

# World of World Modeling

Shane Gu 顾世(World)翔

Senior Staff RS, Google DeepMind

Feb 25, 2026

Stanford CS234

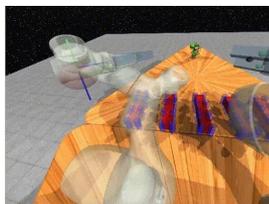
# Profile

- Japan-born Chinese Canadian, with PhD (2018) from UK and Germany
- Employments: Google Brain (2015-2022) → OpenAI (2023) → Google DeepMind (2023-present)
- Philosophy: Move on to a new field if **(impact remaining) / (top talents)** ratio becomes low

$$p(y_1, \dots, y_k) = \Gamma(k) \left( \prod_{i=1}^k \exp(x_i) \frac{y_i^{\tau}}{y_i^{\tau}} \right) \left( \sum_{i=1}^k \exp(x_i) \frac{y_i^{\tau}}{y_i^{\tau}} \right)^{-k} \tau^{k-1} \prod_{i=1}^k y_i^{-1}$$

$$= \Gamma(k) \tau^{k-1} \left( \sum_{i=1}^k \exp(x_i) / y_i^{\tau} \right)^{-k} \prod_{i=1}^k (\exp(x_i) / y_i^{\tau+1})$$

Gumbel-Softmax (2016)



Google Brain Dexterity  
Moonshot Co-Lead (2019)



OpenAI Japan entry (2023/02)



ChatGPT Post-Training (2023)

Rank	Delta	Model	Arena Elo
1	↑	Gemini-1.5-Pro-0014	1418
1	↓	GPT-4o-2024-08-13	1388
1	↑	Gemini-Advanced-0014	1380
2	↑	Yi-Large-0.5-720	1365
2	↑	Gemini-1.5-Flash-0014	1346
2	↑	Claude-3.5-Sonnet	1346

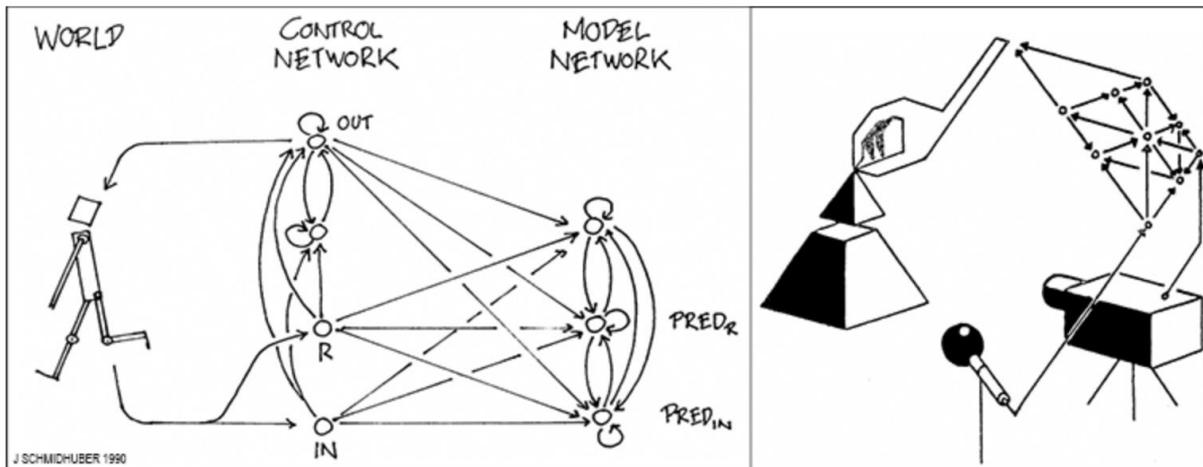
Gemini Post-Training  
Multilinguality Lead (2024-25)



Gemini Thinking / RL (2025-)  
Veo / Genie (2026-)



# What is a world model?



[Jürgen Schmidhuber](#) (Dec 2020, updated 2025)  
Pronounce: You\_again Shmidhoobuh

[AI Blog](#)  
Twitter: [@SchmidhuberAI](#)

## 1990: Planning & Reinforcement Learning with Recurrent World Models and Artificial Curiosity



Jürgen Schmidhuber:

**1990:** Introduced RNN-based world models for planning.

**1990:** Proposed high-dimensional, multi-vector reward signals.

**1990:** Formulated deterministic policy gradients for RNNs.

**1990:** Invented adversarial artificial curiosity (foundation of GANs).

**1991:** Introduced neural network distillation and computational consciousness.

**2004-2005:** Applied world models to physical AI and self-healing robots.

# What is a world model?



## World Models

(2018)

David Ha<sup>1</sup> Jürgen Schmidhuber<sup>2,3</sup>

At each time step, our agent receives an **observation** from the environment.

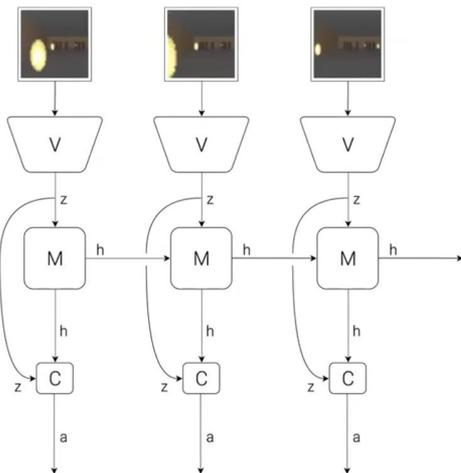
### World Model

The **Vision Model (V)** encodes the high-dimensional observation into a low-dimensional latent vector.

The **Memory RNN (M)** integrates the historical codes to create a representation that can predict future states.

A small **Controller (C)** uses the representations from both **V** and **M** to select good actions.

The agent performs **actions** that go back and affect the environment.



Jitendra Malik: “I think the main contribution of this paper is that it introduced the term ‘world models.’ This term caught on, but I’m not a particular fan of it because there was already a good term in use: ‘dynamics model,’ which control theorists started using in 1960. The consequence of using the term ‘world models’ is that we now have so much confusion about which definition of a world model we are discussing.”

2748 Robust and Interactable World Models in Computer Vision

Unlisted



ComputerVisionFoundation Videos

44.7K subscribers

Subscribe

366

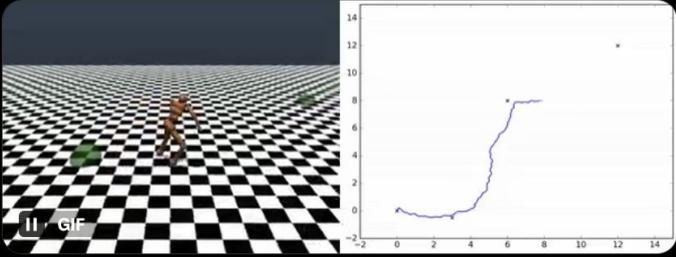


24,050 views Nov 9, 2025

# What is a world model?

**Shane Gu** @shaneguML · Nov 5, 2019

Our work on **model-based RL** (a.k.a. **world models** these days), an impressive first paper from Archit in Google AI residency! We show a single objective can do predictable skill discovery and **model learning** simultaneously! [arxiv.org/abs/1907.01657](https://arxiv.org/abs/1907.01657)



The GIF shows a 3D environment with a checkered floor and a small figure. To the right is a line graph showing performance over time. The graph has a y-axis from -2 to 14 and an x-axis from 0 to 14. The performance starts at 0, dips slightly, then rises sharply to about 8 at x=6, and levels off at 8.

33 118

**Shane Gu** @shaneguML · Apr 28, 2025

Branding is important to inspire researchers



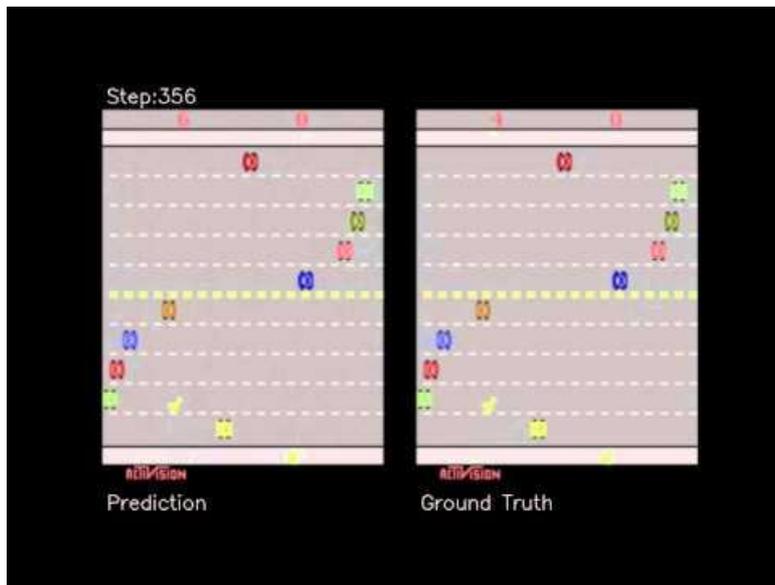
The meme image shows two panels of Drake. The top panel has the text "MODEL-BASED RL" and the bottom panel has the text "WORLD MODEL".

7 9 212 12K



Shane Gu:  
“World model is the ‘model’  
in model-based RL.”

# What is a world model?



---

## Action-Conditional Video Prediction using Deep Networks in Atari Games

---

2015

Junhyuk Oh Xiaoxiao Guo Honglak Lee Richard Lewis Satinder Singh  
University of Michigan, Ann Arbor, MI 48109, USA  
{junhyuk, guoxiao, honglak, rickl, baveja}@umich.edu

Abstract

---

## Unsupervised Learning for Physical Interaction through Video Prediction

---

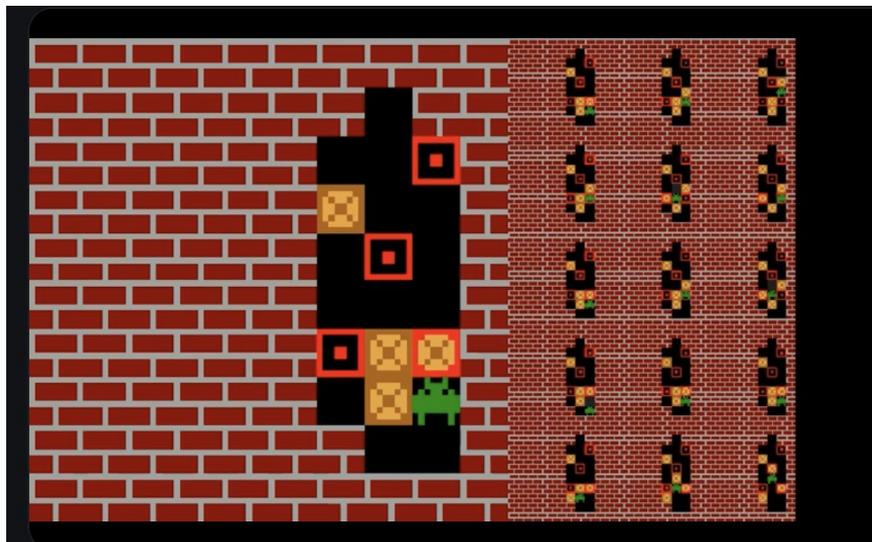
2016

Chelsea Finn\*  
UC Berkeley  
cbfinn@eecs.berkeley.edu

Ian Goodfellow  
OpenAI  
ian@openai.com

Sergey Levine  
Google Brain  
UC Berkeley  
slevine@google.com

# What is a world model?



---

## Imagination-Augmented Agents for Deep Reinforcement Learning

---

2017

Théophane Weber\* Sébastien Racanière\* David P. Reichert\* Lars Buesing  
Arthur Guez Danilo Rezende Adria Puigdomènech Badia Oriol Vinyals  
Nicolas Heess Yujia Li Razvan Pascanu Peter Battaglia  
Demis Hassabis David Silver Daan Wierstra  
DeepMind



## Mastering Atari with Discrete World Models

Danijar Hafner Timothy Lillicrap Mohammad Norouzi Jimmy Ba

ICLR 2021

2021

# Talk Outline

- What is prediction?
  - Solomonoff induction and universal intelligence
  - Causality and OOD generalization
  - Empowerment and 3 levels of predictability maximization
- Forward and inverse world models
  - Shooting and collocation methods
  - Goal-conditioned value functions and (generalized) decision transformers
- Physical and symbolic world models
  - Video models as reasoners
  - Futures of world modeling
- (Forgot) Control as inference: e.g. particle smoothing as optimal control, MCTS vs beam search

**I won't talk about any Gemini / Veo / Genie. Please check out new models! Or ask me questions after the class.**

What is prediction?

# What's the "best" world model?

the "induction" machine: given data, infer the **rules / causes / programs**

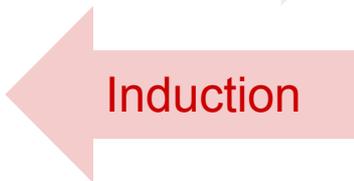
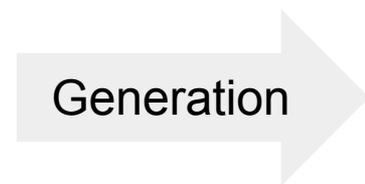
program:  $P$

```
#include <stdio.h>

int main()
{
    printf("Hello World");
    return 0;
}
```

e.g.

$$D_{t+1} = D_t + 2$$



data:  $D_{t-1}, D_t, D_{t+1}, \dots$



$\{1, 3, 5, 7, \dots\}$

# Solomonoff Induction

- Formalization of Bayes theorem in computation theory

$$p(P|D) \propto p(D|P)p(P)$$

- Occam's razor (related: Kolmogorov Complexity, Minimum Description Length)

$$p(P) \propto 2^{-|P|}$$

- Inspired:

C-Test [[Hernandez-Orallo et al1998](#)]

	<i>Sequence Prediction Test</i>	
Complexity	Sequence	Answer
9	a, d, g, j, -, ...	m
12	a, a, z, c, y, e, x, -, ...	g
14	c, a, b, d, b, c, c, e, c, d, -, ...	d

Universal Intelligence [[Legg & Hutter 2007](#)]

$$\sum_{\mu \in E} 2^{-K(\mu)} V_{\mu}^{\pi}$$

environment

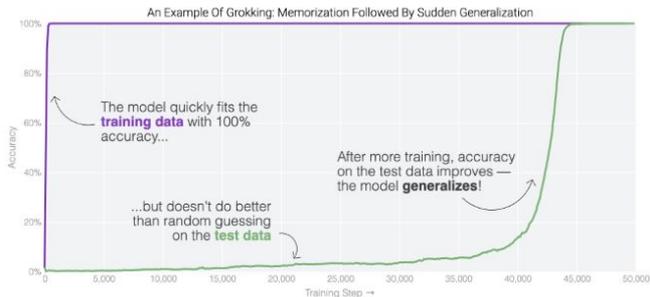
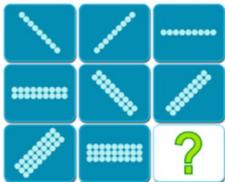


# Solomonoff Induction: Prediction = Understanding

- "The best prediction is inferring the shortest program that reproduces the data." (Ilya Sutskever)
- Prediction = Understanding = Inferring the causal program of a phenomenon.
- Examples: IQ tests, "Aha!" moments, scientific discoveries, deep learning (Grokking), and physics (Phase Transitions).
  - Text prediction yields emotional understanding (OpenAI Sentiment Neuron), while video/audio prediction yields physical and emotional understanding (Google Veo3).
- Despite hardware differences, humans and AI models both build intelligence through prediction—mastering prediction means mastering understanding.

Which number logically follows this series: 4 - 6 - 9 - 6 - 14 - 6 - ...

- 6
- 17
- 19
- 21



April 6, 2017 Publication

## Unsupervised sentiment neuron

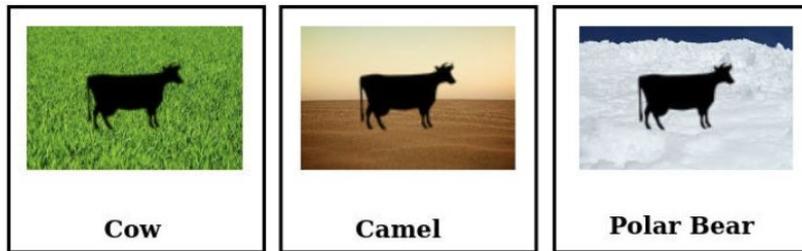
We've developed an unsupervised system which learns an excellent representation of sentiment, despite being trained only to predict the next character in the text of Amazon reviews.

# Causality: Understanding = OOD Generalization

- Invariant Risk Minimization (2019): non-causal spurious data biases prevent true OOD generalization in ML
- Generalization can be achieved by having data coming from different interventions of the true causal graph → “Diversity is all you need” → GPT-3, i.e. train to predict everything.

## Invariant Risk Minimization

Martin Arjovsky, Léon Bottou, Ishaan Gulrajani, David Lopez-Paz



**Neural Network Predictions**

# Empowerment and 3 levels of prediction

- Level 1: Passively fit your world model on world
- Level 2: Actively fit your world model on world
- Level 3: Actively fit world to your world model
- Notations:  $s$  = future of the world,  $z$  = your actions
  - Empowerment: Daniel Polani 2005, Shakir Mohammed 2015

$$\mathcal{I}(s; z) \geq \mathbb{E}_{z \sim p(z), s \sim p(s|z)} [\log q_{\phi}(s|z) - \log p(s)]$$

---

## Variational Empowerment as Representation Learning for Goal-Based Reinforcement Learning

---

# Level 1: Passively fit your world model on world

- Pre-training, Supervised learning, Self-supervised learning, Generative modeling
- Data distribution is **stationary** during predictive training

## How Much Information is the Machine Given during Learning?

Y. LeCun

### ▶ “Pure” Reinforcement Learning (**cherry**)

- ▶ The machine predicts a scalar reward given once in a while.

### ▶ **A few bits for some samples**

### ▶ Supervised Learning (**icing**)

- ▶ The machine predicts a category or a few numbers for each input
- ▶ Predicting human-supplied data
- ▶ **10→10,000 bits per sample**

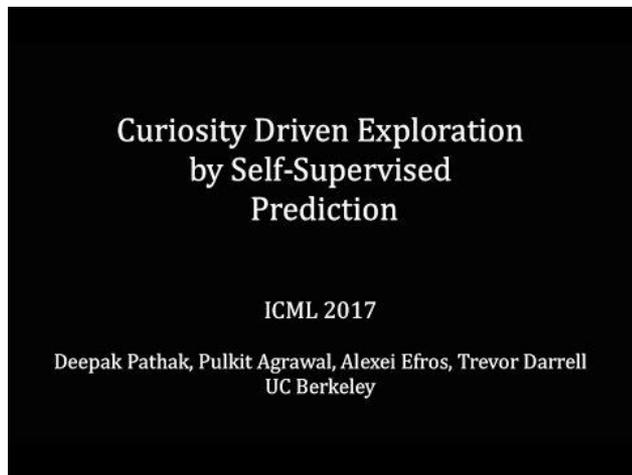
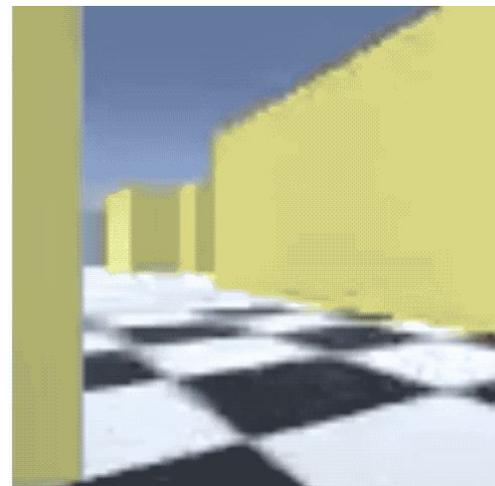
### ▶ Self-Supervised Learning (**cake génoise**)

- ▶ The machine predicts any part of its input for any observed part.
- ▶ Predicts future frames in videos
- ▶ **Millions of bits per sample**



# Level 2: Actively fit your world model on world

- Post-training, DAgger, GAIL, active learning
- Data distribution is **non-stationary** during predictive training



DPM (Stadie 2016) and ICM (Pathak 2017) incentivizes novelty

$$I(S_{t+1}; \Theta | \xi_t, a_t) = \mathbb{E}_{s_{t+1} \sim \mathcal{P}(\cdot | \xi_t, a_t)} [D_{\text{KL}}[p(\theta | \xi_t, a_t, s_{t+1}) || p(\theta | \xi_t)]]$$

$$r'(s_t, a_t, s_{t+1}) = r(s_t, a_t) + \eta D_{\text{KL}}[p(\theta | \xi_t, a_t, s_{t+1}) || p(\theta | \xi_t)]$$

VIME (Houthoofd et al 2017) tackles noisy TV problem

# Level 3: Actively fit world to your world model

- Perform actions in the world to make the world mode predictable wrt you / your own world model
  - Politicians
  - Financial firms
  - X Influencers
- **Difficult objective. Nobody has cracked this yet at scale.**

$$\mathcal{I}(s; z) \geq \mathbb{E}_{z \sim p(z), s \sim p(s|z)} [\log q_\phi(s|z) - \log p(s)]$$

---

**Variational Empowerment as Representation Learning  
for Goal-Based Reinforcement Learning**

---

Jongwook Choi<sup>†1</sup> Archit Sharma<sup>†2</sup> Honglak Lee<sup>1,3</sup> Sergey Levine<sup>4,5</sup> Shixiang Shane Gu<sup>4</sup>

Published as a conference paper at ICLR 2020

---

**DYNAMICS-AWARE UNSUPERVISED DISCOVERY OF  
SKILLS**

Archit Sharma\*, Shixiang Gu, Sergey Levine, Vikash Kumar, Karol Hausman  
Google Brain  
{architsh, shanegu, slevine, vikashplus, karolhausman}@google.com

Forward and inverse world models

# What is a (World) Model?

Forward Models

$$s_{t+1} = F(s_t, a_t)$$

$$a = \Pi(s, F(s, a)) \quad \forall s, a$$

$$a_t = \Pi(s_t, s_{t+1})$$

$$0 = Q(s, a, F(s, a)) \quad \forall s, a$$

$$0 = Q(s_t, a_t, s_{t+1})$$

Inverse Models

\* assume deterministic for simplicity

\*\* will explain why I chose those symbols

# Shooting vs Direct Collocation for Planning

Shooting:

$$\arg \max_{a_{t:t+T}} \sum_{i=t}^{t+T} r(s_i) \quad \text{where} \quad s_{i+1} = F(s_i, a_i)$$

Direction Collocation (Constrained Optimization):

$$\arg \max_{s_{t:t+T}} \sum_{i=t}^{t+T} r(s_i) \quad \text{s.t.} \quad -|A| \leq \Pi(s_i, s_{i+1}) \leq |A|$$

$$\arg \max_{s_{t:t+T}, a_{t:t+T}} \sum_{i=t}^{t+T} r(s_i) \quad \text{s.t.} \quad Q(s_i, a_i, s_{i+1}) = 0$$

# Direct Collocation: Contact Invariant Optimization (CIO)

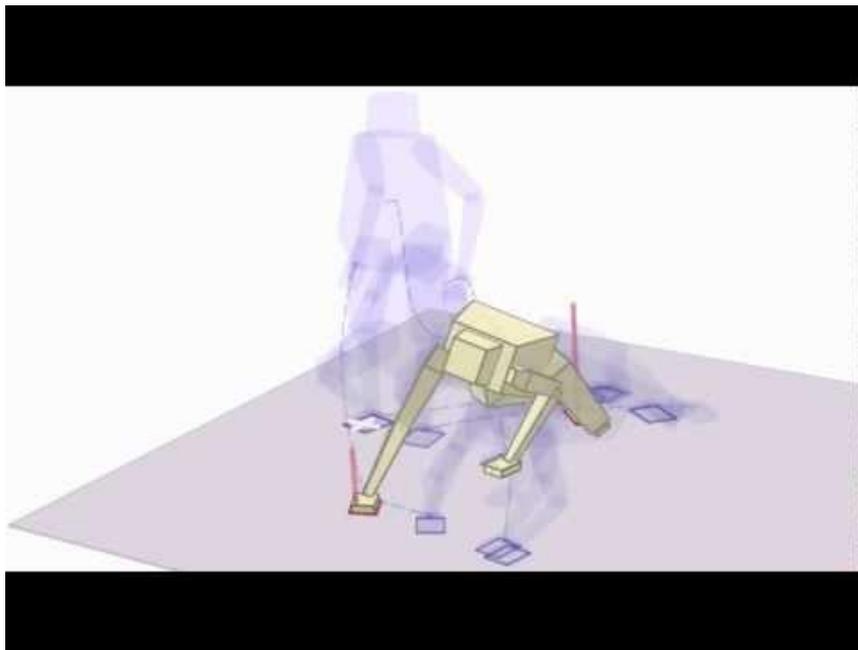
- Constrained optimization wrt dynamics  
constraints := relaxing (hacking) dynamics  
during planning
- “First solve task, then fix physics”

[[Mordatch et al 2012](#)  
([SIGGRAPH](#))]

$$\tau(\mathbf{q}, \dot{\mathbf{q}}, \ddot{\mathbf{q}}) = J(\mathbf{q})^T \mathbf{f} + B\mathbf{u}$$

$$L_{\text{Physics}}(\mathbf{s}) = \sum_t \left\| J_t(\mathbf{s})^T \mathbf{f}_t(\mathbf{s}) + B\mathbf{u}_t(\mathbf{s}) - \tau_t(\mathbf{s}) \right\|^2$$

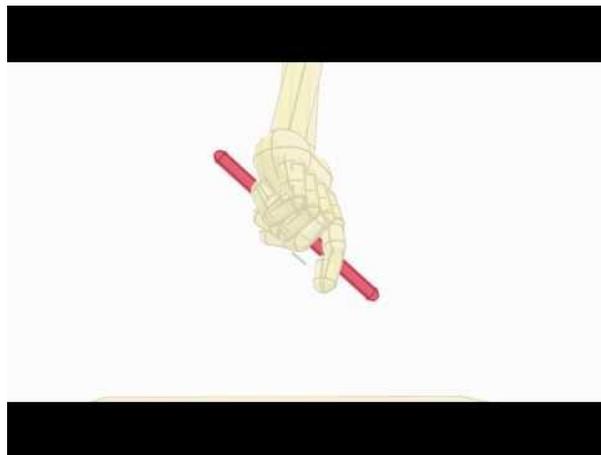
$$L_{\text{CI}}(\mathbf{s}) = \sum_t c_{i,\phi(t)}(\mathbf{s}) (\|\mathbf{e}_{i,t}(\mathbf{s})\|^2 + \|\dot{\mathbf{e}}_{i,t}(\mathbf{s})\|^2)$$



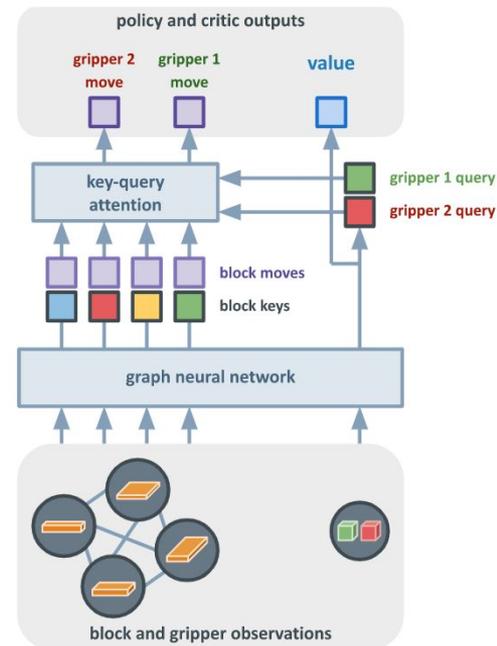
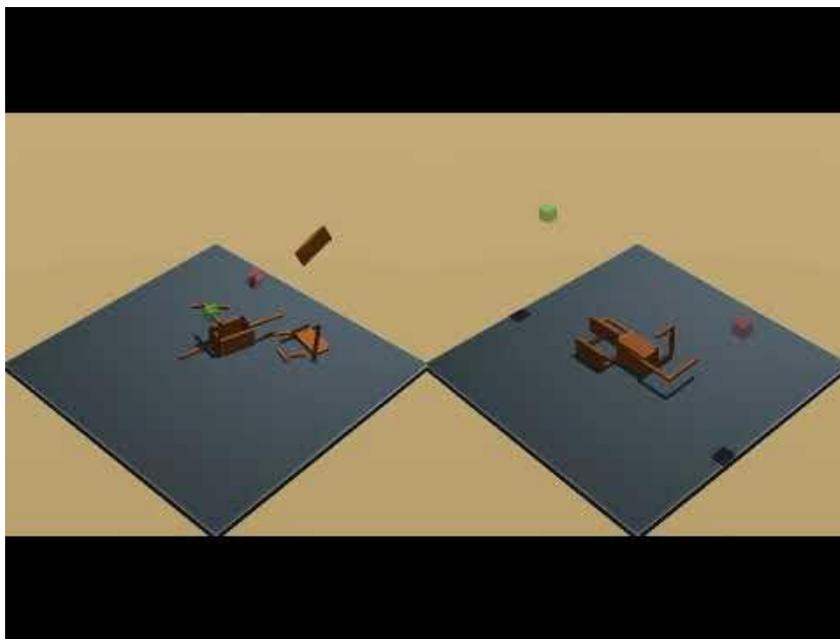
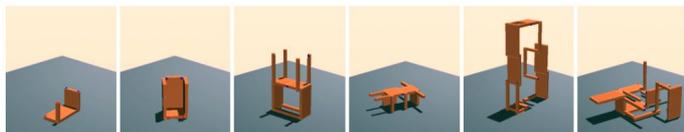
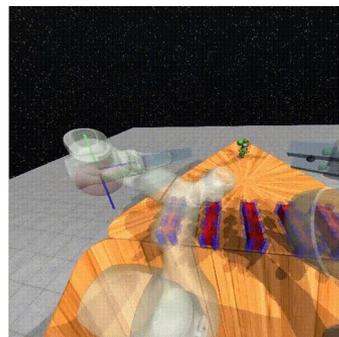
# Direct Collocation: Contact Invariant Optimization (CIO)

- Direct collocation solves contact-rich tasks with minimal reward shaping, i.e. dynamics relaxation provides appropriate “reward shaping”
- **Analogy:** autoregressive video diffusion vs bidirectional video diffusion

[\[Mordatch et al 2012+  
\(SIGGRAPH\)\]](#)



# Dexterity with shooting method is hard

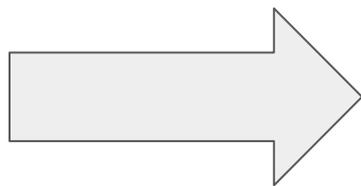


Kamyar Ghasemipour, Byron David, Daniel Freeman, Shixiang Shane Gu, Satoshi Kataoka, Igor Mordatch. "Learning to Assemble with Large-Scale Structured Reinforcement Learning." (2022)

Back to the notations...

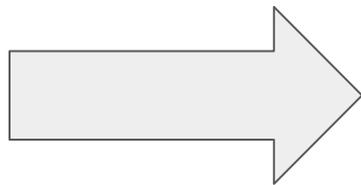
$$s_{t+1} = F(s_t, a_t)$$

$$a_t = \Pi(s_t, s_{t+1})$$



a policy  $|s_{t+1}$

$$0 = Q(s_t, a_t, s_{t+1})$$



a Q-function  $|s_{t+1}$

# Temporal Difference Models (TDMs)

[[Pong et al 2018 \(ICLR\)](#)]

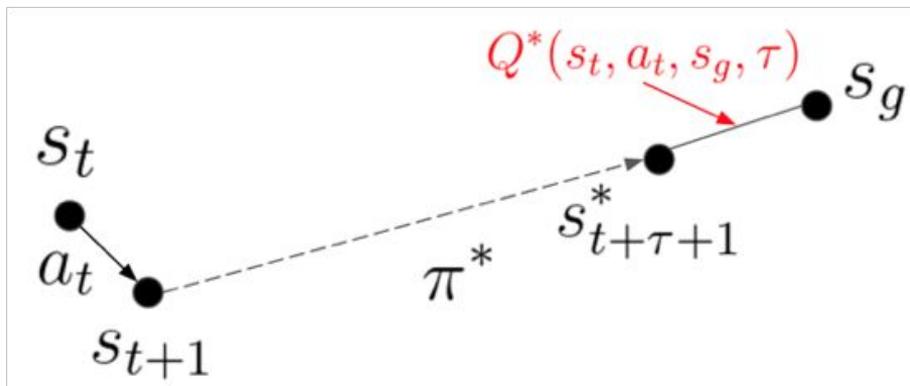
- Goal-conditioned optimal policy or Q-function as an implicit world model (temporally-extended)
- Use hindsight relabeling trick
- **Value functions are world models with different time scale and representation**

$$r_d(s_t, a_t, s_{t+1}, s_g, \tau) = -D(s_{t+1}, s_g) \mathbb{1}[\tau = 0]$$

---

$$Q(s_t, a_t, s_g, \tau) = \mathbb{E}_{p(s_{t+1}|s_t, a_t)}[-D(s_{t+1}, s_g) \mathbb{1}[\tau = 0] + \max_a Q(s_{t+1}, a, s_g, \tau - 1) \mathbb{1}[\tau \neq 0]]$$

---



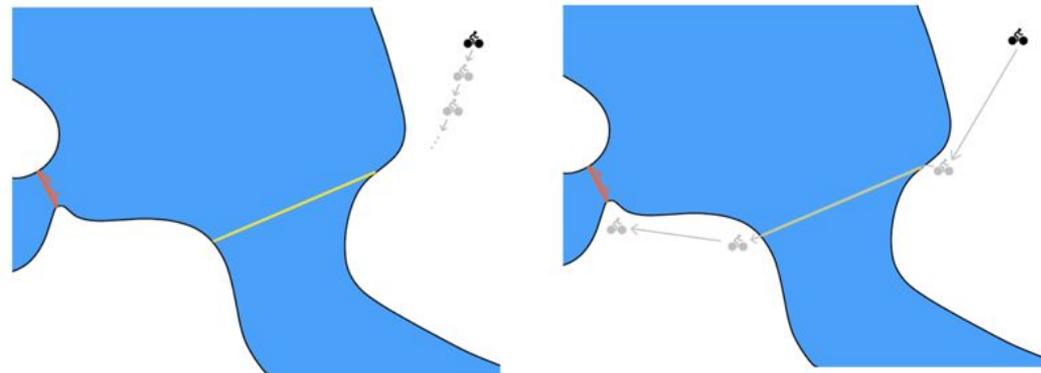
# Direct Collocation with TDMs

[[Pong et al 2018 \(ICLR\)](#)]

- Hierarchical RL with goal-conditioned Q-function + direct collocation

$$a_t = \operatorname{argmax}_{a_{t:K:t+T}, s_{t+K:K:t+T}} \sum_{i=t, t+K, \dots, t+T} r_c(s_i, a_i)$$

such that  $Q(s_i, a_i, s_{i+K}, K-1) = 0 \forall i \in \{t, t+K, \dots, t+T-K\}$



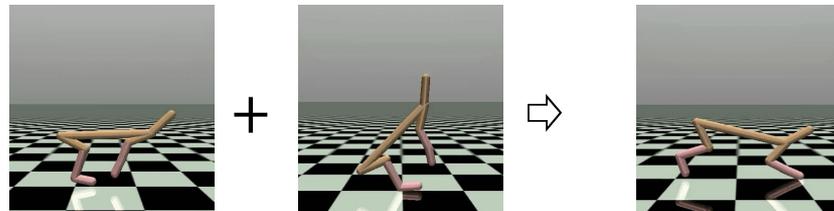
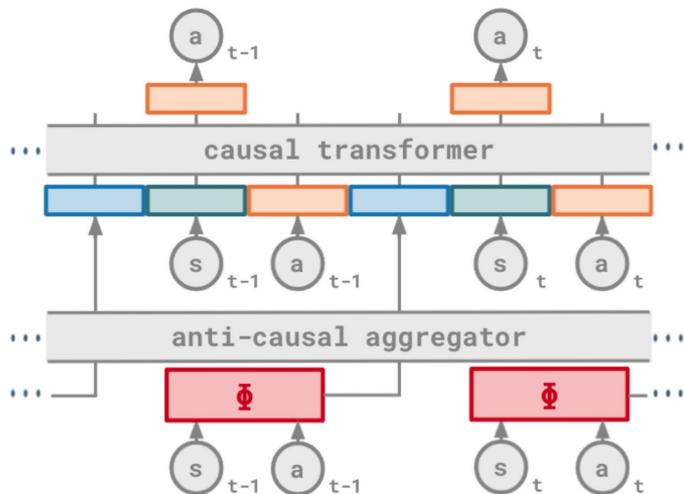
Direction Collocation (Constrained Optimization):

$$\operatorname{arg max}_{s_{t:t+T}} \sum_{i=t}^{t+T} r(s_i) \quad \text{s.t.} \quad -|A| \leq \Pi(s_i, s_{i+1}) \leq |A|$$

$$\operatorname{arg max}_{s_{t:t+T}, a_{t:t+T}} \sum_{i=t}^{t+T} r(s_i) \quad \text{s.t.} \quad Q(s_i, a_i, s_{i+1}) = 0$$

# Generalized Decision Transformers

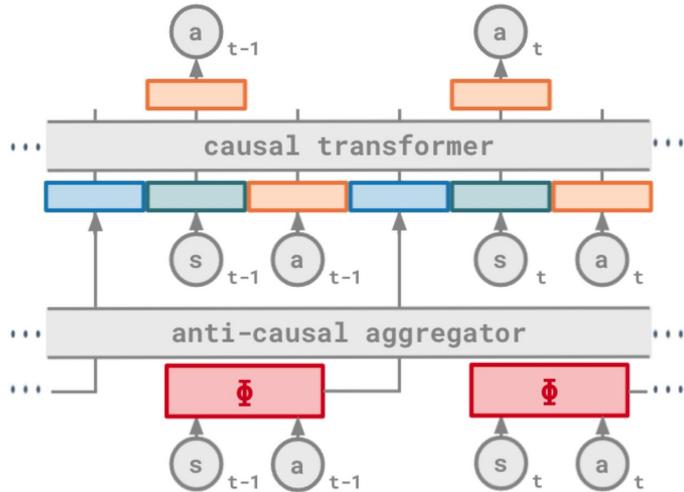
- Training: hindsight BC (behavioral cloning) wrt “any future statistics”
- Test-time: provide unseen “future” and ask it for “generalization”



Method	Algo. Type	Training	$I^{\Phi}(\tau)$	Architectures
Andrychowicz et al. (2017)	RL	Online	$\phi_T$	MLP
Pong et al. (2018)	RL	Online	$\phi_T$	MLP
Chebatar et al. (2021)	RL	Offline	$\phi_T$	CNN
Li et al. (2020)	RL	Online	$\arg \max \sum_t \gamma^t r(s_t, a_t, \cdot)$	MLP
Eysenbach et al. (2020)	BC/RL	On/Offline	$\arg \max \sum_t \gamma^t r(s_t, a_t, \cdot)$	MLP
Lynch et al. (2019)	BC	Offline	$\phi_T$	Stochastic RNN
Ghosh et al. (2021)	BC	Online	$\phi_T$	MLP
Srivastava et al. (2019)	BC	Online	$\sum_t \gamma^t r_t$	Fast Weights
Kumar et al. (2019)	BC	Online	$\sum_t \gamma^t r_t$	MLP
Janner et al. (2021)	BC	Offline	$\sum_t \gamma^t r_t$ or $\phi_T$	Transformer
Duan et al. (2017) <sup>3</sup>	BC	Offline	$\tau$	MLP + LSTM
Generalized DT (ours)	BC	Offline	Any	Transformer
DT (Chen et al., 2021a)	BC	Offline	$\sum_t \gamma^t r_t$	Transformer
Categorical DT (ours) <sup>4</sup>	BC	Offline	histogram( $r_t, \gamma$ )	Transformer
Bi-Directional DT (ours)	BC	Offline	$\tau$	Transformer

# Direct Collocation with GDTs (untested)

- Search for a “reachable future” using the policy function



Shane Gu @shaneguML · Oct 2, 2022

World Model is a causal predictor

Decision Transformer is an anti-causal predictor

Hindsight Experience Replay is the Jedi mind trick to flip the causality

Direction Collocation (Constrained Optimization):

$$\arg \max_{s_{t:t+T}} \sum_{i=t}^{t+T} r(s_i) \quad \text{s.t.} \quad -|A| \leq \Pi(s_i, s_{i+1}) \leq |A|$$

$$\arg \max_{s_{t:t+T}, a_{t:t+T}} \sum_{i=t}^{t+T} r(s_i) \quad \text{s.t.} \quad Q(s_i, a_i, s_{i+1}) = 0$$

Physical and symbolic world models

# 2022: AGI Year 0

- “LLMs can reason” → Symbolic AGI is achievable through LLMs
- “ImagenVideo / DreamFusion” → Physical AGI is achievable through video models

---

## Large Language Models are Zero-Shot Reasoners

---

**Takeshi Kojima**

The University of Tokyo  
t.kojima@weblab.t.u-tokyo.ac.jp

**Shixiang Shane Gu**

Google Research, Brain Team

**Machel Reid**  
Google Research\*

**Yutaka Matsuo**  
The University of Tokyo

**Yusuke Iwasawa**  
The University of Tokyo

### (d) Zero-shot-CoT (Ours)

Q: A juggler can juggle 16 balls. Half of the balls are golf balls, and half of the golf balls are blue. How many blue golf balls are there?

A: **Let's think step by step.**

(Output) *There are 16 balls in total. Half of the balls are golf balls. That means that there are 8 golf balls. Half of the golf balls are blue. That means that there are 4 blue golf balls. ✓*



# 2025: Video model as the missing foundation model

- World consists of symbols, space, and time
- LLMs: reasoning with symbols
- Video models: reasoning in space and time
- Chain-of-Frames (CoF) vs Chain-of-Thoughts (CoT)

Google DeepMind

2025-9-25

## Video models are zero-shot learners and reasoners

Thaddäus Wiedemer<sup>\*1</sup>, Yuxuan Li<sup>1</sup>, Paul Vicol<sup>1</sup>, Shixiang Shane Gu<sup>1</sup>, Nick Matarese<sup>1</sup>, Kevin Swersky<sup>1</sup>, Been Kim<sup>1</sup>, Priyank Jaini<sup>\*1</sup> and Robert Geirhos<sup>\*1</sup>

<sup>1</sup>Google DeepMind

The remarkable zero-shot capabilities of Large Language Models (LLMs) have propelled natural language processing from task-specific models to unified, generalist foundation models. This transformation emerged from simple primitives: large, generative models trained on web-scale data. Curiously, the same primitives apply to today's generative video models. Could video models be on a trajectory towards general-purpose vision understanding, much like LLMs developed general-purpose language

Science & technology | Look at me!

### AI video: more than just “slop”

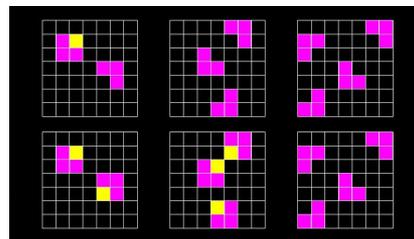
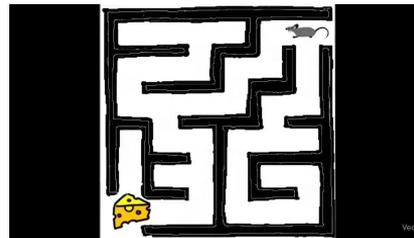
The next big thing in AI may be pictures, not words

Share



PHOTOGRAPH: SORA/OPENAI

Oct 6th 2025 | 4 min read



# 2025: Video model as the missing foundation model

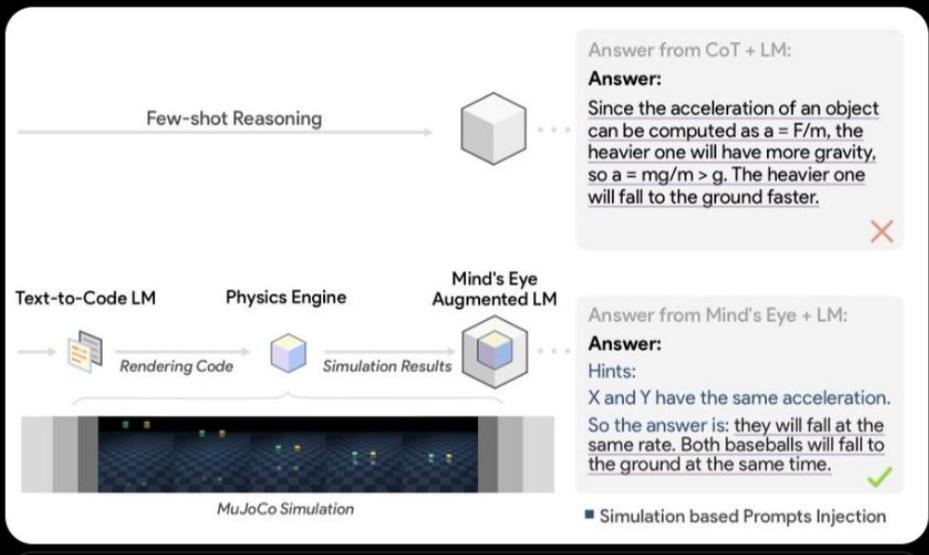


**Shane Gu** ✓ @shaneguML · Oct 12, 2022



Chain-of-thought + **world model** (mental simulator, e.g. mujoco) as a tool = grounded experimenting scientist [arxiv.org/abs/2210.05359](https://arxiv.org/abs/2210.05359)

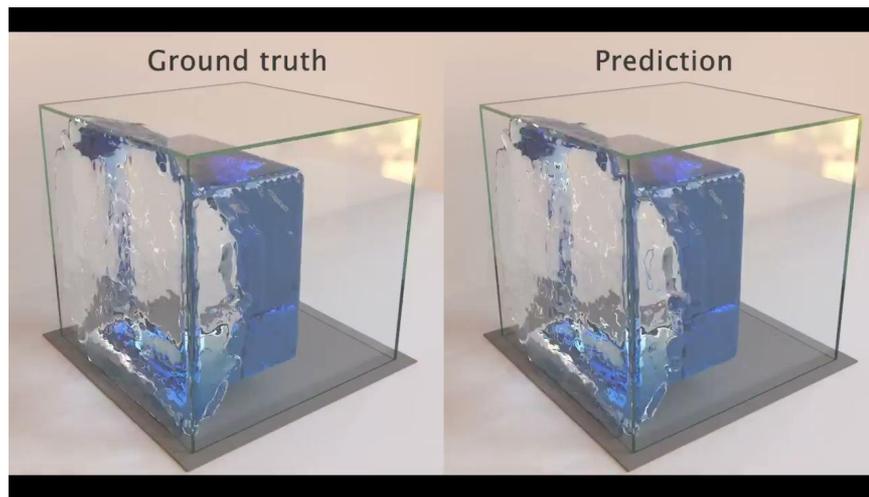
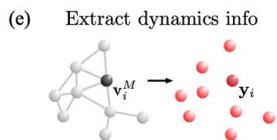
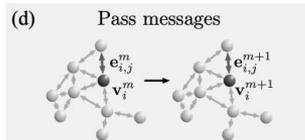
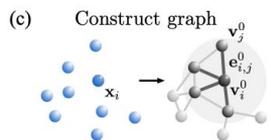
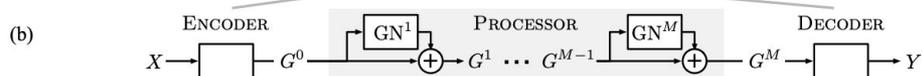
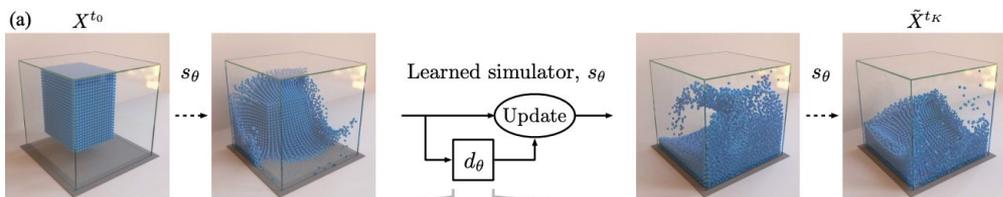
Led by @RuiboLiu! In collaboration w/ @\_jasonwei @denny\_zhou @iamandrewdai et al!



# Older work:

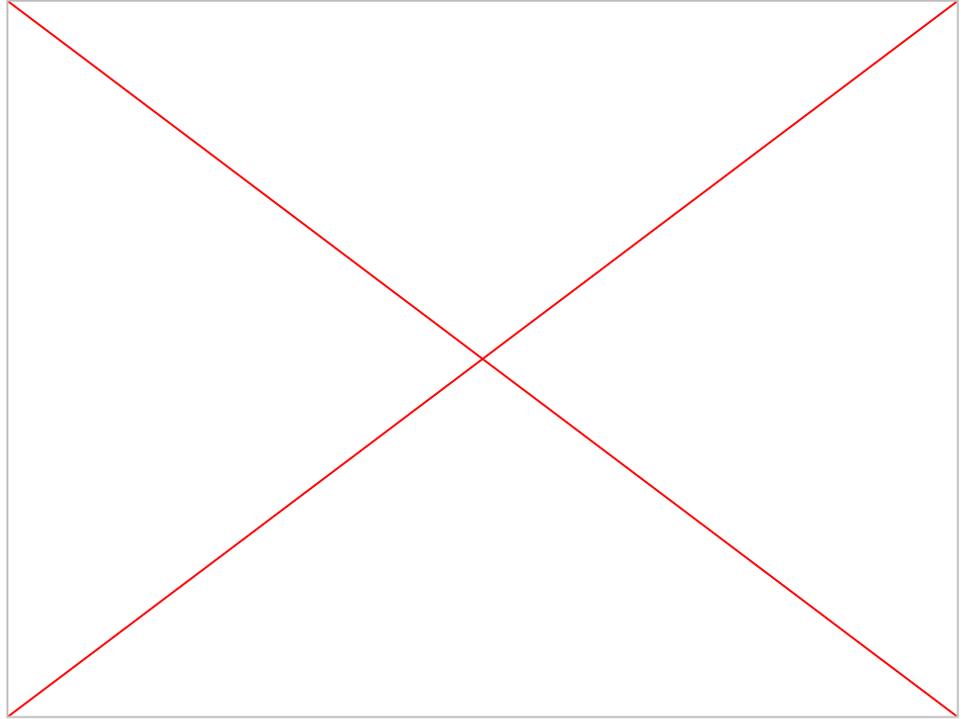
[[Sanchez-Gonzalez et al 2020](#)]

- Particle-based simulation = message passing on graphs
- **We unlikely won't need graph nets / NeRF etc. Just video models.**
  - AlphaFold?



# 2026+: Modeling humanity

Simile: model 8B people



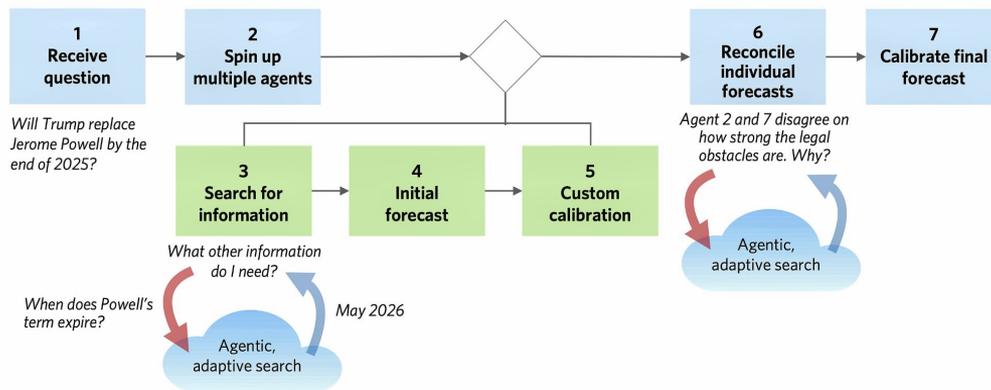
# 2026+: Modeling financial market

License: arXiv.org perpetual non-exclusive license  
arXiv:2511.07678v1 [cs.AI] 10 Nov 2025

## AIA Forecaster: Technical Report

Rohan Alur\*, Bradly C. Stadie\*, Daniel Kang, Ryan Chen, Matt McManus,  
Michael Rickert, Tyler Lee, Michael Federici, Richard Zhu, Dennis Fogerty,  
Hayley Williamson, Nina Lozinski, Aaron Linsky, Jasjeet S. Sekhon

\addrBridgewater AIA Labs  
New York, NY



Thank you! Remember the Level 3.

X (English): shaneguml@

X (Japanese): shanegjp@

Email: shaneguml@gmail.com