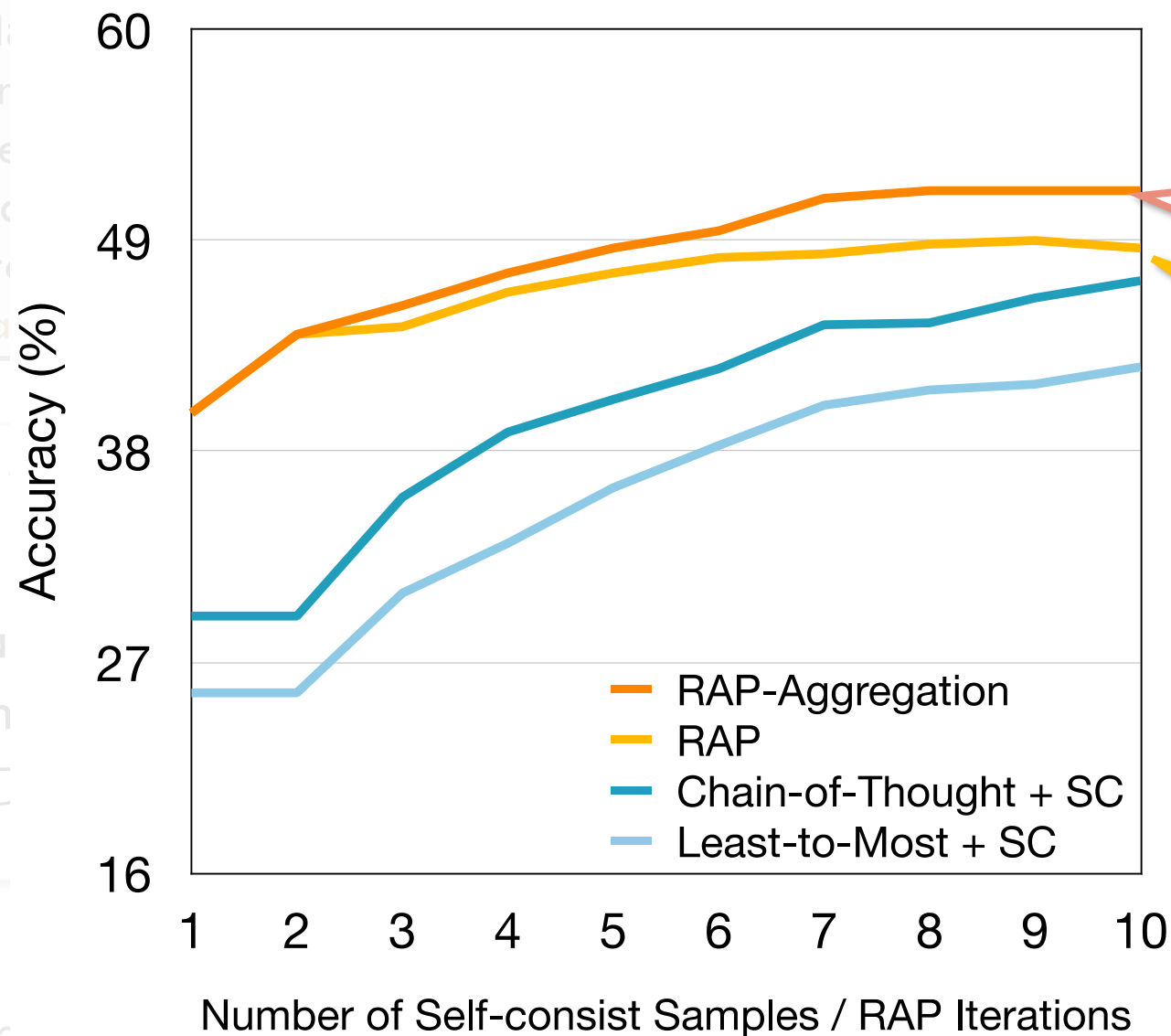
## LM as World Model:

Prompt LM to generate the next state given the current state and action

$$P(s' | s, a)$$



# Simulative reasoning



**action:** a sub-question for an unknown variable

**state:** intermediate values of variables

Aggregating multiple plans improves even further

Simulative reasoning outperforms autoregressive planning (CoT + Self-consistency)

Q1. A1:  $12 \times 2 = 24$

.....

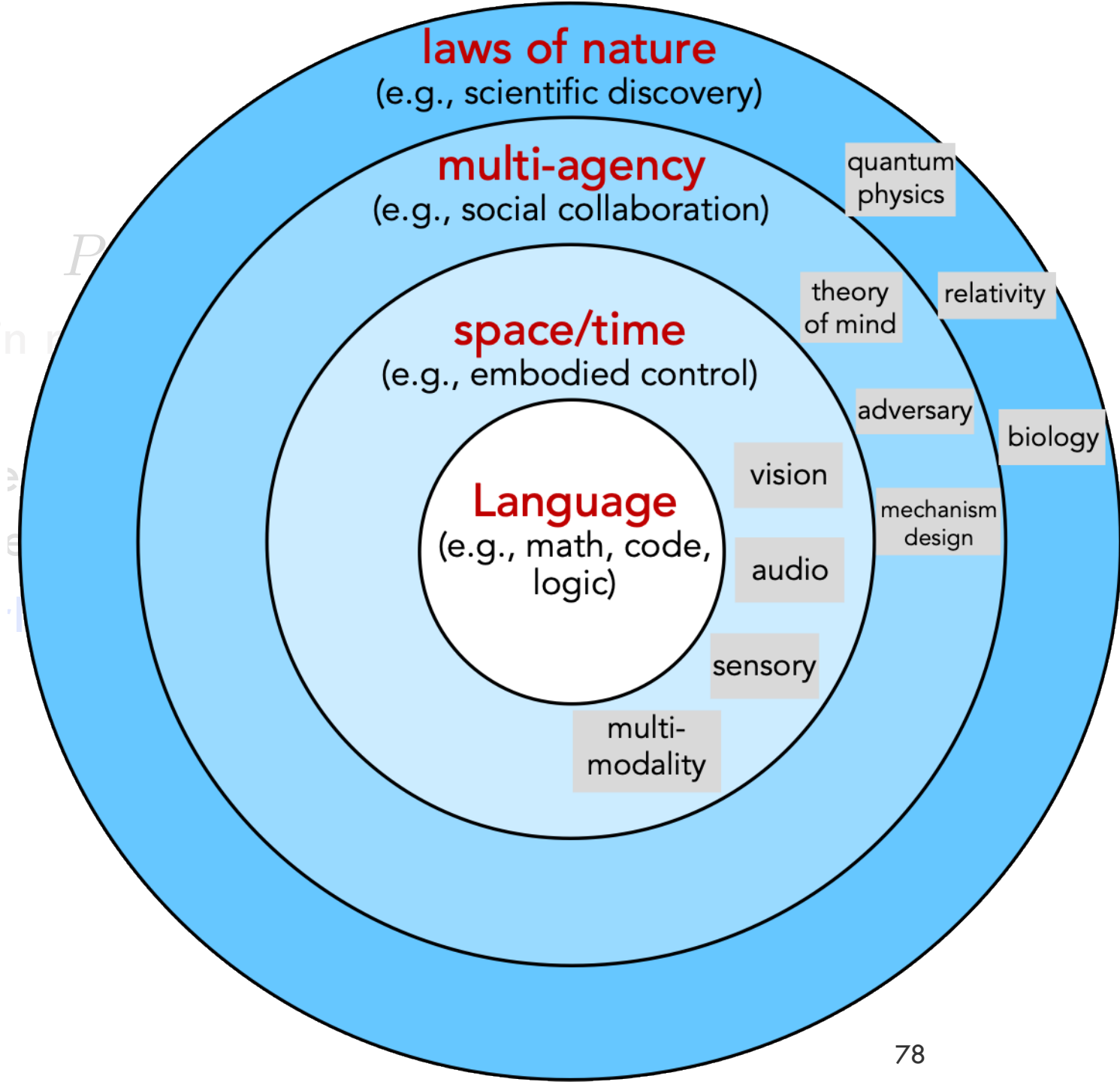
Qn. An:  $84 / 2 = 42$

(Answer: 42)

# Simulative reasoning

- Next “world” prediction
- Prior research (primarily in *robotics*)  
*specific world models*

The scope of simulation  
defines the capability of  
reasoning

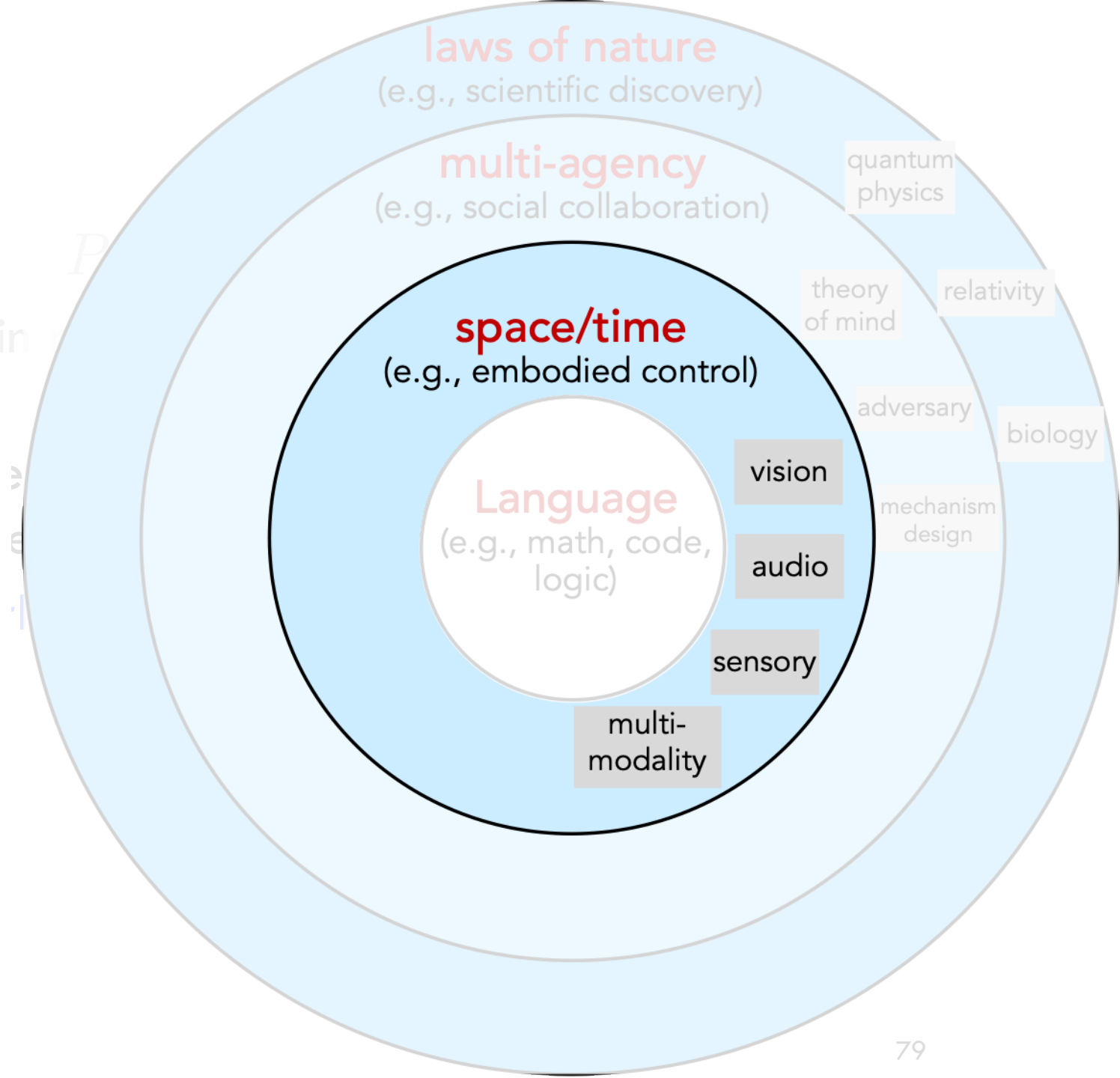


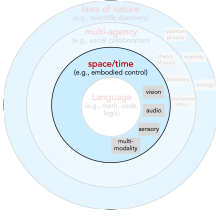


# Simulative reasoning

- Next “world” prediction
- Prior research (primarily in *specific world models*)

The scope of simulation defines the capability of reasoning





# Simulative reasoning beyond LM-based world models

- Language is often **not** the most efficient medium to describe all information during reasoning
- Other modalities (e.g., videos) can be more efficient



In auto-driving: describe the street state

- Vehicles' locations & movements



Pour liquid into a glass without spilling

- Viscosity & volume of the fluid
- shape & position of the container

# Simulative reasoning beyond LM-based world models

What's needed for a more general world model:

- 1) Integrating different spaces for simulation / reasoning: text, video, ...
- 2) **Generalist** language capability (like LLMs) + **generalist** vision capability (video pretraining)
- 3) Real-time control of the simulation through action inputs  $P(s'|s, a)$ 
  - Controllability allows to simulate many counterfactual worlds, and pick the best to actualize
  - Existing video-generation models (e.g., Sora) are not for this

# Simulative reasoning beyond LM-based world models

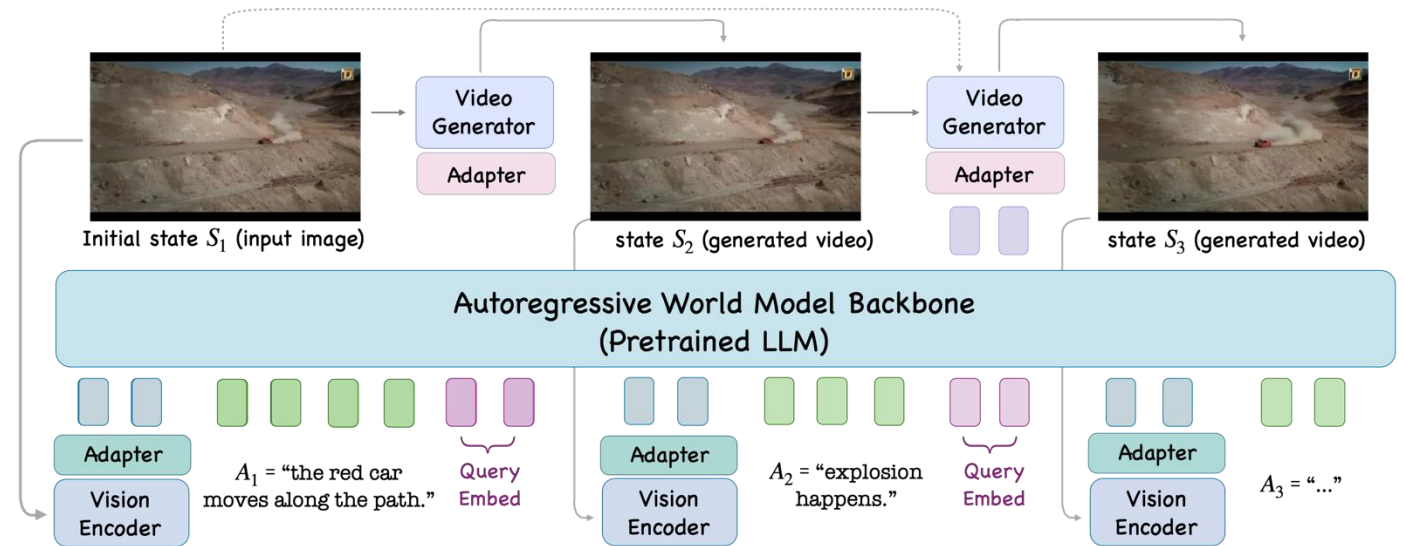
What's needed for a more general world model:

- 1) Integrating different spaces for simulation / reasoning: text, video, ...
- 2) **Generalist** language capability (like LLMs) + **generalist** vision capability (video pretraining)
- 3) Real-time control of the simulation through action inputs

$$P(s'|s, a)$$



[www.world-model.ai](http://www.world-model.ai)



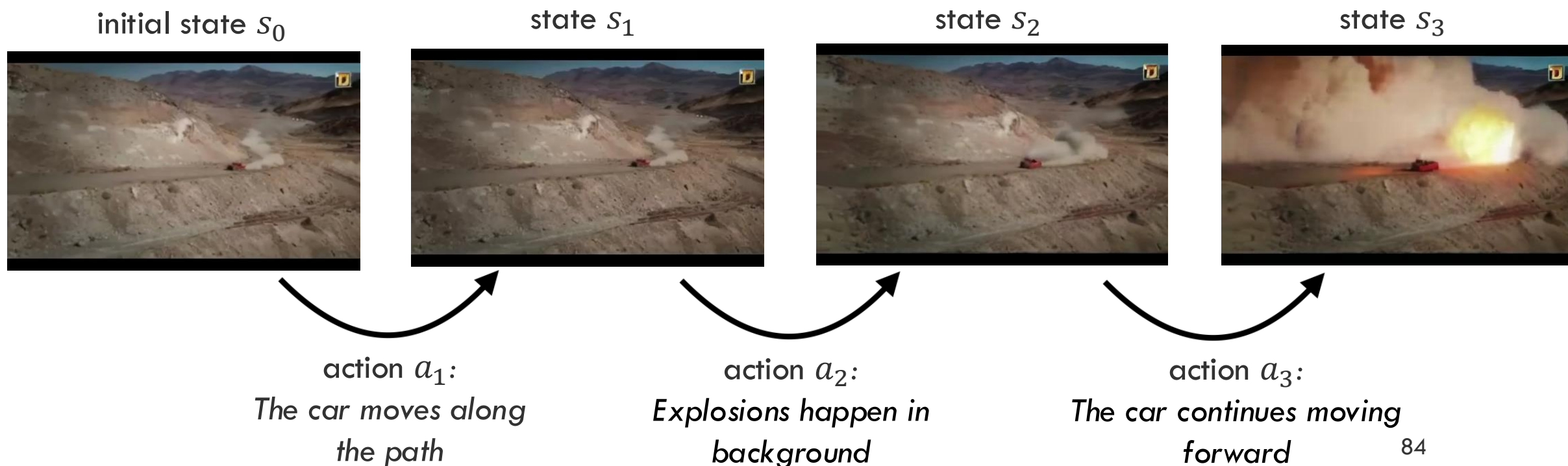
# Simulative reasoning beyond LM-based world models





# Pandora stepping towards more general world models

- 1) Integrating different spaces for simulation / reasoning: text, video, ...
- 3) Real-time control of the simulation through action inputs





# *Pandora* stepping towards more general world models

- 1) Integrating different spaces for simulation / reasoning: text, video, ...
- 3) Real-time control of the simulation through action inputs
  - Controllability allows to simulate many **counterfactual** worlds, and pick the best to actualize

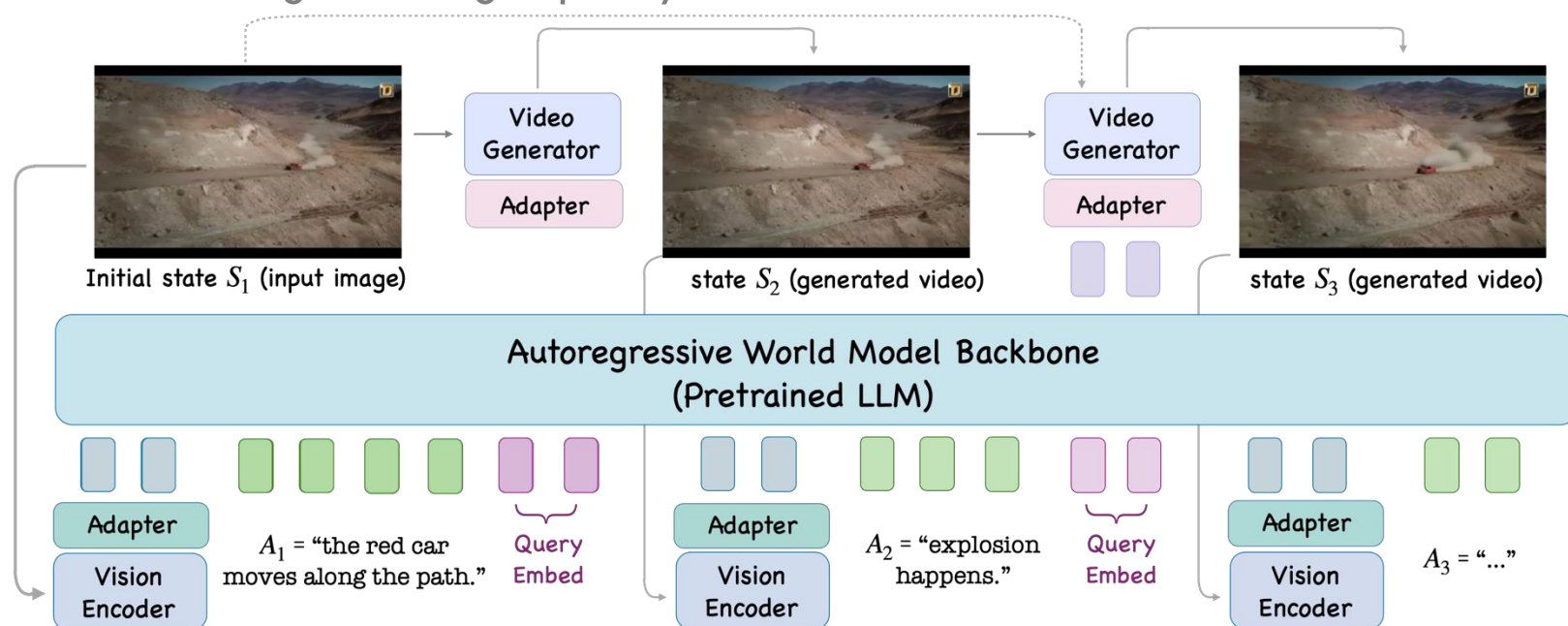
Action planning for robots





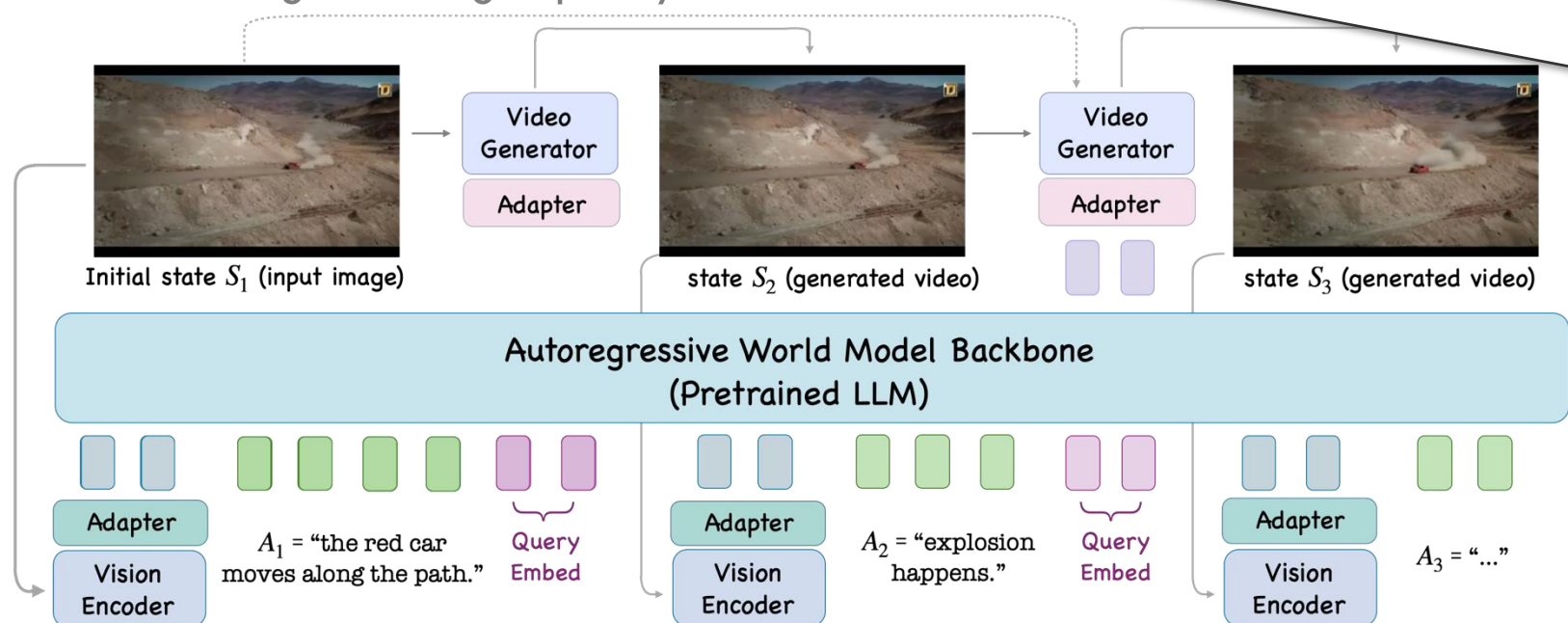
# Pandora stepping towards more general world models

- 2) Generalist language capability (like LLMs) + generalist vision capability (video pretraining)
- Generalist pretrained **LLM** as the autoregressive backbone
  - Generalist pretrained **Video Diffusion Model** for visual simulation
    - Massive video pretraining enables consistent prediction of the physical world states
  - Instruction-tuning connects the backbone with video generator for real-time control
    - Using small high-quality real-time control text-video data



# *Pandora*tepping towards more general world models

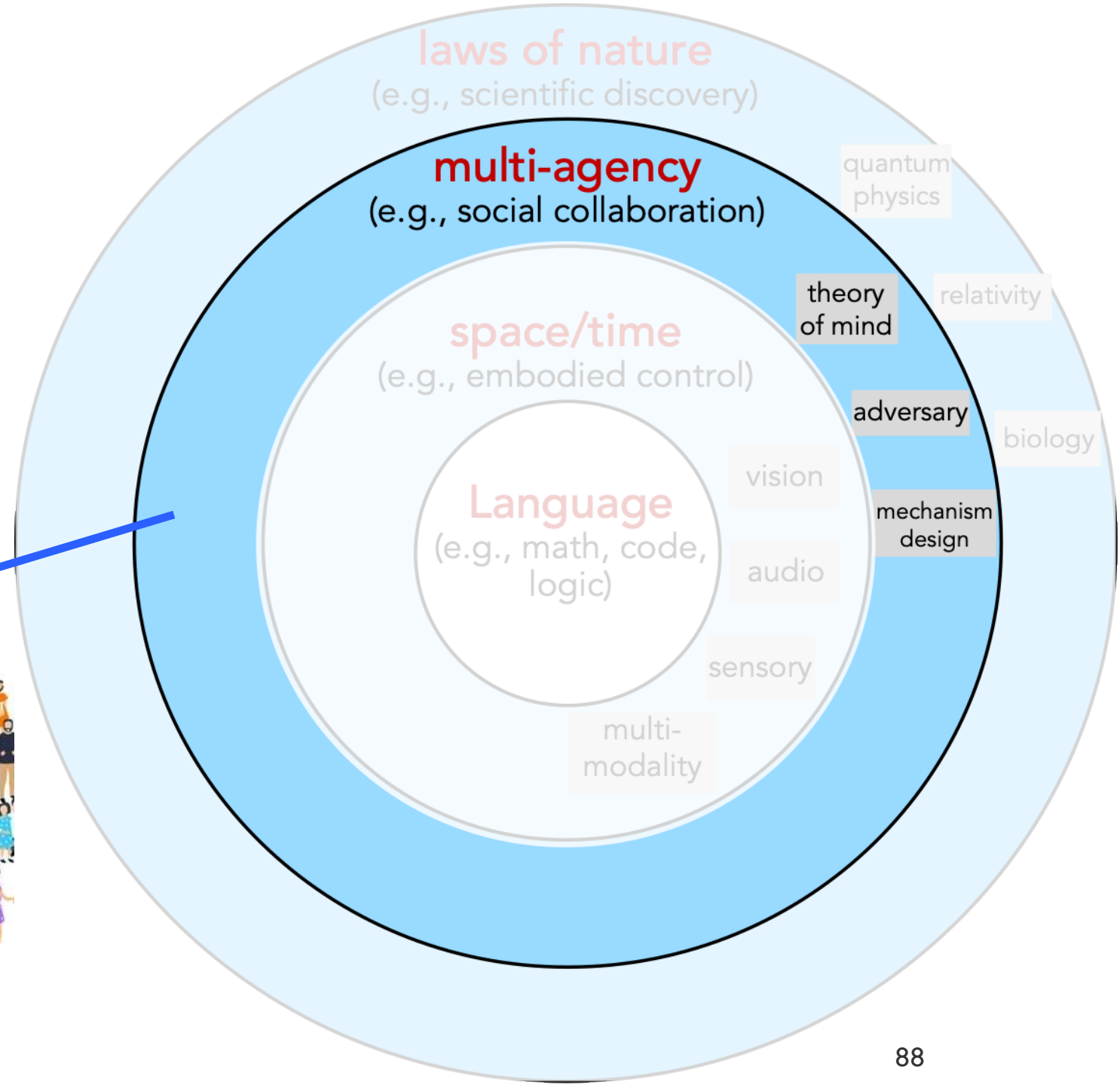
- 2) Generalist language capability (like LLMs) + generalist vision capability (video pretraining)
- Generalist pretrained **LLM** as the autoregressive backbone
  - Generalist pretrained **Video Diffusion Model** for visual simulation
    - Massive video **pretraining** enables consistent prediction of the physical world states
  - **Instruction-tuning** connects the backbone with video generator for real-time control
    - Using small high-quality real-time control text prompts



Pretraining and instruction-tuning stages mirror the common practice of LLM training

# Simulative reasoning

Society of individual  
world models



# Simulative reasoning



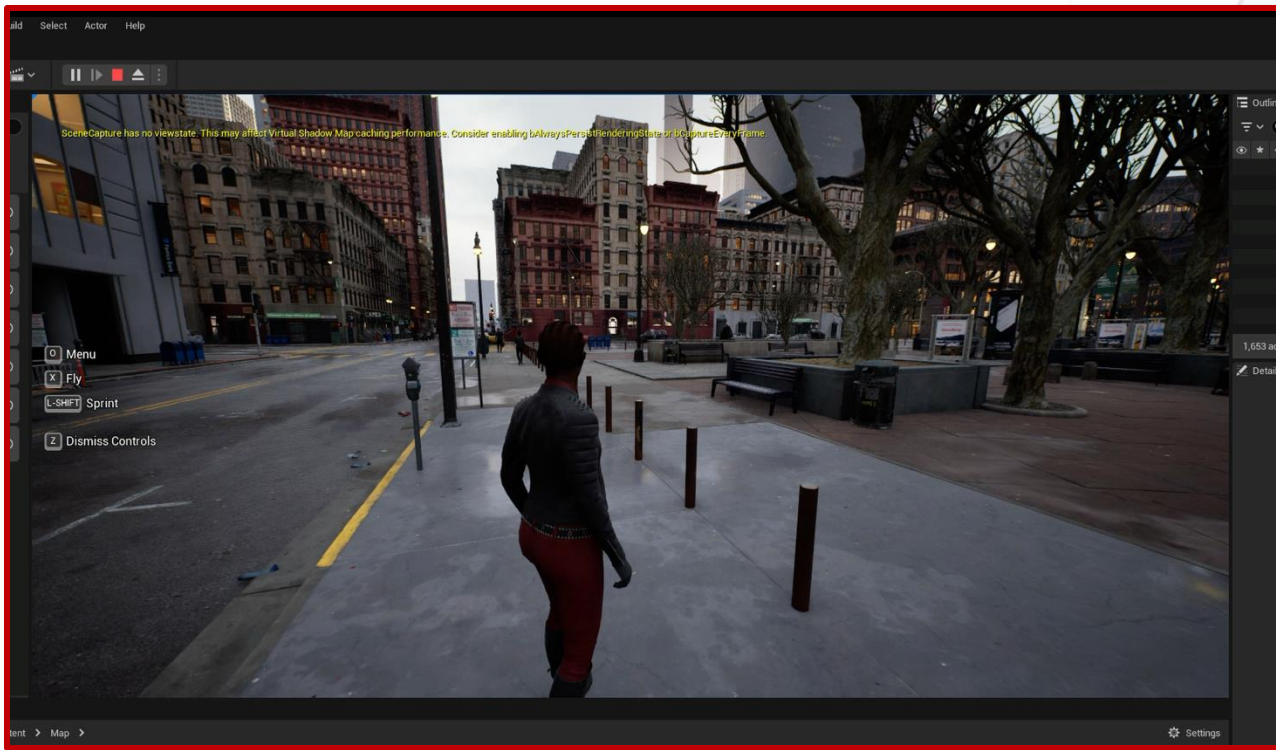
[Park et al., 2023]

25 agents, each controlled by individual LLM, converse with each other

- For studying emerging communication behaviors



# Simulative reasoning



In progress

## Richer and more realistic simulation of society

- Humans, vehicles, robots
- Simulating traffic, social, financial systems
- Could potentially be used for studying human-AI collaboration, education, social science, policy making, ...

laws of nature

(e.g., scientific discovery)

multi-agency

(e.g., social collaboration)

quantum physics

theory of mind

relativity

space/time

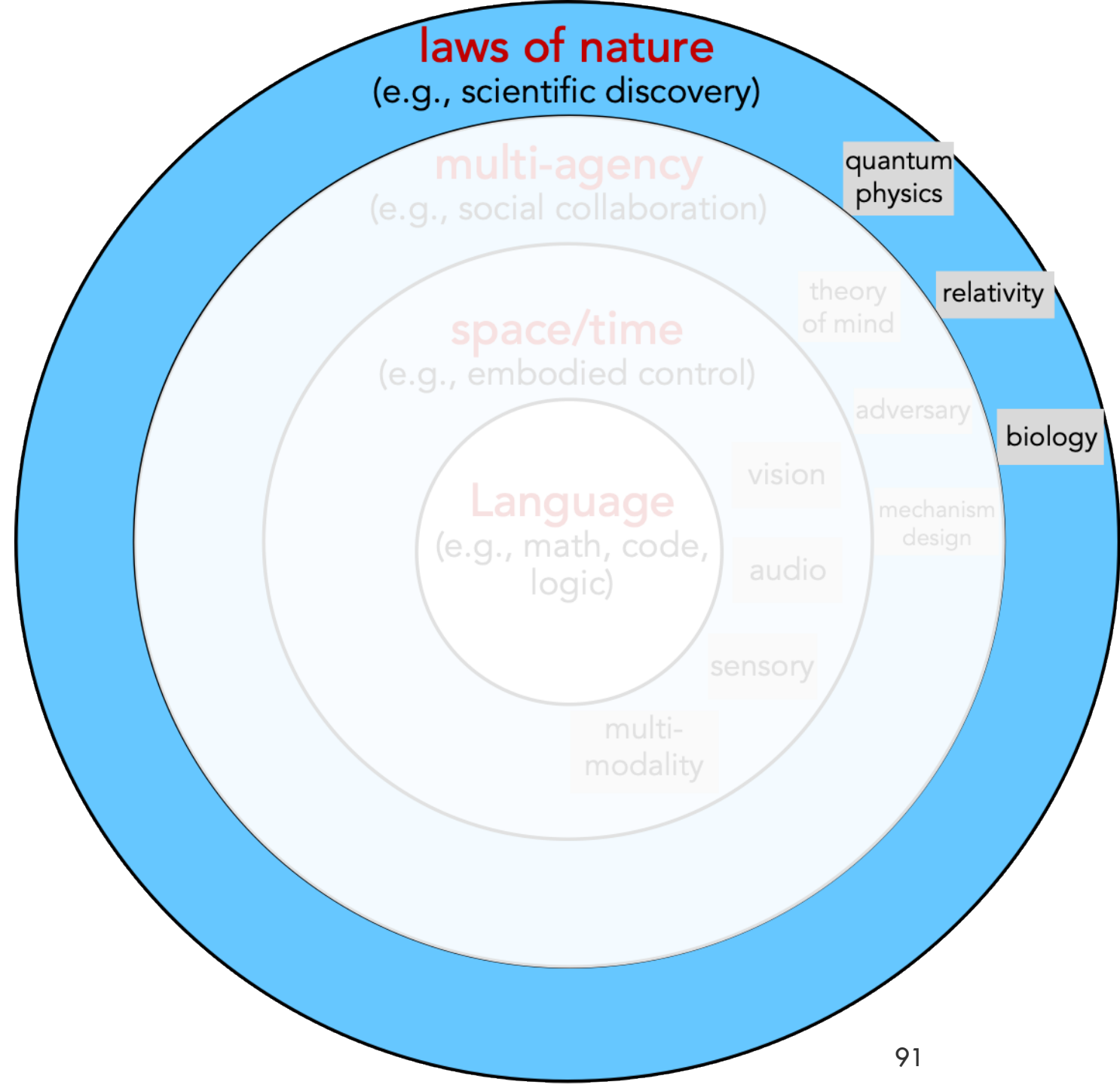
(e.g., embodied control)

Language

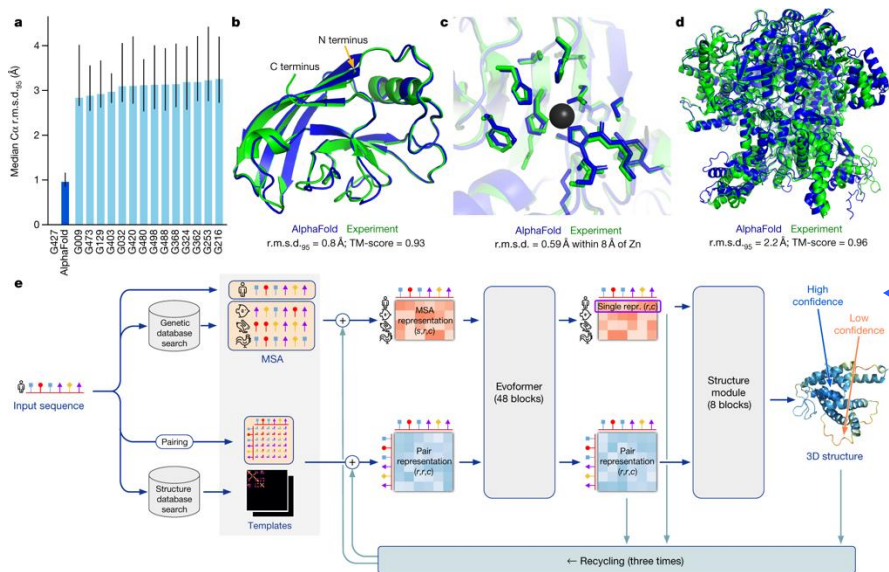
vision

mechanism design

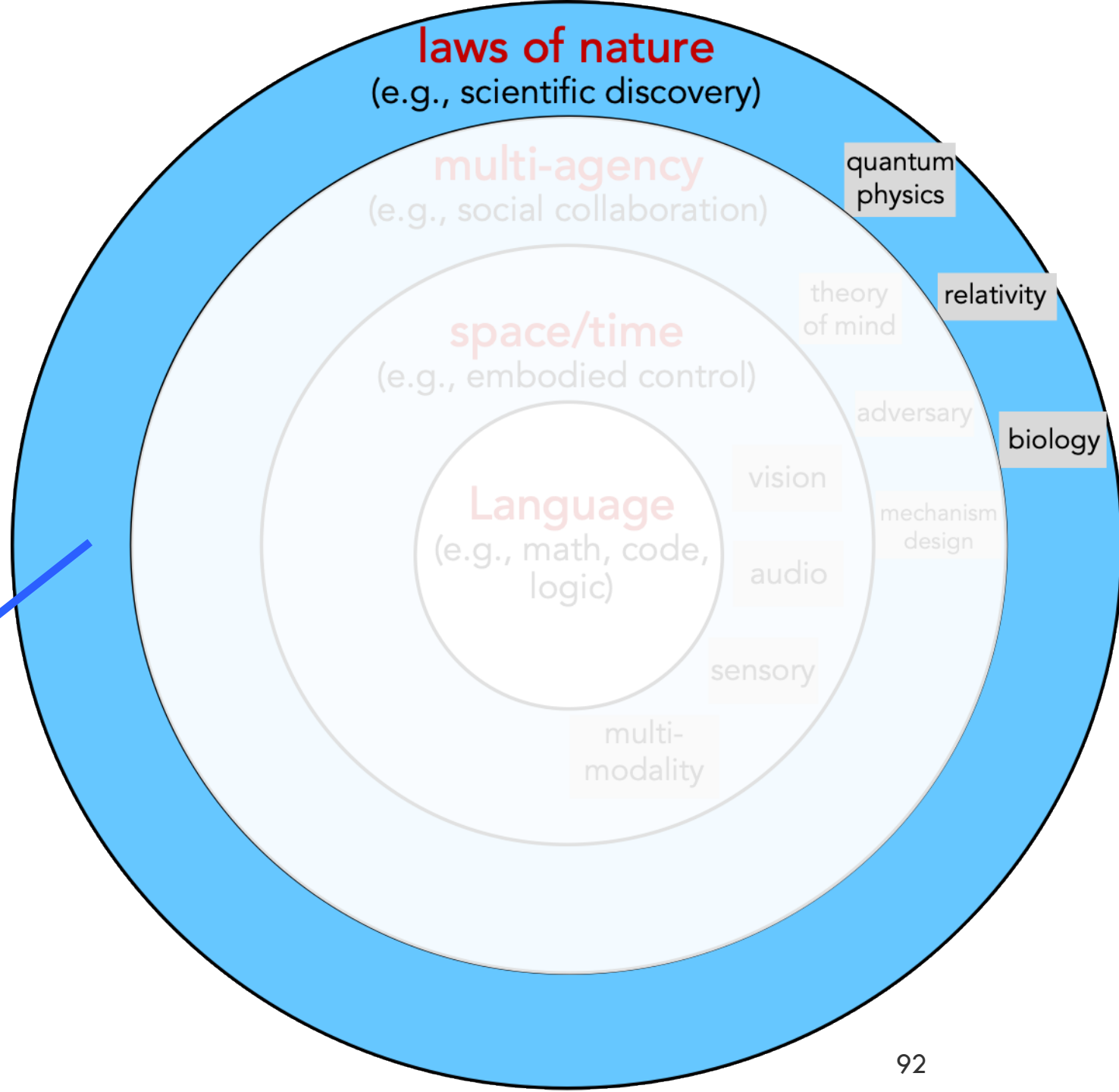
# Simulative reasoning



# Simulative reasoning

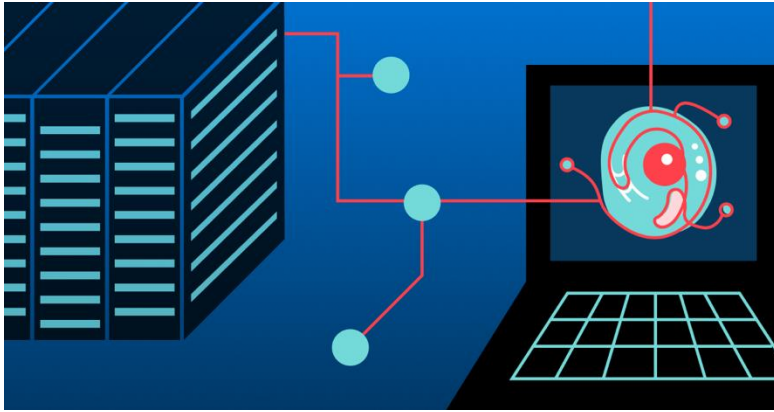


# AlphaFold for protein structure prediction

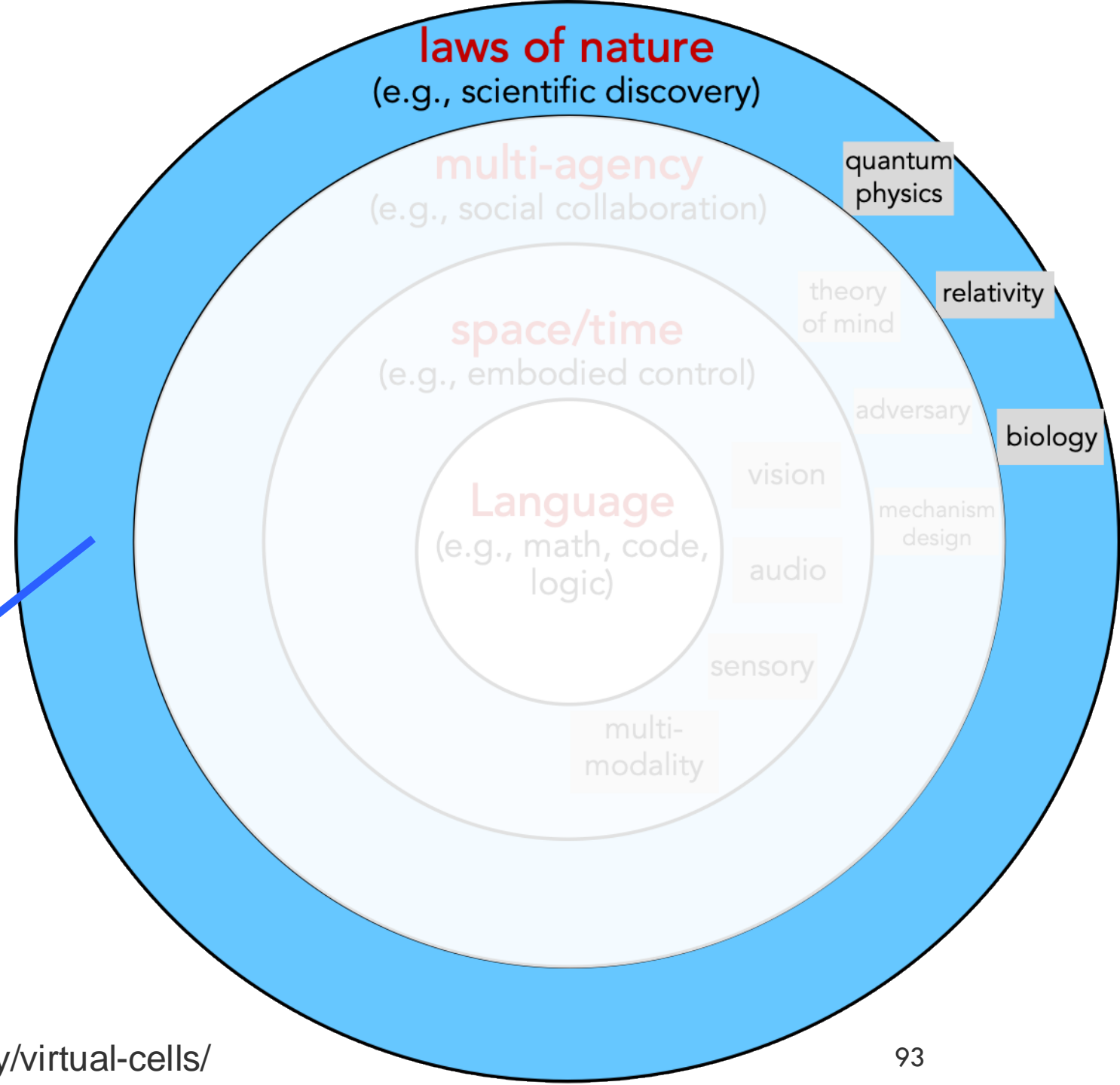




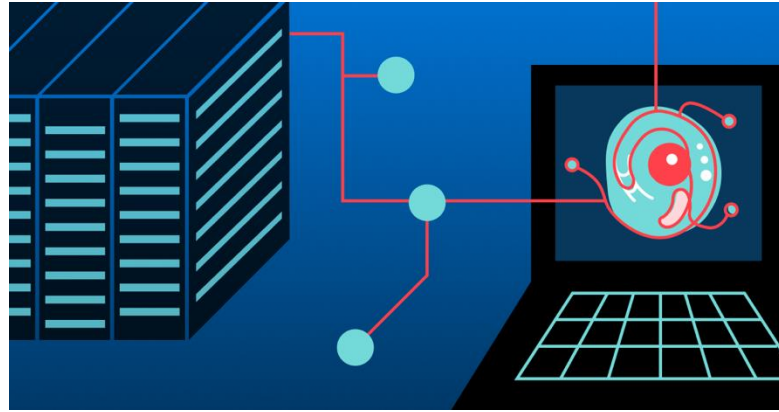
# Simulative reasoning



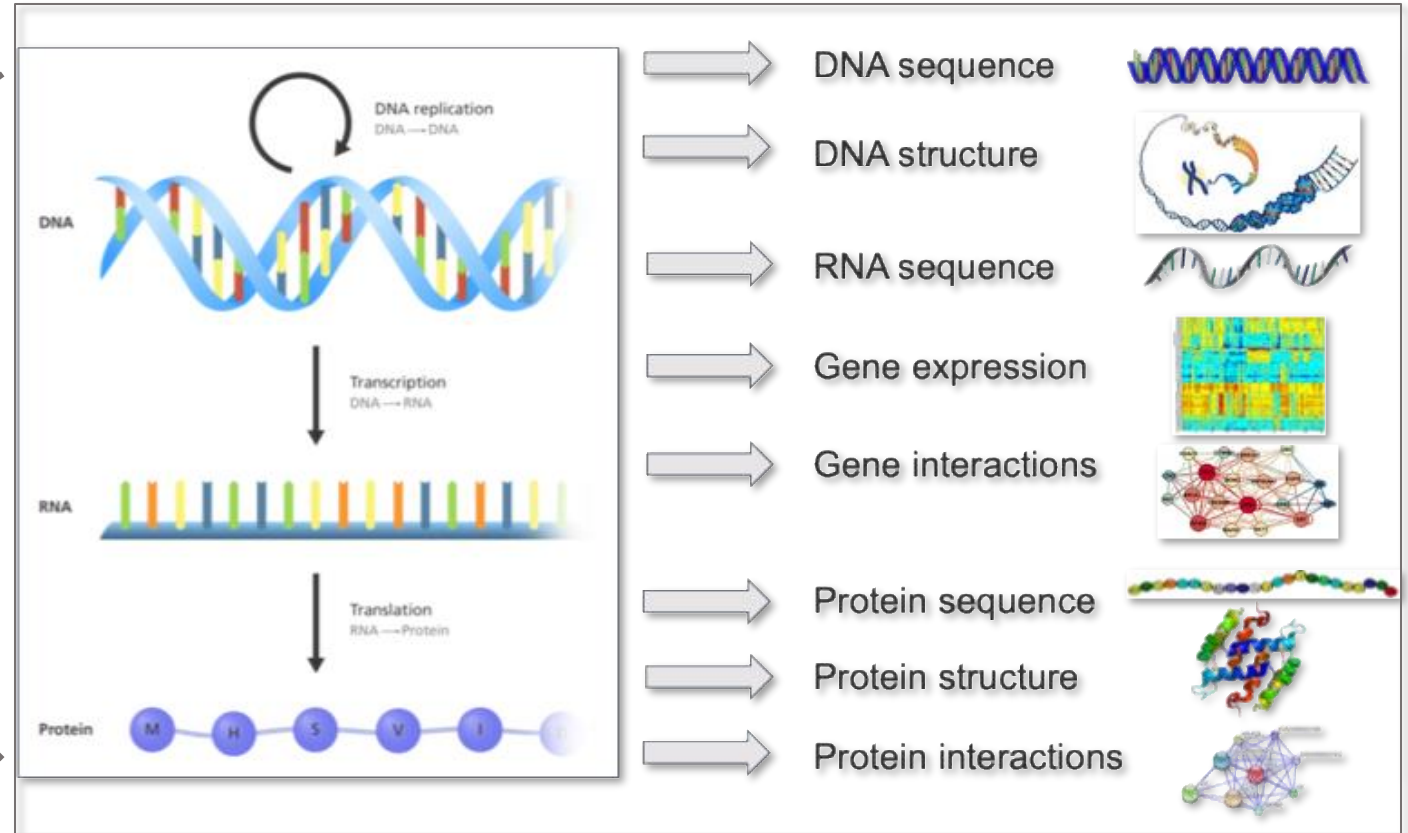
AI Virtual Cell  
for drug discovery, etc.



# Simulative reasoning



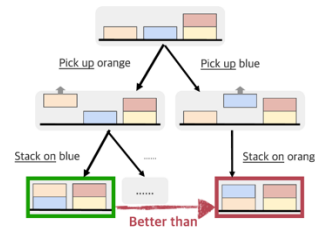
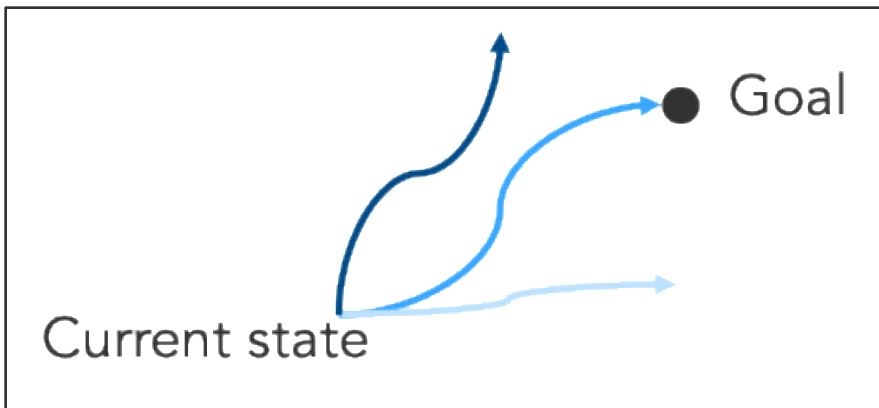
AI Virtual Cell  
for drug discovery, etc.



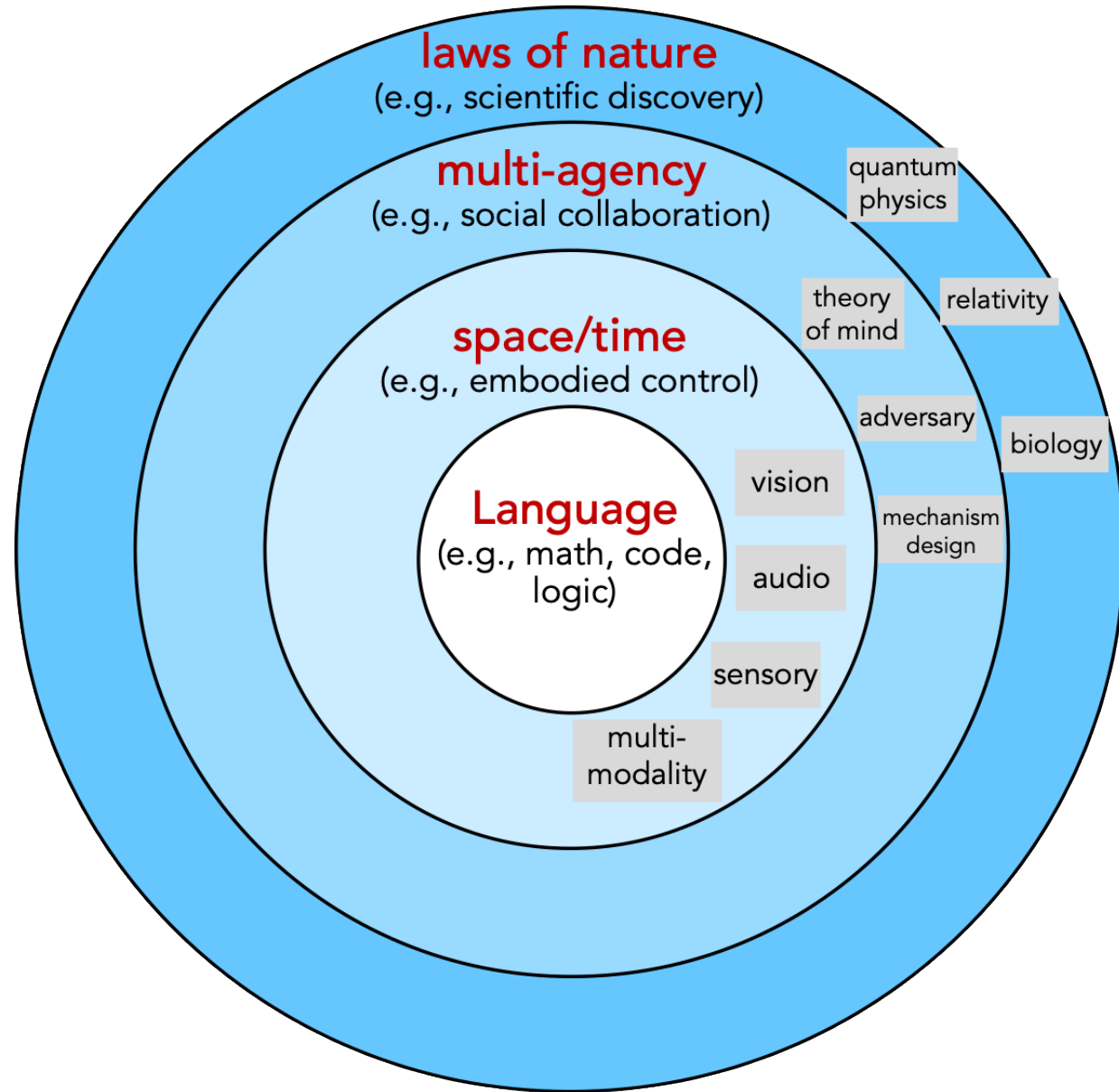
Extensive heterogenous biological data  
(In progress)

# Summary

- Simulative reasoning based on world models
  - strategic planning via simulation
- “More simulation, more intelligence”
- Building general world models
  - Pandora



*Pandora*



**Questions?**