# The Objects of Our Curiosity
## Intrinsic Motivation, Intuitive Physics and Self-Supervised Learning

NeurIPS Workshop: Modeling the Physical World
*2018.12.07*

Daniel Yamins

Stanford Neurosciences Institute
Stanford Artificial Intelligence Laboratory
Departments of Psychology and Computer Science
Stanford University

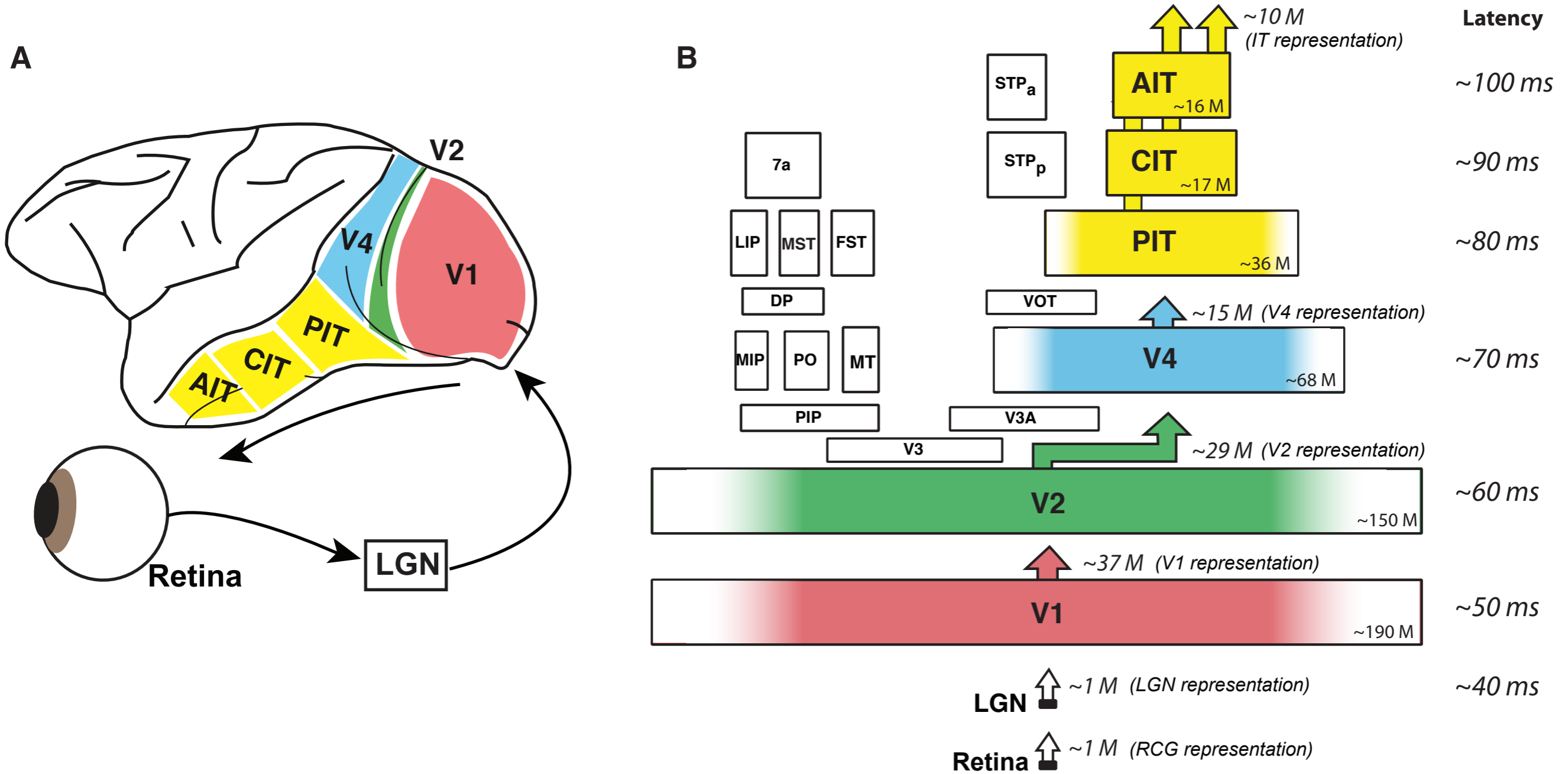# Our work is founded on two mutually reinforcing goals:



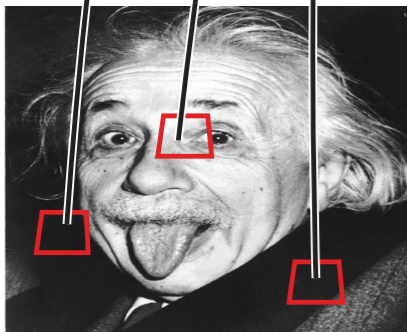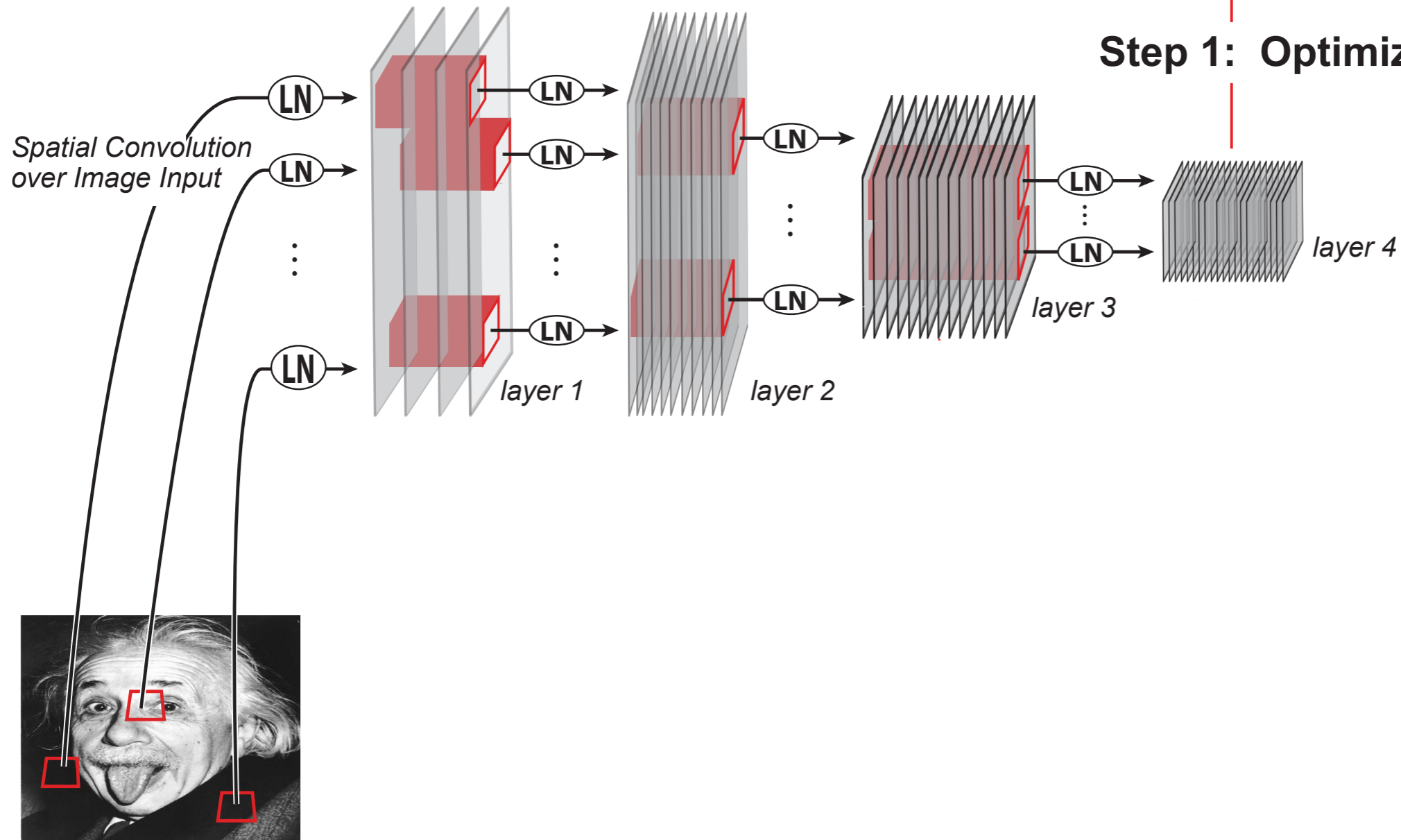*Understanding Brains & Cognition*

*Better AI/ML techniques*

# Computational Models of the Visual System

The primate visual system as a hierarchical, convolutional neural network:



*Adapted from DiCarlo et al. 2012*

# Computational Models of the Visual System



Visual Recognition Task

**Step 1: Optimize for Task**

*Spatial Convolution over Image Input*

LN

layer 1

layer 2

layer 3

layer 4

# Computational Models of the Visual System

Visual Recognition Task

Step 1: Optimize for Task

Step 2: Compare to Neural Data

Spatial Convolution over Image Input

LN

layer 1

layer 2

layer 3

layer 4

100ms Visual Presentation

V1

V2

V4

IT

To our knowledge best (in terms of neural prediction) feedforward model is a ~12-layer CNN

To our knowledge best (in terms of neural prediction) feedforward model is a ~12-layer CNN



...trained on ImageNet Categorization.

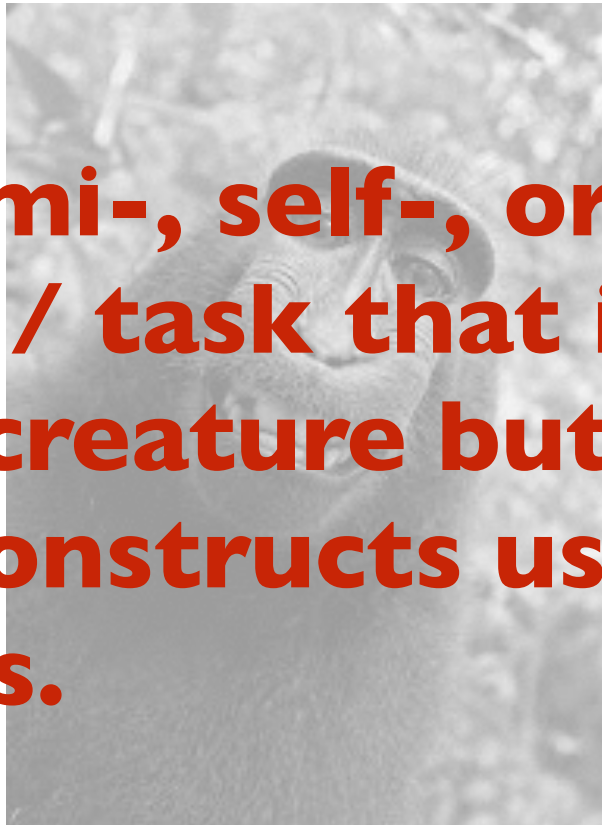There's just no way that these creatures receive millions of high-level semantic labels during learning.



ImageNet is a pretty effective proxy, but just obviously deeply wrong.

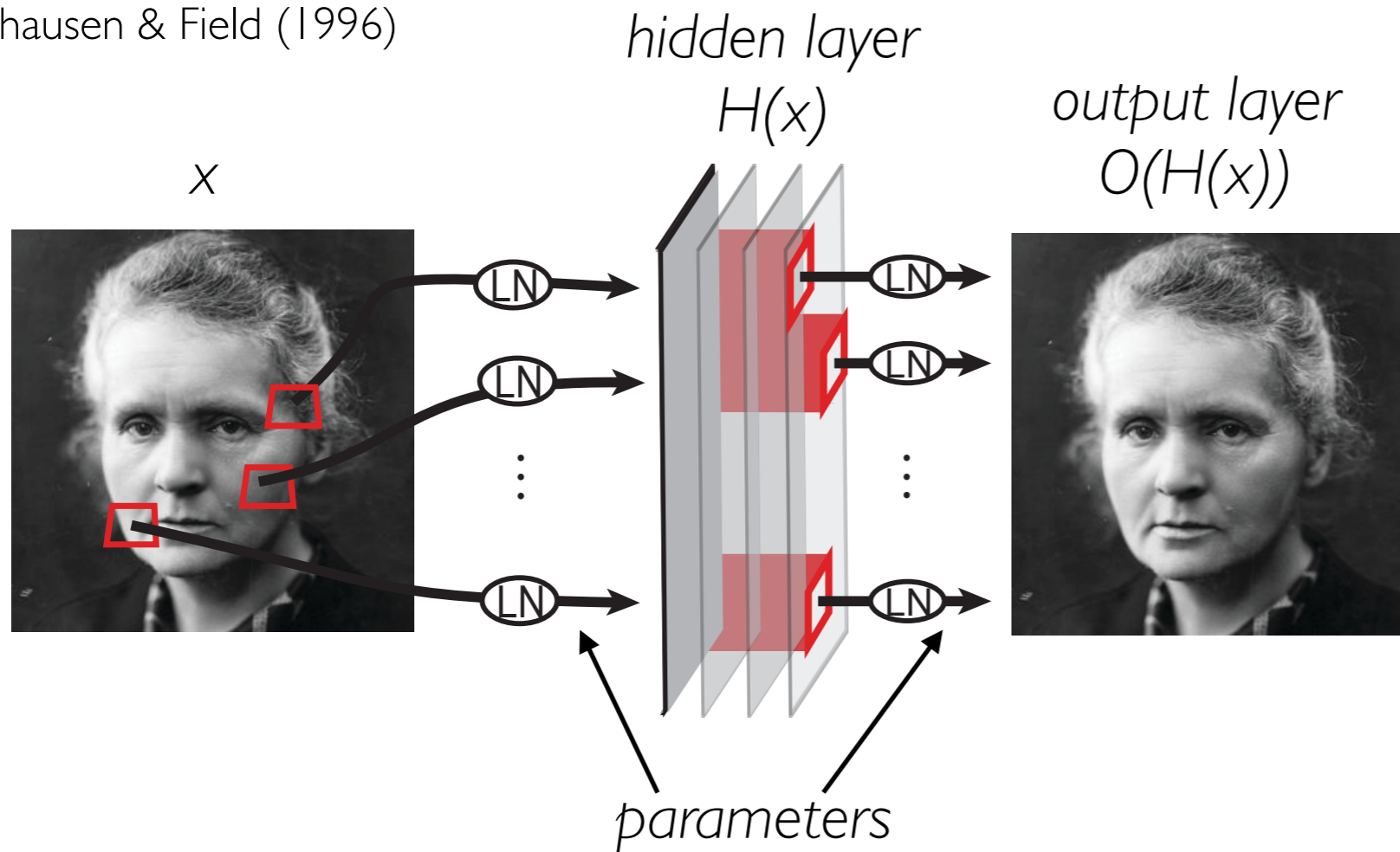There's just no way that these creatures receive millions of high-level semantic labels during learning.

**Must find some sort of semi-, self-, or unsupervised loss function / task that is "realistically costly" to the creature but is sufficiently powerful that it constructs useful representations.**

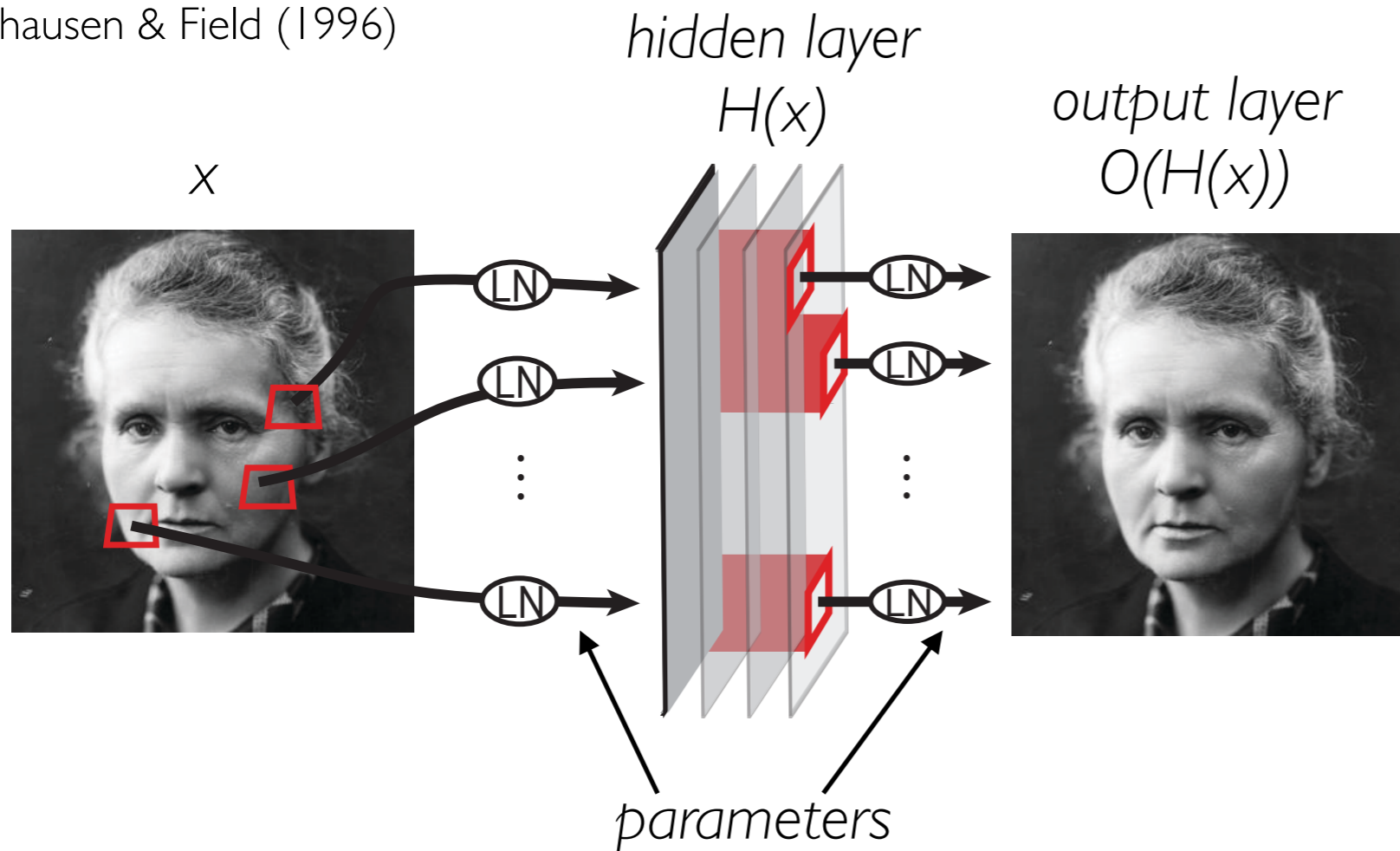ImageNet is a pretty effective proxy, but just obviously deeply wrong.

# Self-supervised learning

Olshausen & Field (1996)



*x*

*hidden layer*
*H(x)*

*output layer*
*O(H(x))*

*parameters*

# Self-supervised learning

Olshausen & Field (1996)



$$L(x) = |x - O(H(x))|^2 + \lambda \cdot |H(x)|$$

reconstruction
loss

complexity
penalty

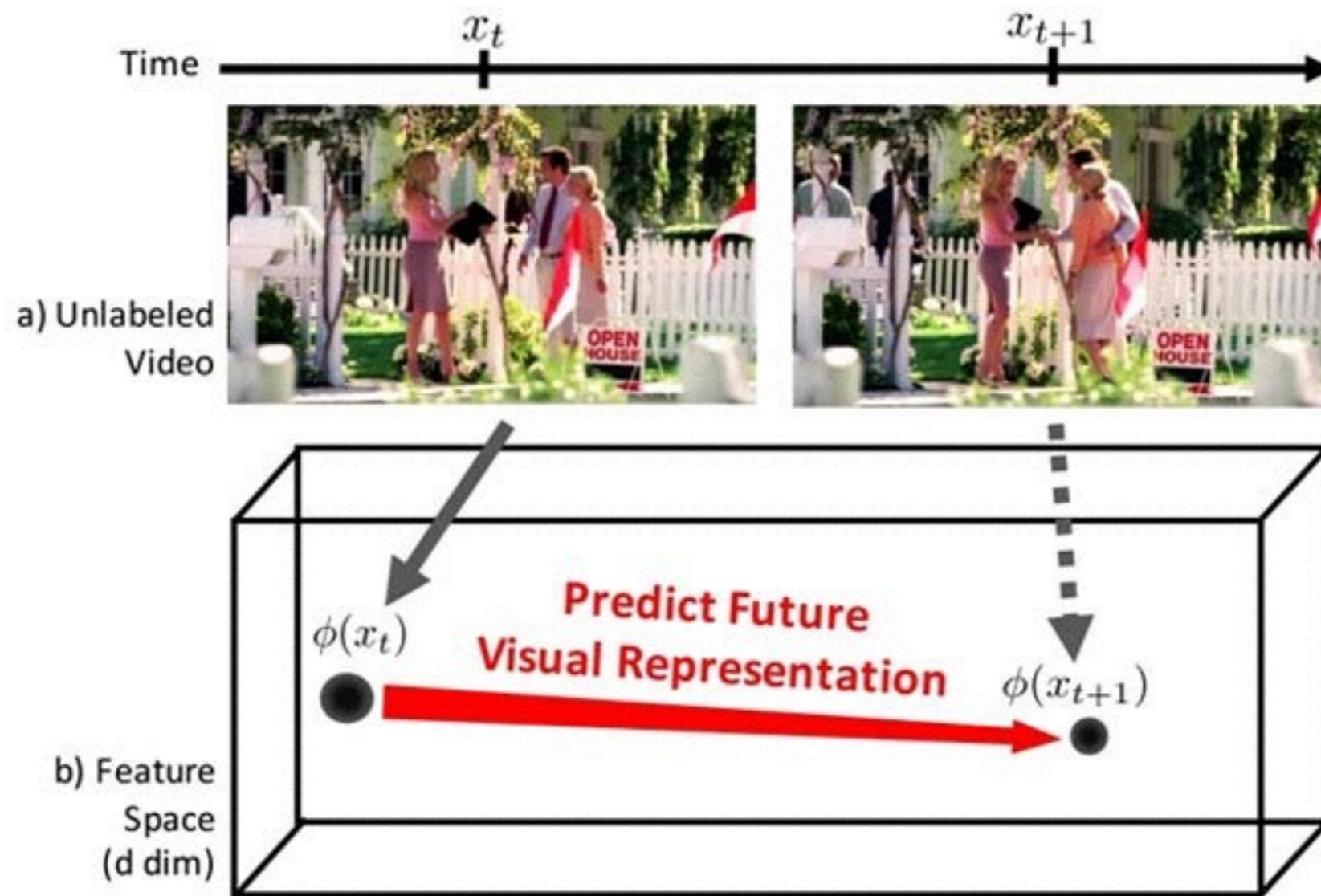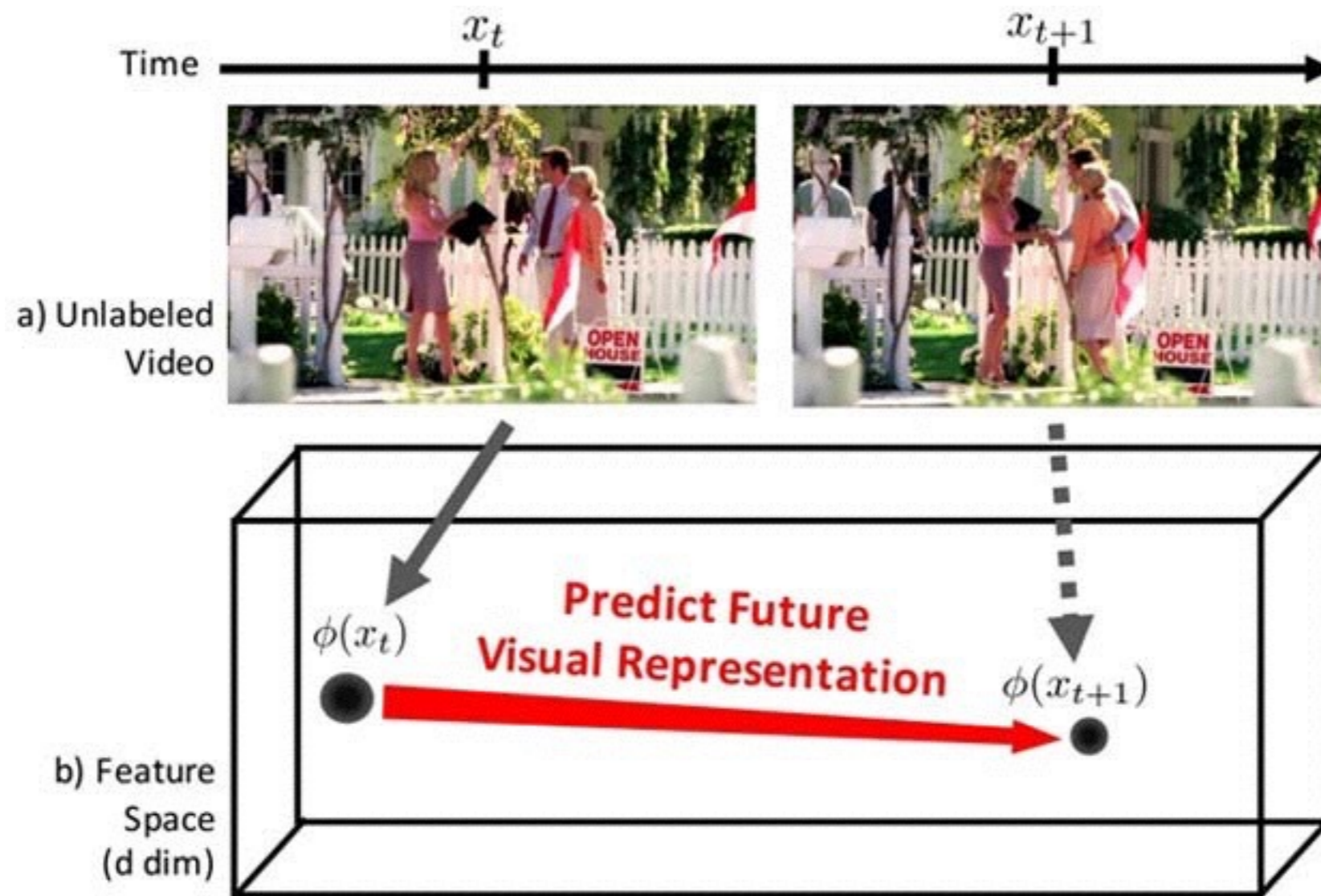# Self-supervised learning

Olshausen & Field (1996)

*x*

*hidden layer*
*H(x)*

*output layer*
*O(H(x))*



LN

LN

LN

LN

LN

LN

*parameters*

*(to some extent)*

$$L(x) = |x - O(H(x))|^2 + \lambda \cdot |H(x)|$$

reconstruction
loss

complexity
penalty

# Self-supervised learning

*Dynamics might give richer signal*

# Self-supervised learning

*Dynamics might give richer signal*



$$L(x) = |x_{t+1} - \text{Decode}(\text{Encode}(x_t))|^2 + \lambda \cdot \text{Penalty}(\text{Encode}(x_t))$$

# Self-supervised learning

*Dynamics might give richer signal*      *…but most passive video sequences are **quite boring***



*fairly trivial features*

$$L(x) = |x_{t+1} - \text{Decode}(\text{Encode}(x_t))|^2 + \lambda \cdot \text{Penalty}(\text{Encode}(x_t))$$

Children learn through **play.**

How does this work?

# Self-supervised learning

*Give agent some kind of volition to take actions*



$$L(x) = |x_{t+1}^{\text{action}} - \text{Decode}(\text{Encode}(x_t))|^2 + \lambda \cdot \text{Penalty}(\text{Encode}(x_t))$$

# Self-supervised learning

*Give agent some kind of volition to take actions ... but now the agent will be **lazy***



$$L(x) = |x_{t+1}^{\text{action}} - \text{Decode}(\text{Encode}(x_t))|^2 + \lambda \cdot \text{Penalty}(\text{Encode}(x_t))$$

# Self-supervised learning

*Give agent some kind of volition to take actions … but now the agent will be **lazy***



$$L(x) = |x^{\text{action}}{}_{t+1} - \text{Decode}(\text{Encode}(x_t))|^2 + \lambda \cdot \text{Penalty}(\text{Encode}(x_t))$$
$$+ \text{Intrinsic Motivation}$$

*Environment*

Environment

Agent

Perception

*Environment*

*Agent*

*Perception*

Environment

Agent

Perception

Action

*Environment*

*Agent*

**Perception**

World

Model

Environment

**Perception**

Agent
*learns to predict the environment*

World
Model

*Action*

Self
Model

*Environment*

**Perception**

**Action**

*Agent*

World
Model

Self
Model

*learns to predict
the environment*

*learns to predict
own world-model*

# A curiosity principle:

The **self-model** directs the agent toward **interesting** actions — the ones that the **world-model** doesn't yet fully understand

# A curiosity principle:

The **self-model** directs the agent toward **interesting** actions — the ones that the **world-model** doesn't yet fully understand

# A curiosity principle:

The **self-model** directs the agent toward ***interesting*** actions — the ones that the **world-model** doesn't yet fully understand



Environment

Agent

**Perception**

World Model

*learning*

*learns to predict the environment*

**Action**

Self Model

*curiosity*

*learns to predict own world-model*

# A curiosity principle:

The **self-model** directs the agent toward **interesting** actions — the ones that the **world-model** doesn't yet fully understand

Nick Haber    Damian Mrowca    Stephanie Wang    Fei-Fei Li

NIPS 2018

# Learning to Play With Intrinsically-Motivated, Self-Aware Agents

Nick Haber[1,2,3,*], Damian Mrowca[4,*], Stephanie Wang[4], Li Fei-Fei[4], and
Daniel L. K. Yamins[1,4,5]

Departments of Psychology[1], Pediatrics[2], Biomedical Data Science[3], Computer Science[4], and Wu
Tsai Neurosciences Institute[5], Stanford, CA 94305

{nhaber, mrowca}@stanford.edu

## Abstract

Infants are experts at playing, with an amazing ability to generate novel structured behaviors in unstructured environments that lack clear extrinsic reward signals. We seek to mathematically formalize these abilities using a neural network that implements curiosity-driven intrinsic motivation. Using a simple but ecologically naturalistic simulated environment in which an agent can move and interact with objects it sees, we propose a "world-model" network that learns to predict the dynamic consequences of the agent's actions. Simultaneously, we train a separate explicit "self-model" that allows the agent to track the error map of its world-model. It then uses the self-model to adversarially challenge the developing world-model. We demonstrate that this policy causes the agent to explore novel and informative interactions with its environment, leading to the generation of a spectrum of complex behaviors, including ego-motion prediction, object attention, and object gathering. Moreover, the world-model that the agent learns supports improved performance on object dynamics prediction, detection, localization and recognition tasks. Taken together, our results are initial steps toward creating flexible autonomous agents that self-supervise in realistic physical environments.

Agent ("baby") can (a) swivel its head

(b) move around the room

(c) apply forces to objects

…and receives back images of what happened, given action

Model has two pieces: **(1) World-Model**

$$\Lambda_{\psi_t}$$

$$\omega_{\theta_t}$$

Model has two pieces:   **(1) World-Model**

**(2) Self-Model**

*Goal of **world-model**:*
*"Post-dict" the action taken given past and future states and actions*

*Goal of **world-model**:*
*"Post-dict" the action taken given past and future states and actions*



## ID
*Inverse Dynamics*

*Goal of **self-model**: Predict errors ("loss") of World-Model*

*Goal of **self-model**: Predict errors ("loss") of World-Model*

# Learning to Play - Model overview



Goal of **world-model** network is to predict consequences of actions

Goal of **self-model** network is to predict errors of world-model (''self-aware'')

*Goal of **self-model**: Predict errors ("loss") of World-Model*



*Sample 1000x actions and choose the one that maximizes the World-Model loss*

$$\pi(a) \sim \exp(\beta \sigma_\Lambda(a))$$

**Policy mechanism**

# Learning to Play — Adversarial Policy

Action choice: self-model is **adversarial** to world-model ("curious intrinsic motivation")



Goal of **world-model** network is to predict consequences of actions

Goal of **self-model** network is to predict errors of world-model ("self-aware")

Place agent in room with a single object.

# Self-supervised learning



curious policy          random policy

At first, agent totally ignores the object, and focuses on learning **ego-motion**…

# Self-supervised learning



But the curious agent eventually "gets bored" of ego-motion prediction and starts to focus on the object!

# Self-supervised learning



curious policy     random policy

Training Loss

Ego motion learning   *    Emergence of object attention   o    Object interaction learning

Training Steps

Testing Loss

Ego-motion only (easy)

Object Interactions (hard)

Fraction of Frames

Object Presence

Emergence of:
(a) ego-motion understanding

(b) object attention

(c) improved world-model

Simple navigation and planning behavior emerges …



If an object is not in view, the agent turns to find one…

# Self-supervised learning

Simple navigation and planning behavior emerges …



If an object is not in view, the agent turns to find one…

… if an object is too far to touch, the agent moves toward one.

Simple navigation and planning behavior emerges …



If an object is not in view, the agent turns to find one…

… if an object is too far to touch, the agent moves toward one.

… and once the agent is close to an object, it stays close and interacts with it.

# Self-supervised learning

Simple navigation and planning behavior emerges …



Moreover, substantially improved transfer learning accuracy:

(a) object detection (present or not): ~8% vs ~40% accuracy

(b) object position: ~6px vs ~4px error

(c) object recognition (among 16 geometries): ~12% vs ~30% accuracy

# Self-supervised learning



When multiple objects are present, the agent at first recapitulates its behavior with a single object …

# Self-supervised learning



… but then discovers the interest of bringing objects together.

# Self-supervised learning



Object recognition in testing (one object per image):  ~16%  vs ~40% accuracy

… especially large gain compared to training in single-obj case

Glossing over a key problem:

The above ideas rely on having the agent solve a dynamics prediction problem about the world.

*Start with some data H…*

$$H$$

*…Create input data X from H…*

$$H$$

$$\xi$$

$$X$$

*…Create output data Y from H…*

$$H$$

$$\xi \qquad \eta$$

$$X \qquad Y$$

*…Predict Y from X.*

$$H$$

$$\xi_t \qquad \eta_t$$

$$X \cdots\cdots\cdots\cdots\rightarrow Y$$

$$\omega_t$$

*Examples:*

*1) Forward future prediction*

*state(t), action(t) $\Longrightarrow$ state(t+1)*

$$H$$

$$\xi_t \qquad \eta_t$$

$$X \cdots\cdots\cdots\cdots\cdot\rightarrow Y$$

$$\omega_t$$

*Examples:*

*1) Forward future prediction*

*state(t), action(t) $\Longrightarrow$ state(t+1)*

*2) inverse dynamics prediction*

*state(t), state(t+1) $\Longrightarrow$ action(t)*

$$H$$
$$\xi_t \quad \eta_t$$
$$X \quad \cdots\cdots\!\!\!\!\!\longrightarrow Y$$
$$\omega_t$$

*Examples:*

*1) Forward future prediction*

*state(t), action(t) $\implies$ state(t+1)*

*2) inverse dynamics prediction*

*state(t), state(t+1) $\implies$ action(t)*

*1) is hard, because … pixel prediction is hard!*

Obvious idea: just predict future pixels

Finn et. al (2016)



PredRNN(2017) ; Wang (2018) ; among many others

Pixel prediction is hard.

t=1      t = 2



Two blue objects in a room

# Intuitive Physics as Underlying Goal

Pixel prediction is hard.

t=1     t = 2



Two blue objects in a room

Objects acted on and camera moves

Pixel prediction is hard.

t=1    t = 2



Two blue objects in a room

Objects acted on and camera moves

t=3    t = 4    t=5    t=6

Ground
truth

# Intuitive Physics as Underlying Goal

Pixel prediction is hard.

t=1        t = 2        Two blue objects in a room

Objects acted on and camera moves

t=3        t = 4        t=5        t=6

Ground
truth

Prediction

$$H$$

$$\xi_t \qquad \eta_t$$

$$X \cdots\cdots\cdots\rightarrow Y$$
$$\omega_t$$

*Examples:*

*1) Forward future prediction*

*state(t), action(t)* $\implies$ *state(t+1)*

*2) inverse dynamics prediction*

*state(t), state(t+1)* $\implies$ *action(t)*

*1) is hard, because ... pixel prediction is hard!*

*2) is mostly what we did in the work described above because it's easier ...*

$$H$$

$$\xi_t \qquad \eta_t$$

$$X \cdots\cdots\cdots\cdots\rightarrow Y$$

$$\omega_t$$

*Examples:*

*1) Forward future prediction*

*state(t), action(t) $\Longrightarrow$ state(t+1)*

*2) inverse dynamics prediction*

*state(t), state(t+1) $\Longrightarrow$ action(t)*

*1) is hard, because … pixel prediction is hard!*

*2) is mostly what we did in the work described above because it's easier …*

*BUT DEGENERATE!*

$H$

$\xi_t$ $\eta_t$

$X$ $\cdots\cdots\cdots\rightarrow$ $Y$

$\omega_t$

*possibly ill-defined*

THE DREADED
**WHITE-NOISE PROBLEM**

*Examples:*

*1) Forward future prediction*

*state(t), action(t)* $\implies$ *state(t+1)*

*2) inverse dynamics prediction*

*state(t), state(t+1)* $\implies$ *action(t)*

*1) is hard, because … pixel prediction is hard!*

*2) is mostly what we did in the work described above because it's easier …*

*BUT DEGENERATE!*

$H$

$\xi_t$        $\eta_t$

$X$ $\cdots\cdots\cdots\cdots\cdots\rightarrow$ $Y$

$\boxed{\omega_t}$

*possibly ill-defined*

*THE DREADED*

**WHITE-NOISE PROBLEM**

*Examples:*

*1) Forward future prediction*

*state(t), action(t)* $\implies$ *state(t+1)*

*2) inverse dynamics prediction*

*state(t), state(t+1)* $\implies$ *action(t)*

*1) is hard, because … pixel prediction is hard!*

*2) is mostly what we did in the previous work because it's easier …*

*BUT DEGENERATE!*

*Ex: pushing down on an object*

$$H$$

$$\xi_t \qquad \eta_t$$

$$X \cdots\cdots\cdots\cdots\cdots Y$$

$$\boxed{\omega_t}$$

*THE DREADED*

**WHITE-NOISE PROBLEM**

*Examples:*

*1) Forward future prediction*

*state(t), action(t) $\Longrightarrow$ state(t+1)*

*...prediction*

*$\Longrightarrow$ action(t)*

> **Conclusion:**
> we cannot escape having to do better
> future prediction
> — so let's attack the problem directly.

*1) is hard, because … pixel prediction is hard!*

*2) is mostly what we did in the previous work because it's easier …*

*BUT DEGENERATE!*

*Ex: pushing down on an object*

**Damian Mrowca\*** **Chengxu Zhuang\*** Eli Wang Nick Haber Fei-Fei Li Josh Tenenbaum

# Flexible Neural Representation for Physics Prediction

Damian Mrowca[1,*], Chengxu Zhuang[2,*], Elias Wang[3,*], Nick Haber[2,4,5], Li Fei-Fei[1], Joshua B. Tenenbaum[7,8], and Daniel L. K. Yamins[1,2,6]

Department of Computer Science[1], Psychology[2], Electrical Engineering[3], Pediatrics[4] and Biomedical Data Science[5], and Wu Tsai Neurosciences Institute[6], Stanford, CA 94305
Department of Brain and Cognitive Sciences[7], and Computer Science and Artificial Intelligence Laboratory[8], MIT, Cambridge, MA 02139

{mrowca, chengxuz, eliwang}@stanford.edu

t-n ...

-1 ... t+k

## Abstract

Humans have a remarkable capacity to understand the physical dynamics of objects in their environment, flexibly capturing complex structures and interactions at multiple levels of detail. Inspired by this ability, we propose a hierarchical particle-based object representation that covers a wide variety of types of three-dimensional objects, including both arbitrary rigid geometrical shapes and deformable materials. We then describe the Hierarchical Relation Network (HRN), an end-to-end differentiable neural network based on hierarchical graph convolution, that learns to predict physical dynamics in this representation. Compared to other neural network baselines, the HRN accurately handles complex collisions and nonrigid deformations, generating plausible dynamics predictions at long time scales in novel settings, and scaling to large scene configurations. These results demonstrate an architecture with the potential to form the basis of next-generation physics predictors for use in computer vision, robotics, and quantitative cognitive science.

# Discovering the proper latent space for physical prediction…

*"Encoding"*         *"Physics"*         *"Rendering"*



**Encoder**    **t → t+1**    **Decoder**

Liz Spelke

Experimental results with infants: **object permanence** present very early, perhaps by 3 months.

Liz Spelke

Experimental results with infants: **object permanence** present very early …

Cognition

Volume 20, Issue 3, 1985, Pages 191-208

ELSEVIER

COGNITION

Object permanence in five-month-old infants ☆

Renée Baillargeon ☺*, Elizabeth S. Spelke *, Stanley Wasserman *

⊞ Show more

https://doi.org/10.1016/0010-0277(85)90008-3

Get rights and content

# Intuitive Physics as Underlying Goal



Liz Spelke

Experimental results with infants: **object permanence** present very early …

Conv2d structures, even with RNNs, have trouble with object permanence.

Liz Spelke

Experimental results with infants: **object permanence** present very early, perhaps by 3 months.

Conv3d structures are better for object permanence, but very inefficient: hard to achieve high resolution.

# Spatial convolutions are not ideal for physics propagation



*"Derendering"*      *"Physics"*      *"Rendering"*

$S_t$     $S_{t+1}$

$F_{2D}$

$P_{2D-INV}$

*Mass transport problem*

$P_{2D}$

$3D{\rightarrow}2D$     $2D{\rightarrow}3D$

$I_{t-n \dots t}$     $P_{3D-INV}$     $P_{3D}$     $I_{t+1}$

*Back completion problem*

$F_{3D}$

$S_t$     78     $S_{t+1}$

Experimental results with infants: **object permanence** present very early, perhaps by 3 months.

Alternative to spatially-uniform priors are **graph-based** priors

Liz Spelke



Relational Networks
(Battaglia et. al., 2016)

Liz Spelke

Experimental results with infants: **object permanence** present very early, perhaps by 3 months.

Alternative to spatially-uniform priors are **graph-based** priors

… still local and convolutional, just on the graph.



Final CNN feature maps

RN

object

Conv.

Object pair with question $g_\theta$-MLP

$f_\phi$-MLP

small

Element-wise sum

What size is the cylinder that is left of the brown metal thing that is left of the big sphere?

what size is … sphere

LSTM

Relational Networks (Battaglia et. al., 2016)

Relational Networks
(Battaglia et. al., 2016)

Neural Physics Engine
(Chang et. al., 2016)

(a)

$o_1$

neighborhood

$o_2$

$o_3$

$o_4$

$t - 1$

$t$

$v_1$

$v_3$ $v_2$

$v_4$

$t + 1$

(b)  NPE applied on object 3

$o_3$

$o_2$

$e_{2,3}$ $o_3$

$o_3$

$o_4$

$e_{4,3}$

$\Sigma$

$v_3$

(c)  NP applied on object 3

$o_3$

$o_2$

$o_4$

$\Sigma$

$v_3$

(d)  LSTM applied on object 3

$o_4$

0

$o_2$

0

$o_3$

1

$v_3$

# Intuitive Physics as Underlying Goal

Complex Scenes

Complex Materials

# Describe objects through complex graphs:



| Cube | Cuboid | Pyramid | Flat Pyramid |
| --- | --- | --- | --- |
| Octahedron | Prism | Cylinder | Ellipsoid |
| Sphere | Mentos | Stick | Bowl |
| Cone | Pentagon | Domino | Torus |
| Duck | | Bunny | Teddy |

In fact, describe whole scenes.



Plane

Bowl

Random Plane

Stairs

Slope

Half-Pipe

In fact, describe whole scenes.



Plane          Bowl          Random Plane

Stairs          Slope          Half-Pipe

$$G = \langle N, E \rangle$$ scene graph

In fact, describe whole scenes.



Plane

Bowl

Random Plane

Stairs

Slope

Half-Pipe

$$G = \langle N, E \rangle \quad \text{scene graph}$$

N = nodes corresponding to particles comprising objects

E = edges corresponding to relationships between particles

In fact, describe whole scenes.



Plane

Bowl

Random Plane

Stairs

Slope

Half-Pipe

$$G = \langle N, E \rangle \quad \text{scene graph}$$

N = nodes corresponding to particles comprising objects

E = edges corresponding to relationships between particles

edges are labelled by vector capturing bond characteristics

Of course, humans don't think about all the particles at once all the time.

Of course, humans don't think about all the particles at once all the time.



$$G \longmapsto G_H$$

$G_H$ = dynamic "hierarchicalization" of underlying scene graph

*(right now computed via k-means)*

Of course, humans don't think about all the particles at once all the time.



$$G \longmapsto G_H$$

$G_H$ = dynamic "hierarchicalization" of underlying scene graph

*(right now computed via k-means)*

graph convolution ➝ hierarchical graph convolution

$\phi^{L2A}$     graph conv. leaves to ancestors

$\phi^{WS}$     graph conv. with siblings

$\phi^{A2D}$     graph conv. ancestors to descendants

# Intuitive Physics as Underlying Goal



$\phi^{L2A}$   graph conv. leaves to ancestors

$\phi^{WS}$   graph conv. with siblings

$\phi^{A2D}$   graph conv. ancestors to descendants

$\eta$   module composing these three operations from one up-down cycle, adding physical effects

# Intuitive Physics as Underlying Goal



$\phi^{L2A}$    graph conv. leaves to ancestors

$\phi^{WS}$    graph conv. with siblings

$\phi^{A2D}$    graph conv. ancestors to descendants

$\eta$    module composing these three operations from one up-down cycle, adding physical effects

Hierarchical graph convolution propagates interactions efficiently

Hierarchical Relational Network (HRN):



…generates momentum updates (**P**) from hierarchical graph state (**G**).

Hierarchical Relational Network (HRN):



…generates momentum updates (**P**) from hierarchical graph state (**G**).

Network learns to interpret graph structure (including meaning of material-vector edge labels)…

Deformable cone bouncing off a flat floor

# Deformable cone bouncing off a flat floor



# Stanford bunny

# Deformable box bouncing off an incline

Deformable box bouncing off an incline

Ground Truth

Prediction

Multiple rigid objects colliding

Ground Truth

Prediction

# rigid sphere rolling out of rigid bowl

# rigid sphere rolling out of rigid bowl



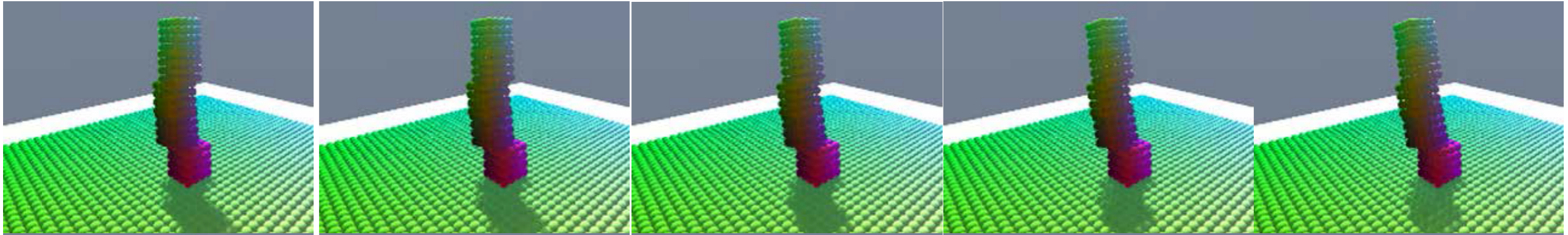# floppy teddybear bouncing off floor and recovering
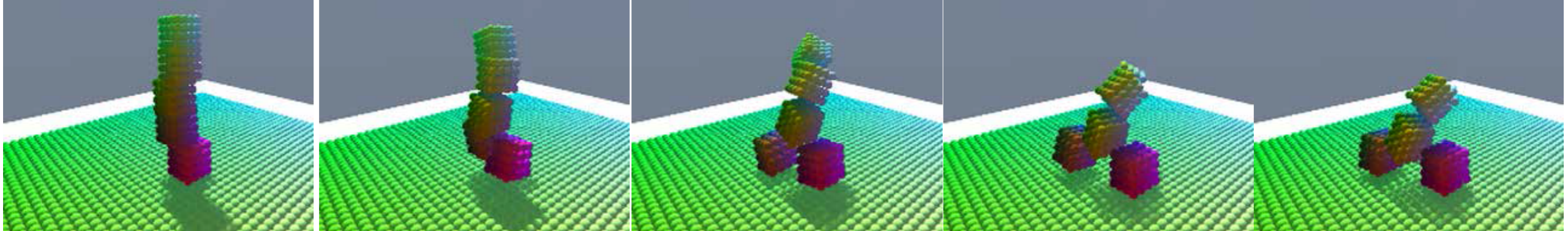
# knocking over an unstable block tower



*in GT the tower does fall, but prediction falls too fast . . .*

# knocking over an unstable block tower



**Ground truth**

**Prediction**
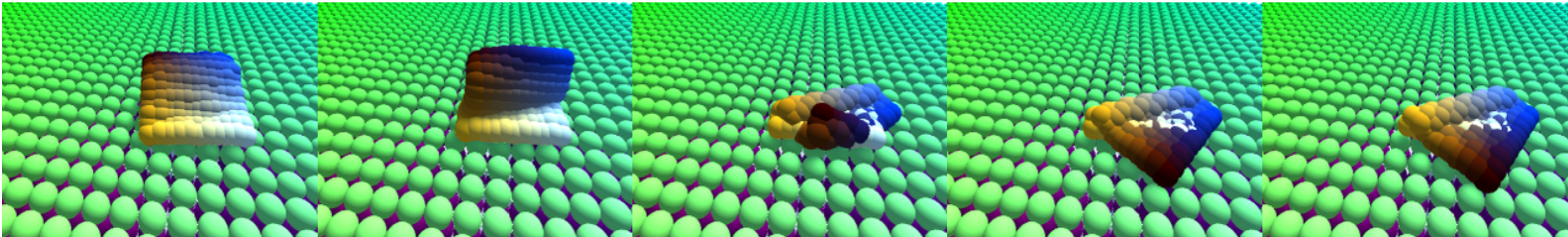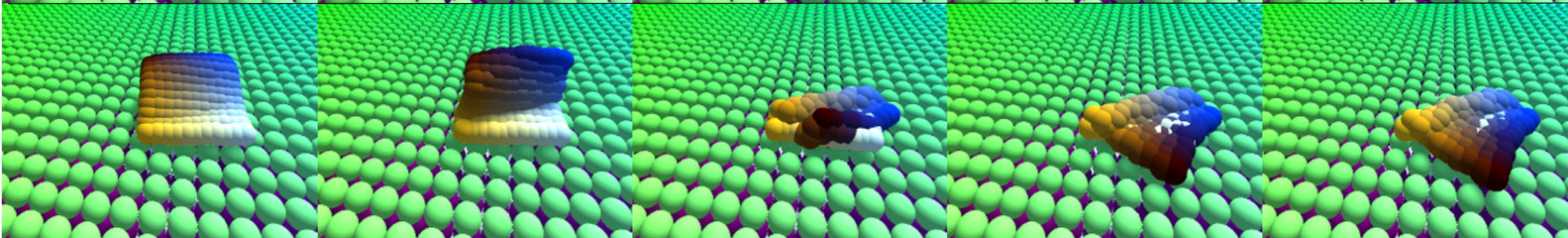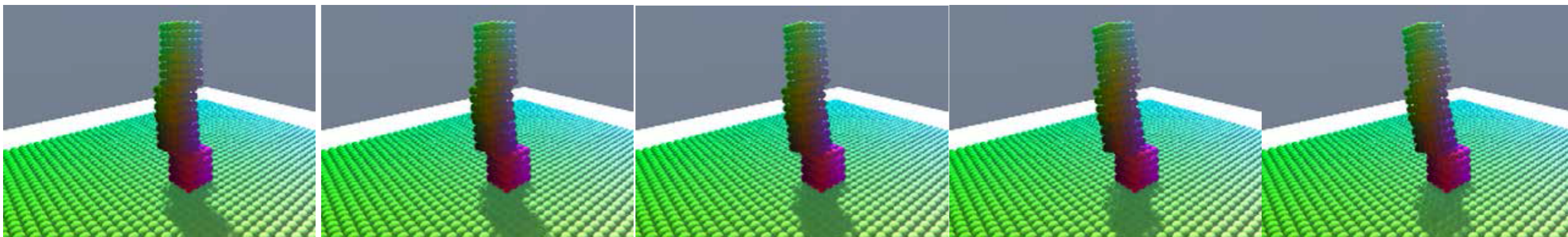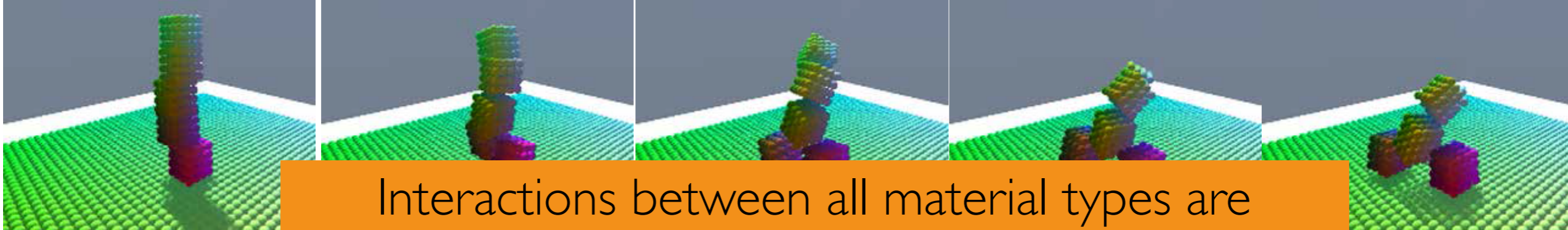
# Folding Cloth



**Ground Truth**

**Prediction**

t+1   t+3   t+5   t+7   t+9

# knocking over an unstable block tower
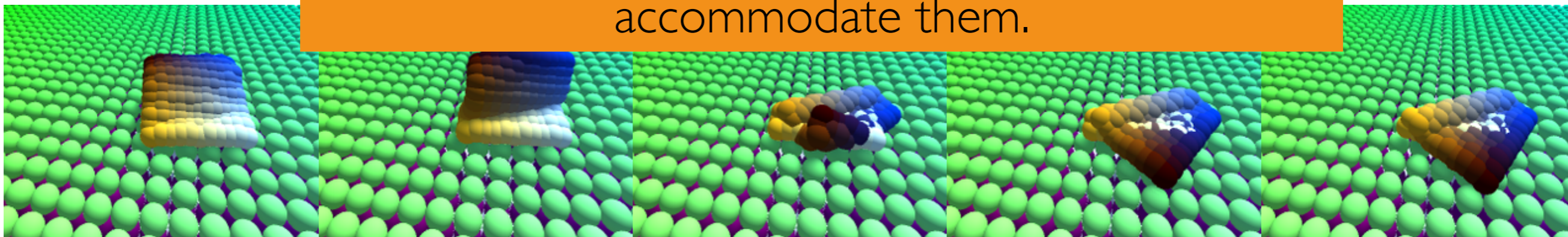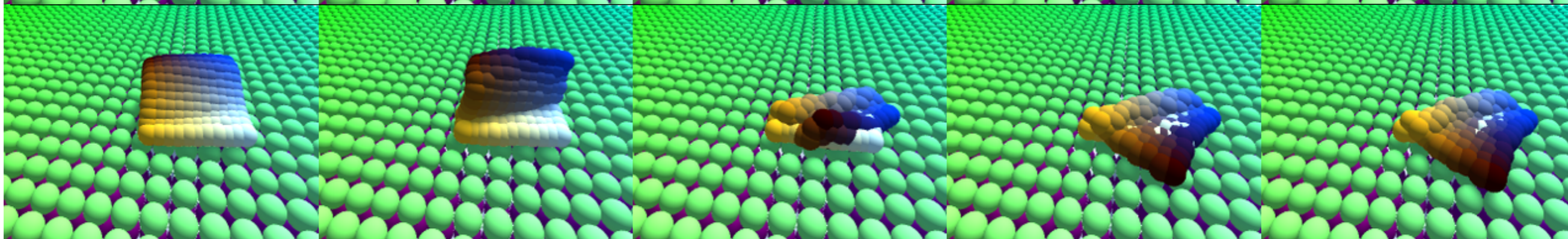


Interactions between all material types are possible — as well as non-uniform materials — since the edge-labelled graph structure can accommodate them.
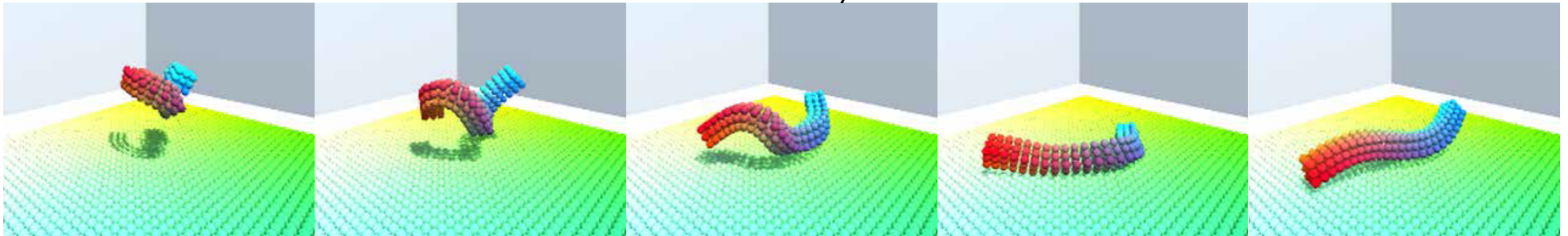
Ground truth

Prediction

Ground Truth

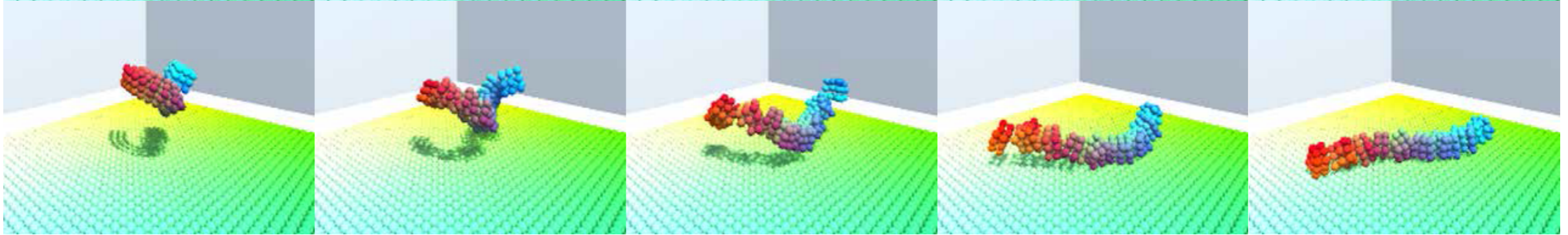Prediction

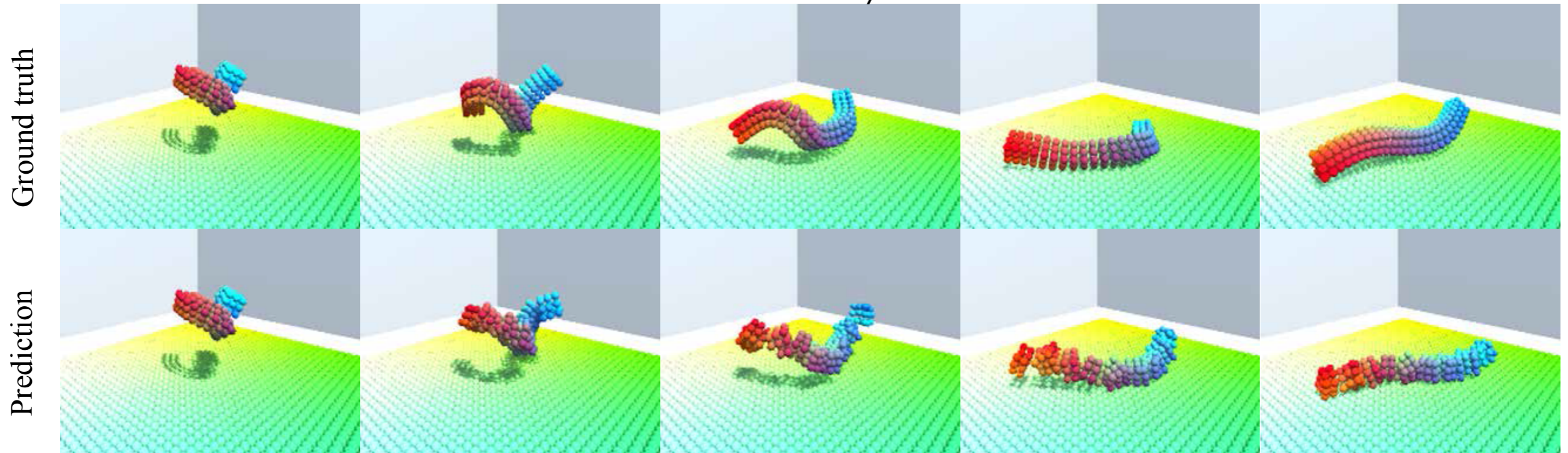t+1          t+3          t+5          t+7          t+9

# Challenges:



slinky

*shape is not preserved super well over long rollouts…*

Challenges:

slinky



*shape is not preserved super well over long rollouts…*

Easy to impose simple shape conversation rules — in a "per material" way. (e.g. rigid different than cloth different than soft-body)

# Challenges:
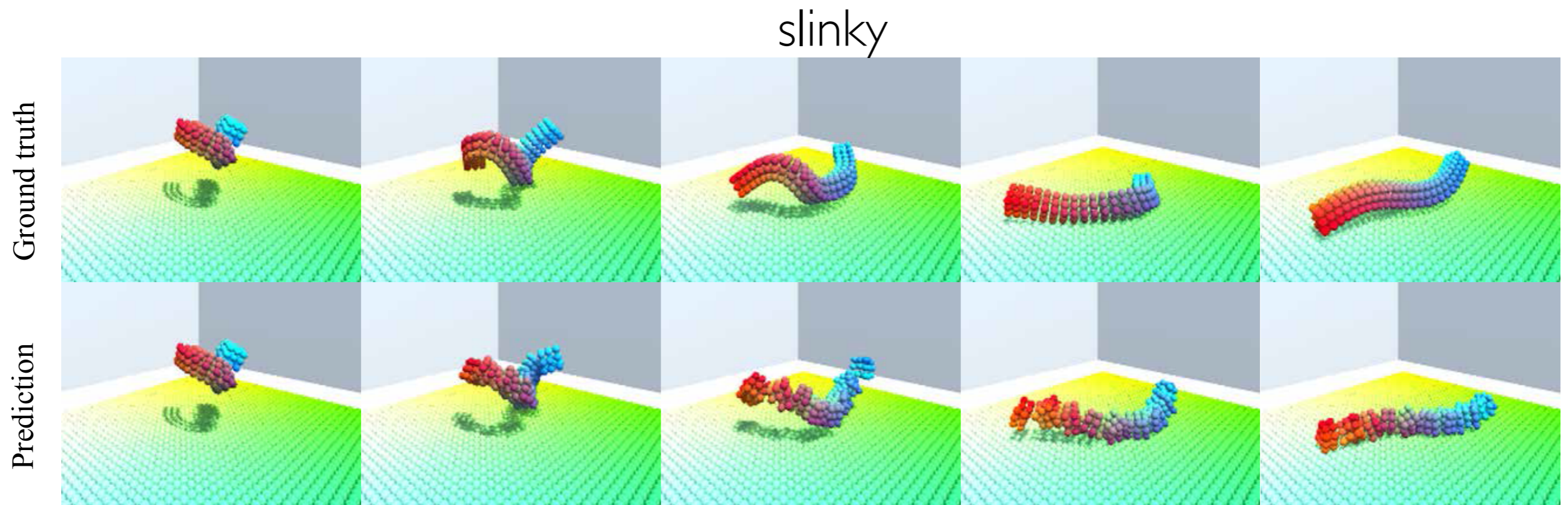
slinky



*shape is not preserved super well over long rollouts…*

Easy to impose simple shape conversation rules — in a "per material" way.  (e.g. rigid different than cloth different than soft-body)

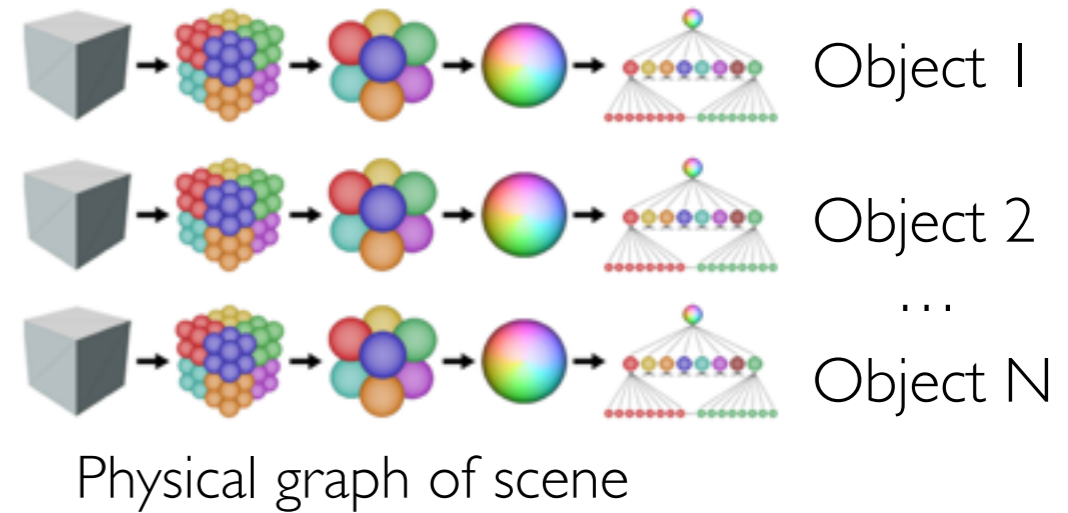…less easy to understand how to do this in material-agnostic way.

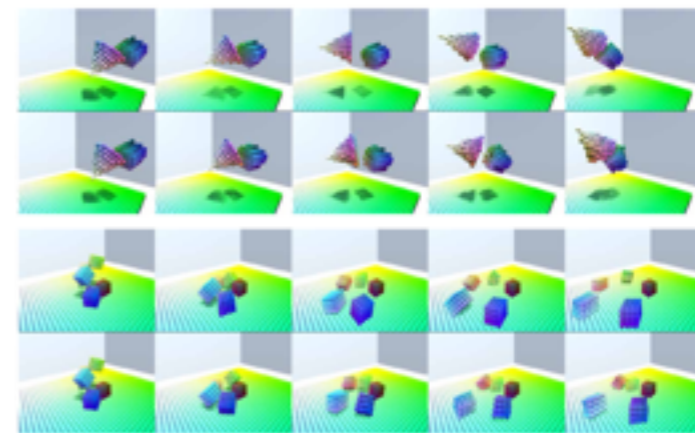# Challenges:

Extracting the graph description from video.



ConvRNN

Object 1

Object 2

...

Object N

Physical graph of scene

Scene at time T = 0, 1, ..., t

Hierarchical Relation Network

Rendering

Scene at time T = t+1, t+2, ...

| t+1 | t+2 | t+3 | t+4 | t+5 |

Predictions of graph in future

ConvRNN

Object 1

Object 2

...

Object N

Physical graph of scene

Hierarchical Relation Network

Predictions of graph in future

Retina ⇧ ~1 M (RCG representation)

LGN ⇧ ~1 M (LGN representation)

V1 ~37 M (V1 representation) ~190 M

V2 ~29 M (V2 representation) ~150 M

V3

PIP

MIP

DP

PO

MT

LIP

MST

FST

7a

V3A

VOT

V4 ~15 M (V4 representation) ~68 M

PIT ~36 M

CIT ~17 M

AIT ~16 M

STP_p

STP_a

~10 M (IT representation)

~40 ms

~50 ms

~60 ms

~70 ms

~80 ms

~90 ms

~100 ms

Latency

t+1

t+2

t+3

t+4

t+5

# Human-centered feedback loop

# Thanks!