

Neuro-Symbolic Embodied AI

Chuang Gan



MIT-IBM
Watson
AI Lab



My 15-Month-Old Daughter (Carolyn)



Recognize objects and substances

→ Plate

→ Hole

Reason from abstract knowledge

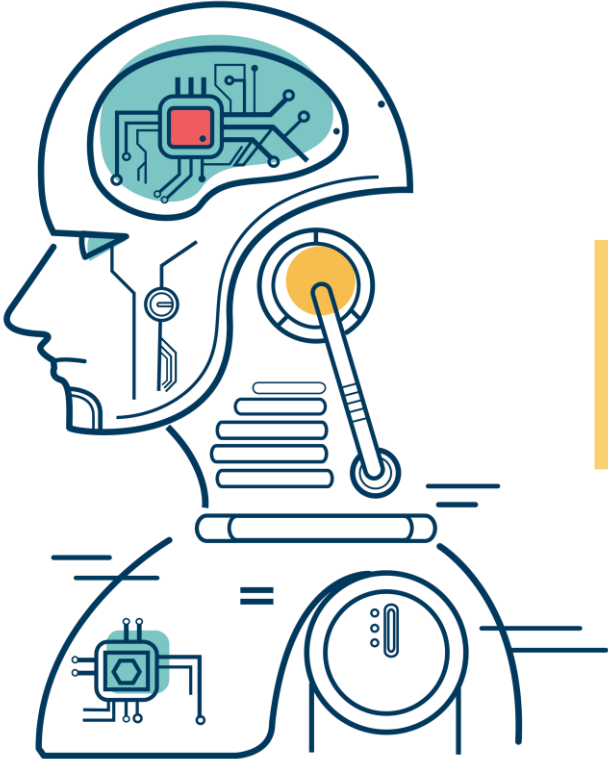
Not Fit

Fit

Use world models for inference and planning

→ Try to put the plate into the correct hole

How About Machine Intelligence Today?



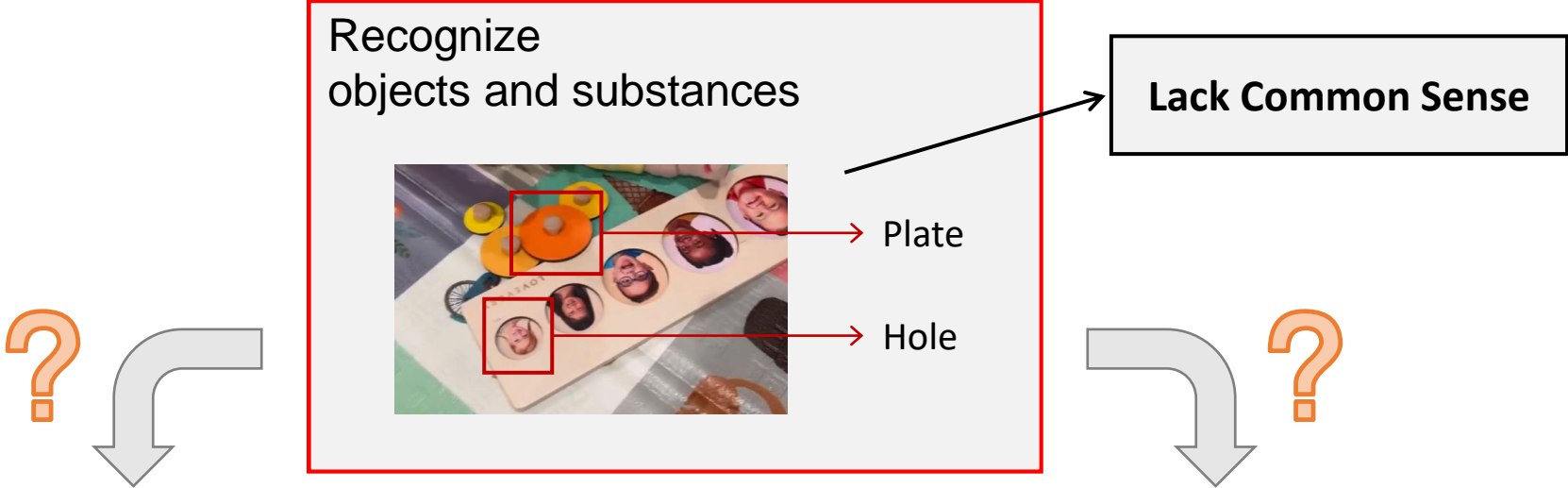
Machine

Does machine have the intelligence at the the level of Carolyn?



Carolyn

Most Visual Intelligence Now



Reason from abstract knowledge

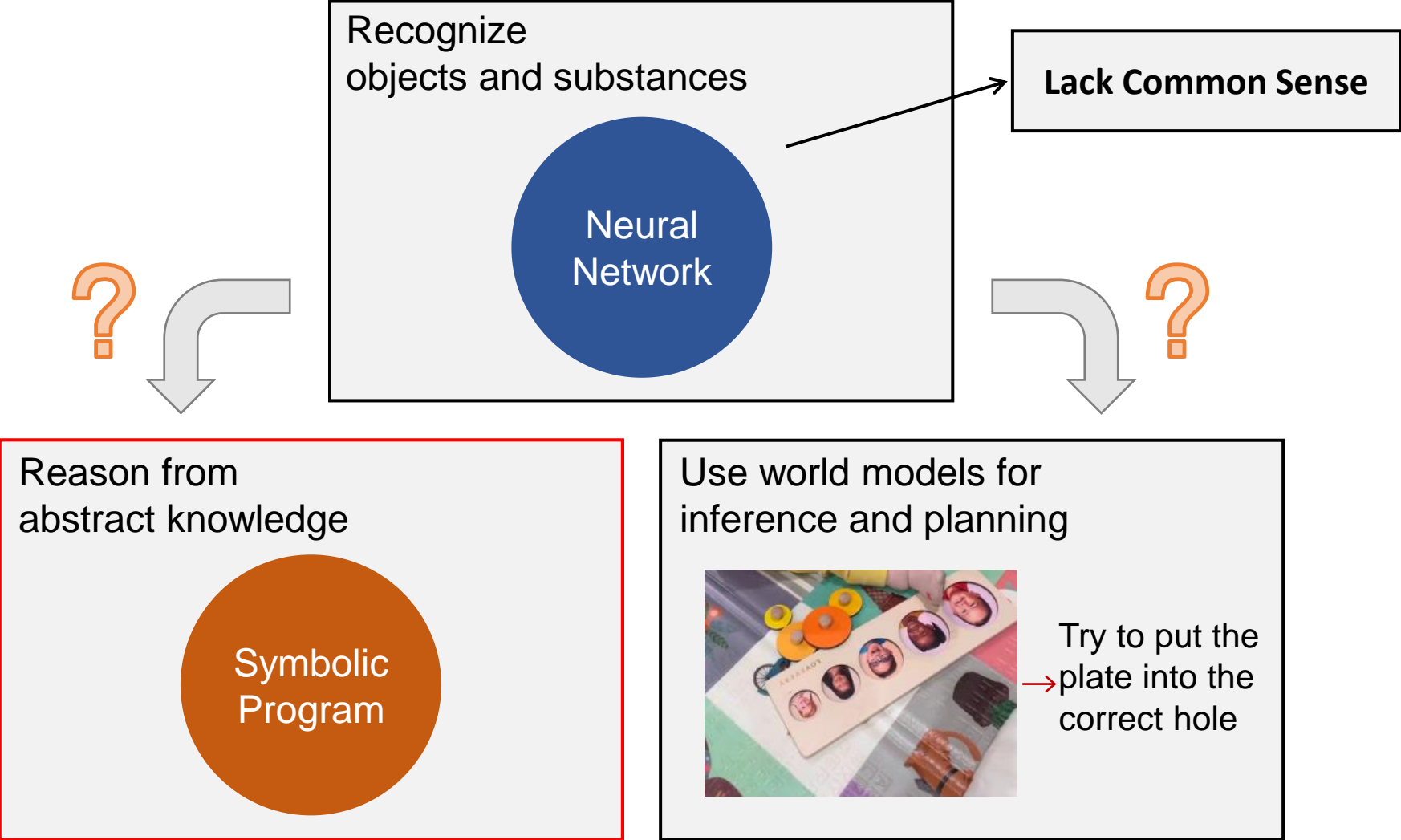
Not Fit

Fit

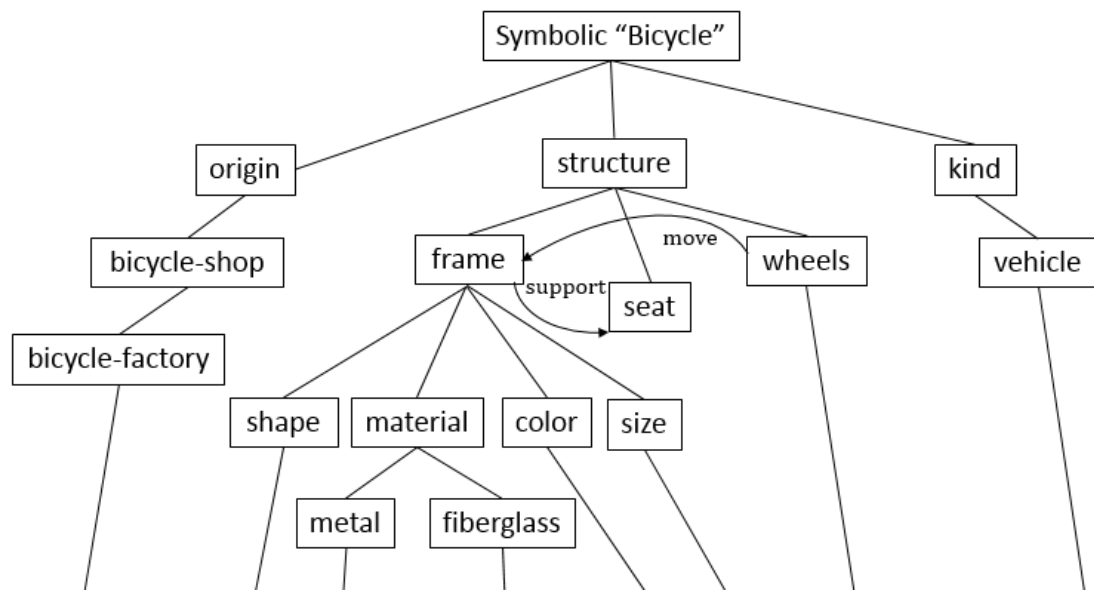
Use world models for inference and planning

Try to put the plate into the correct hole

Symbolic Program

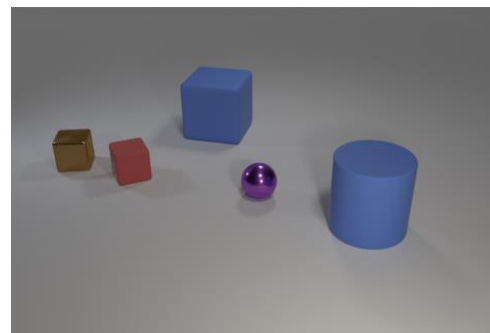


Related Work: Symbolic AI



A symbolic representation of bicycle [1]

Q: How many cubes are behind the cylinder?
A: 3



```
1. filter_shape(scene, cylinder)
   ↓
2. relate(behind)
   ↓
3. filter_shape(scene, cube)
   ↓
4. count(scene)
```

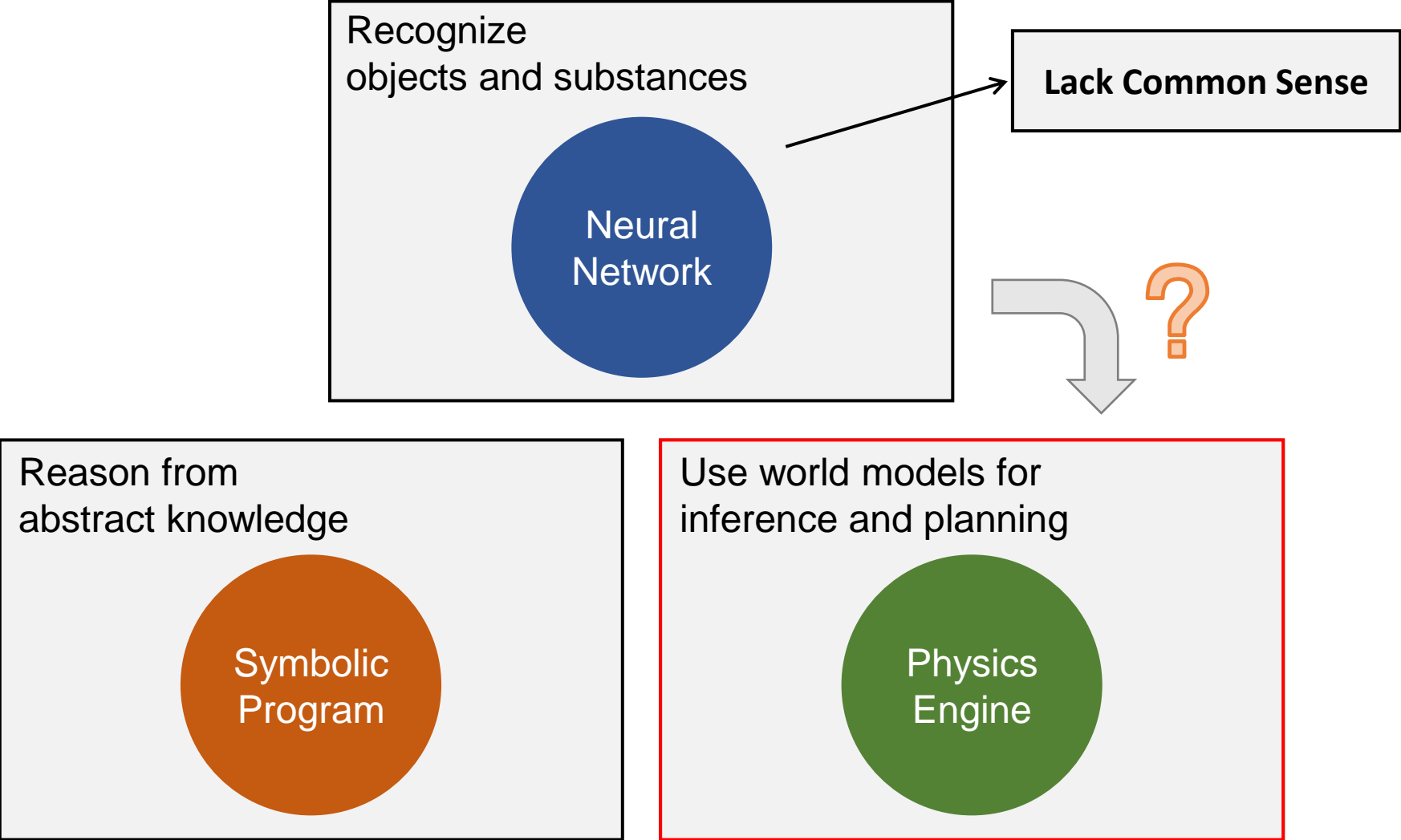
Program

A symbolic program of CLEVR [2]

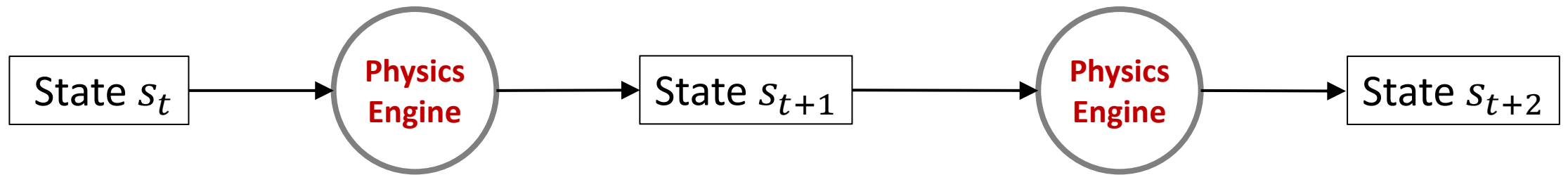
[1] Minsky, M., & Winston, P. H. (1990). Logical vs. Analogical or Symbolic vs. Connectionist or Neat vs. Scruffy” and “Excerpts from the Society of Mind. In Artificial Intelligence at MIT, Expanding Frontiers (Vol. 1).

[2] Johnson J, Hariharan B, Van Der Maaten L, et al. Clevr: A diagnostic dataset for compositional language and elementary visual reasoning. Proceedings of the IEEE conference on computer vision and pattern recognition.

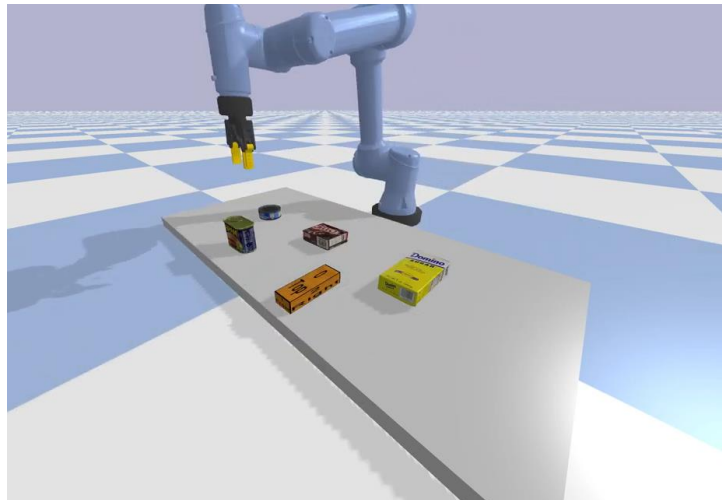
Physics Engines



Related Work: Physics Engines



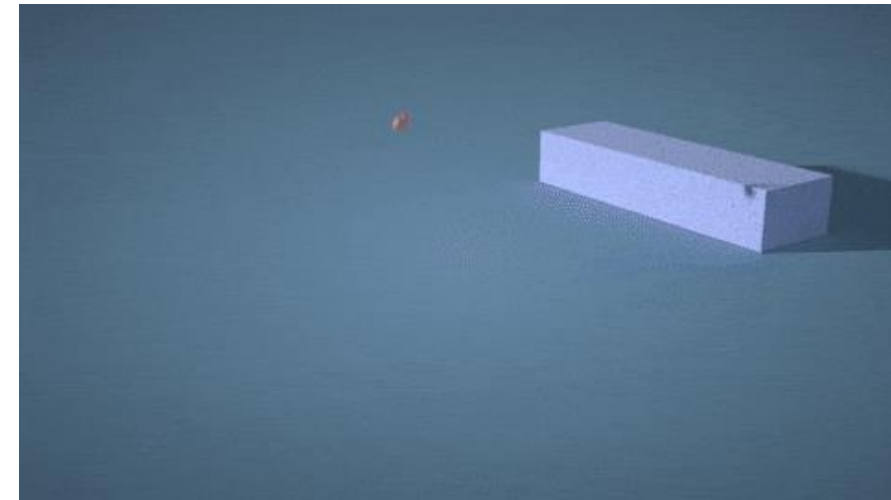
PyBullet



NVIDIA-Flex



Taichi Programming Language

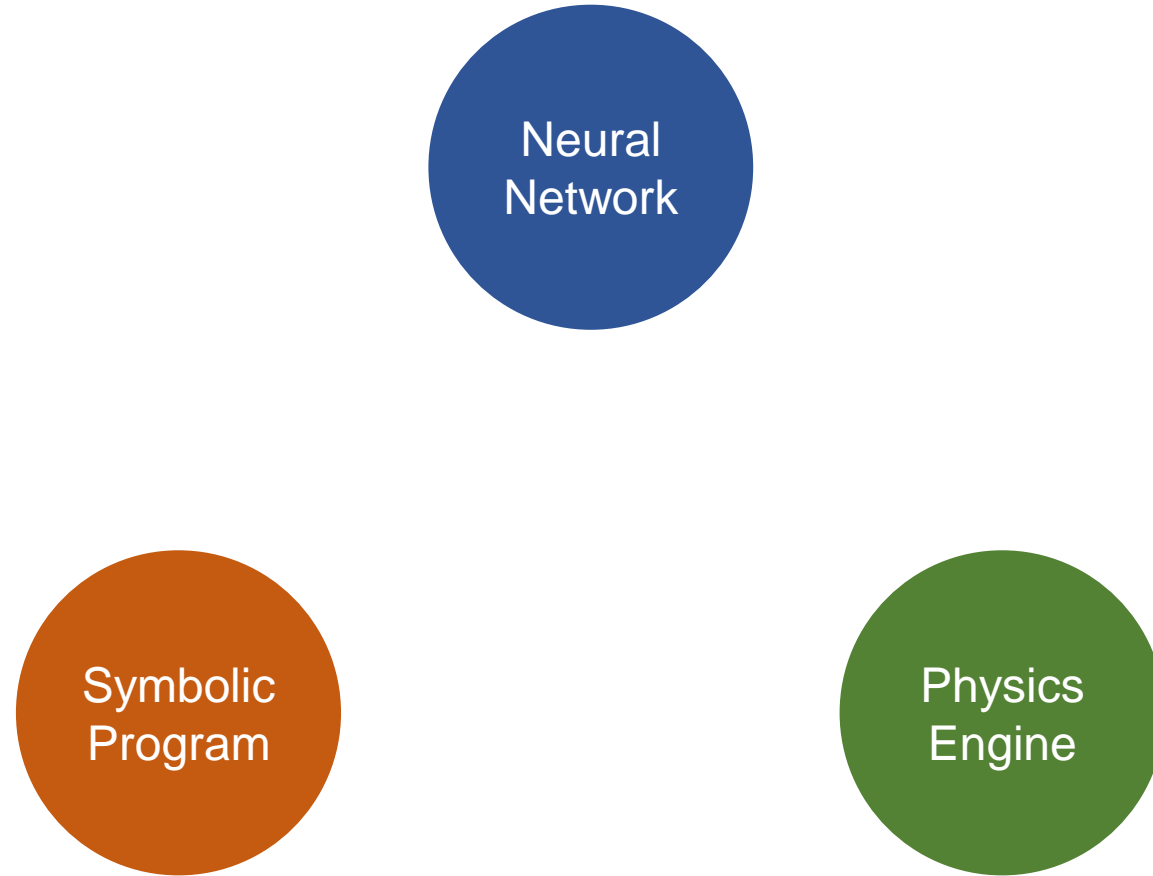


PyBullet: <https://www.youtube.com/watch?v=8e-KjBUakqY>

NVIDIA-Flex: <https://youtu.be/1o0Nuq71gl4>

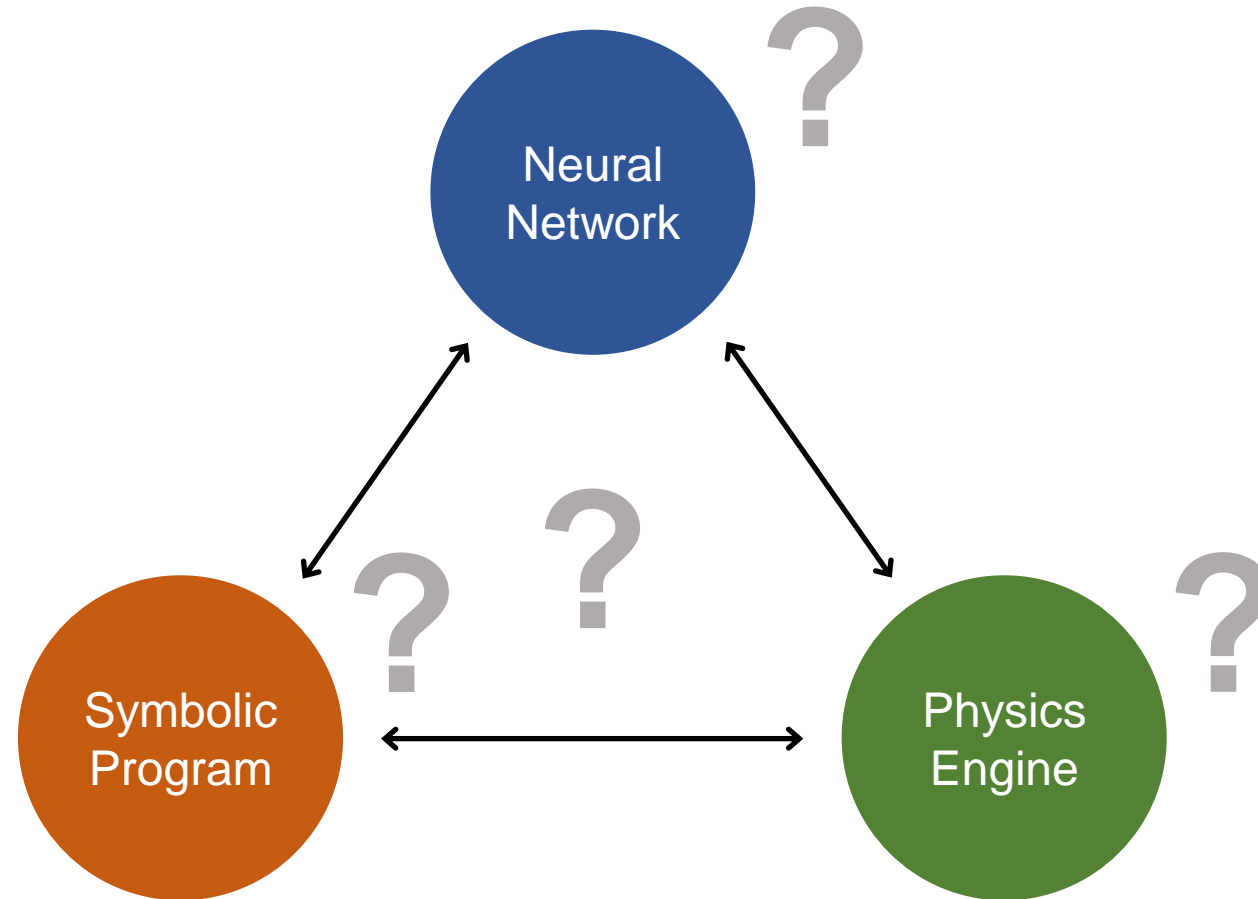
Taichi: <https://github.com/taichi-dev/quantaichi>

How to Integrate Them to Advance Visual Intelligence?



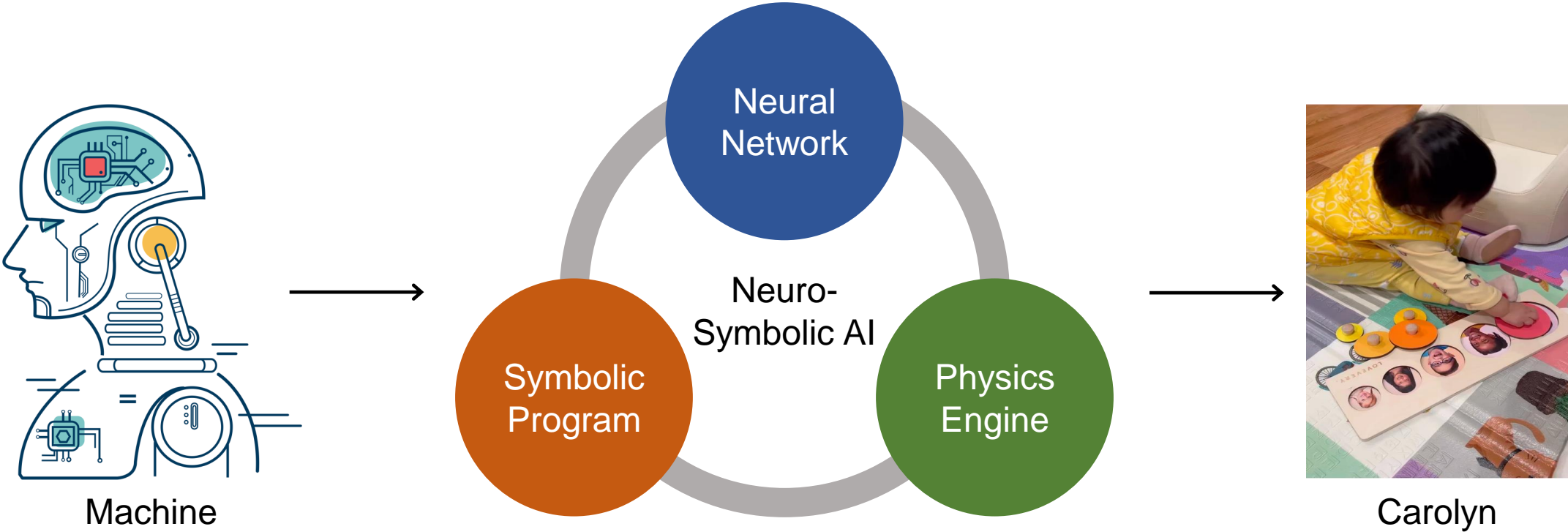
How to Integrate Them to Advance Visual Intelligence?

What is the role of each one?



How to combine each strength?

My Research: Neuro-Symbolic AI for Visual Intelligence

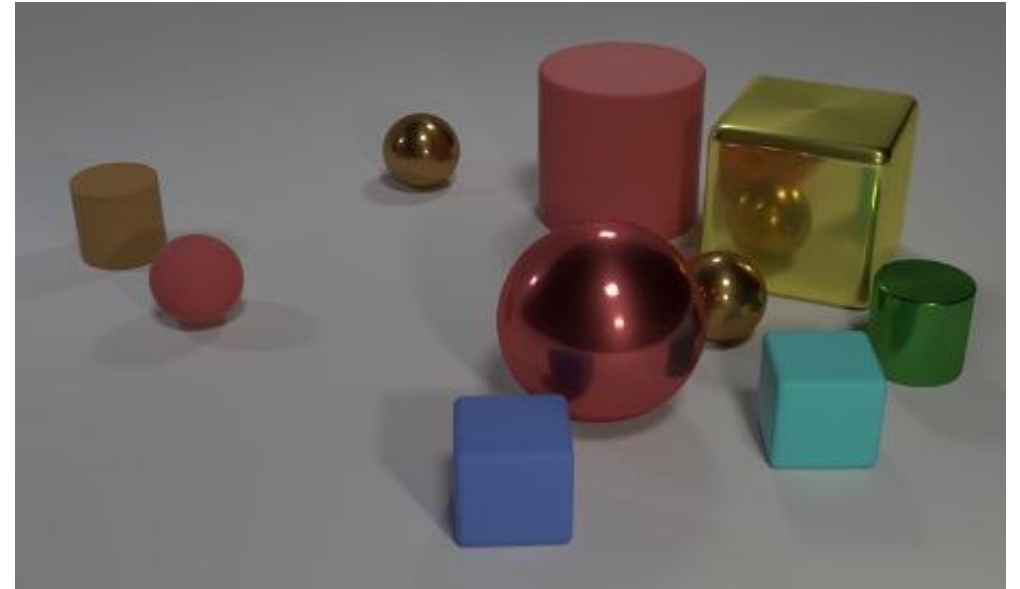


Task: Visual Reasoning



Q: What color is the fire hydrant?

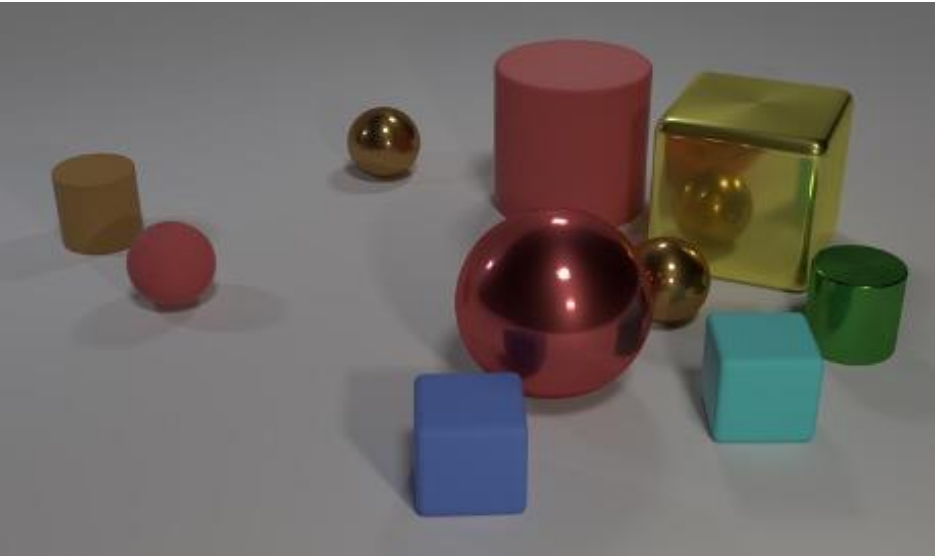
A: Yellow



Q: Are there an equal number of large things and metal spheres?

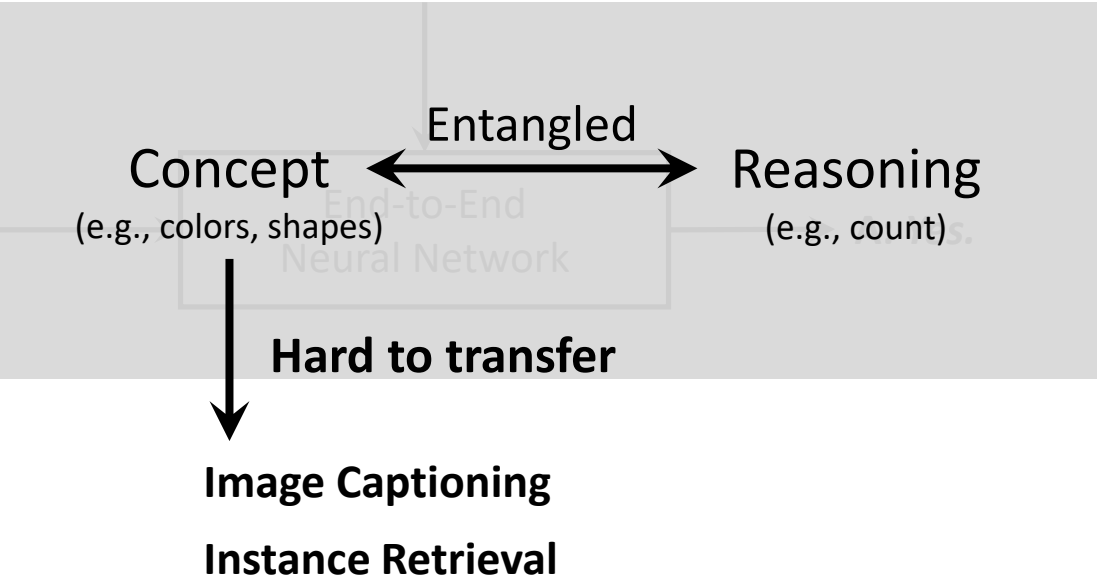
A: Yes

Prior Work: End-to-End Visual Reasoning



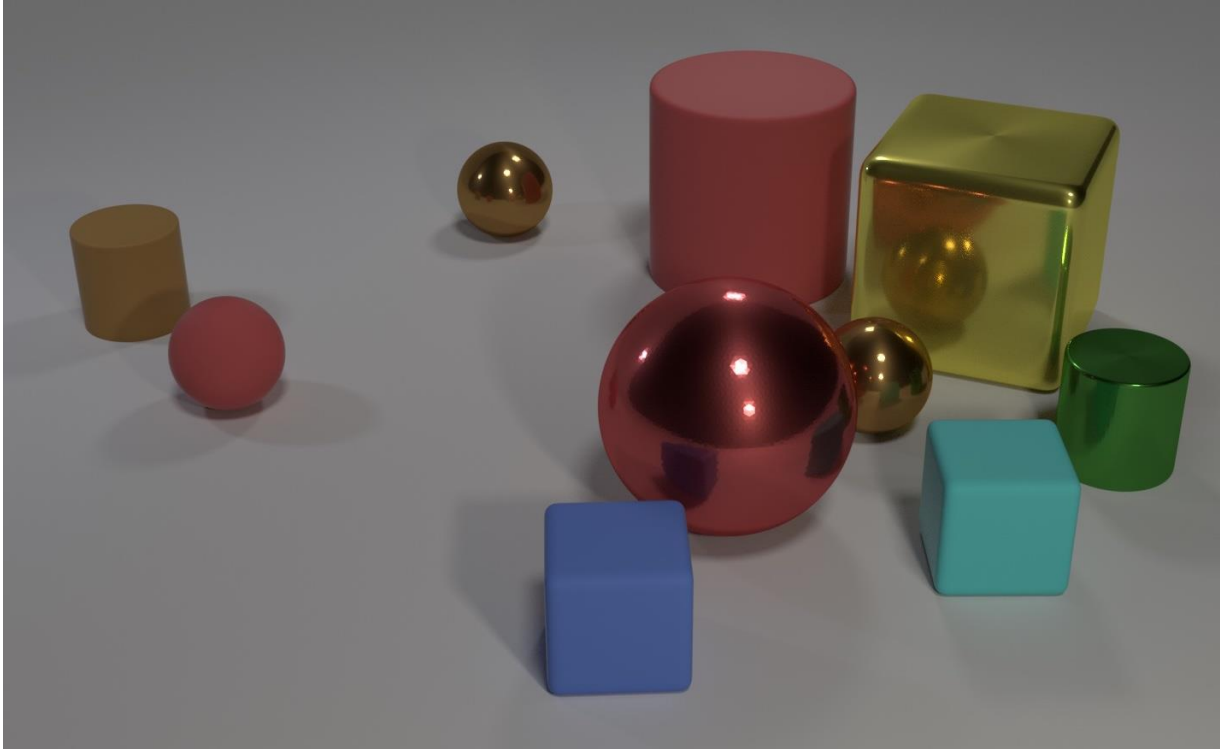
Visual Question Answering

Q: Are there an **equal number of large things** and **metal spheres**?



Agrawal et al. VQA, 2015. Johnson et al. CLEVR 2017.
Johnson et al. CLEVR 2017, Andreas et al. NMN, 2016. Johnson et al. IEP, 2017. Perez et al. FiLM, 2018. Hudson & Manning. MAC, 2018. Hu et al. Stack-NMN, 2018. Mascharka et al. TbD, 2018.

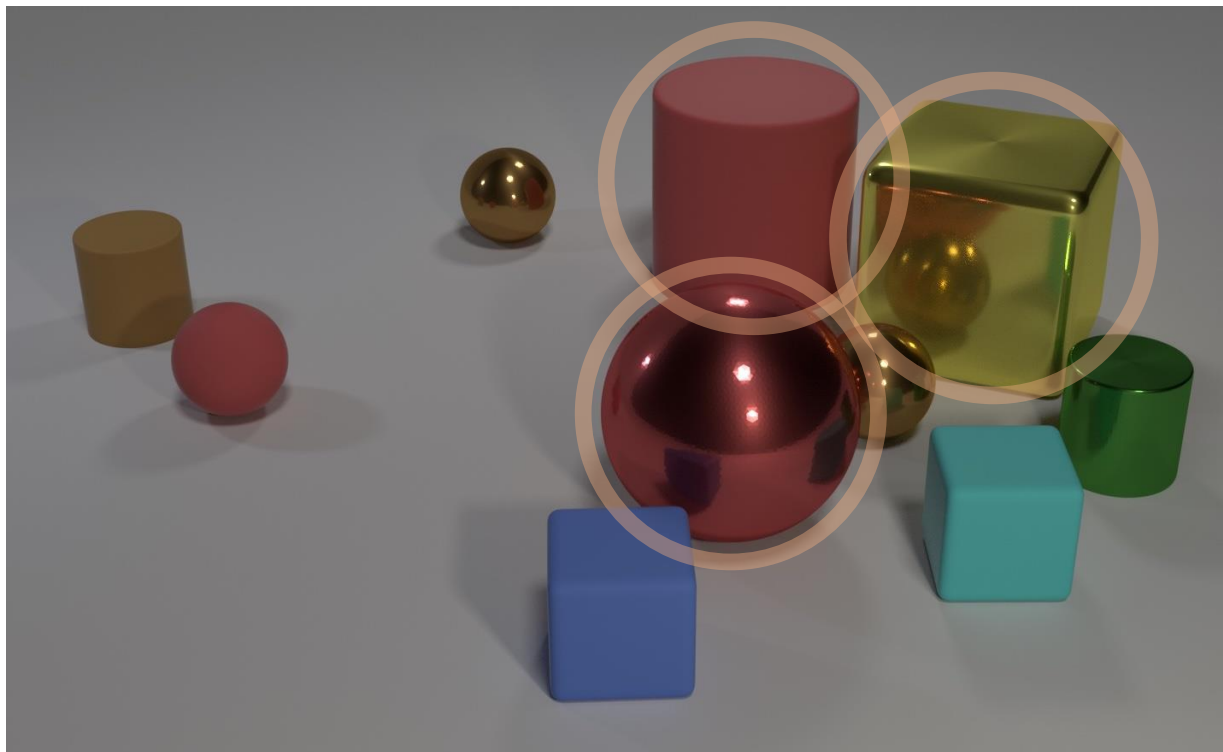
How Do Humans Reason From A Visual Scene?



Question: *Are there an equal number of large things and metal spheres?*

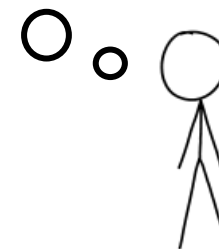


How Do Human Reason From A Visual Scene?

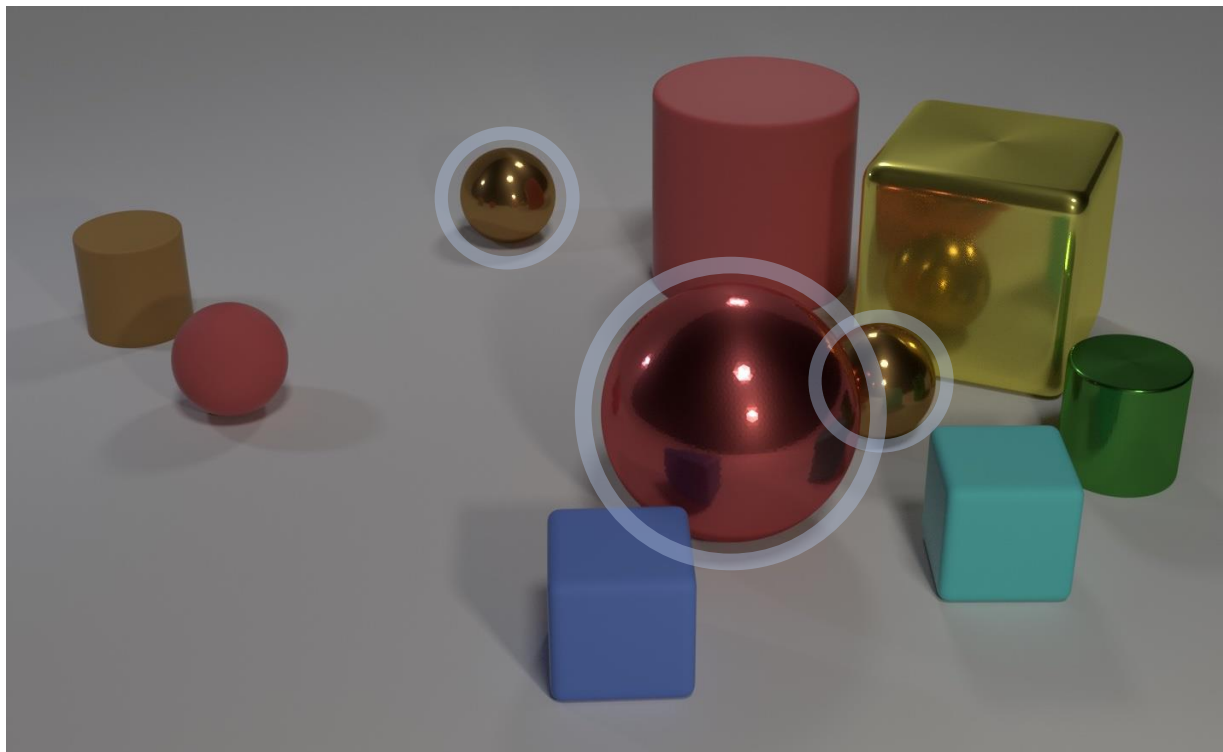


Question: *Are there an equal number of **large things** and metal spheres?*

3 large things!



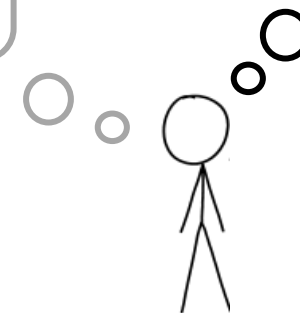
How Do Human Reason From A Visual Scene?



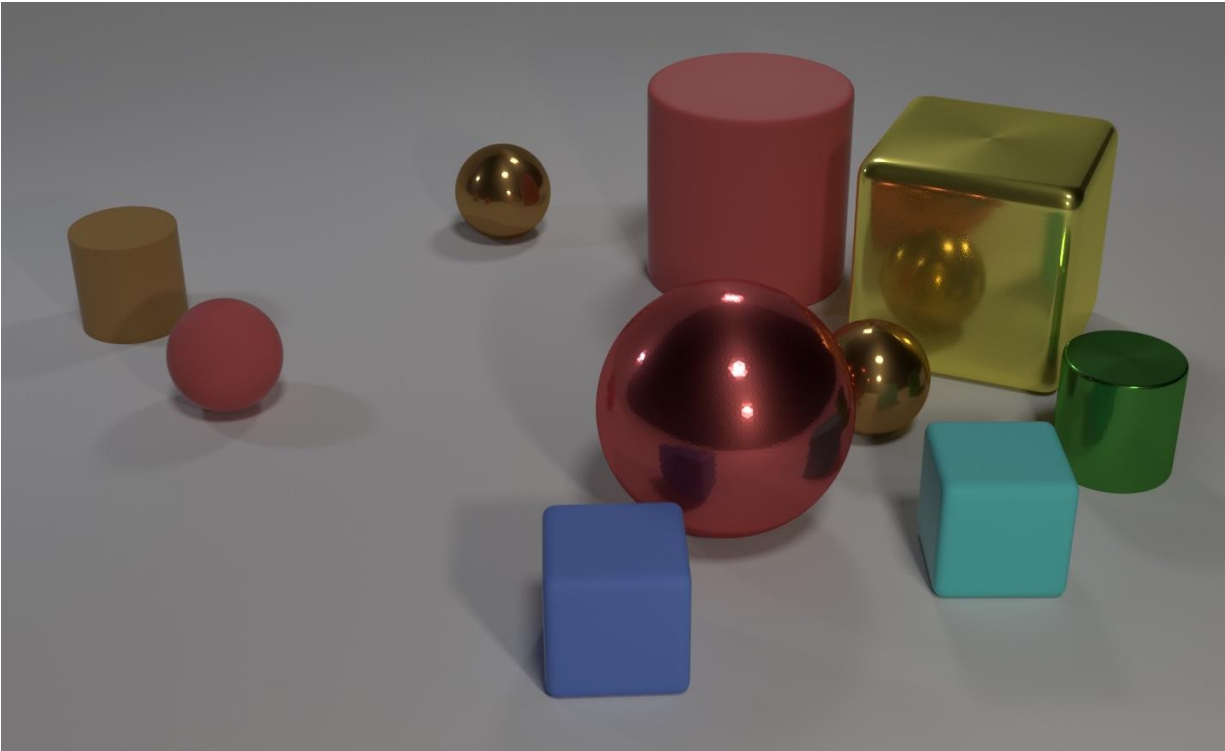
Question: *Are there an equal number of large things and **metal spheres**?*

3 large things!

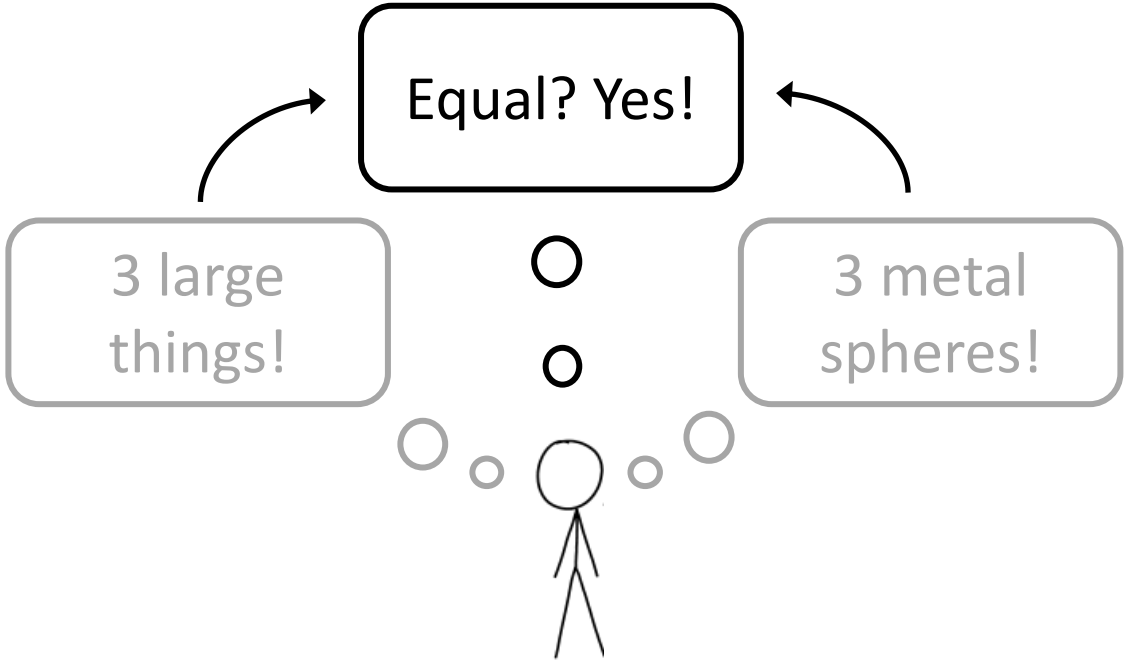
3 metal spheres!



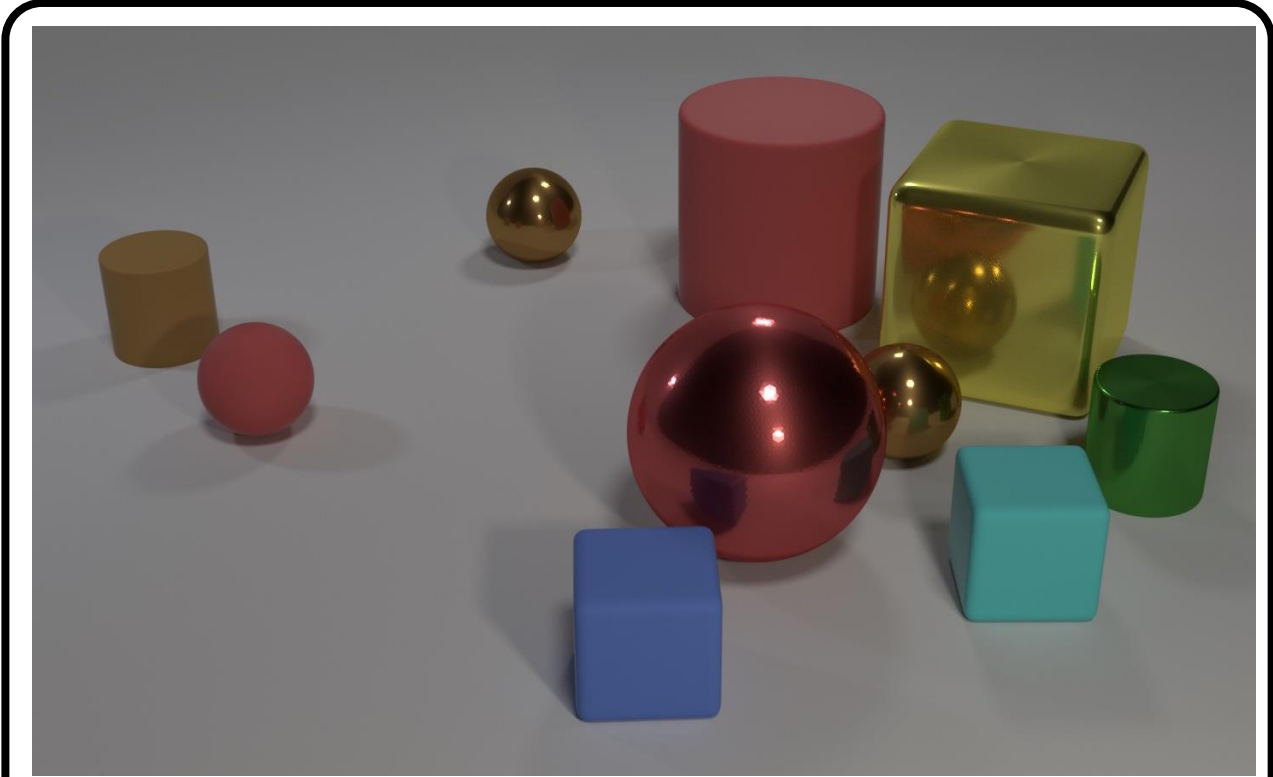
How do human reason from a visual scene?



Question: *Are there an **equal number** of large things and metal spheres?*



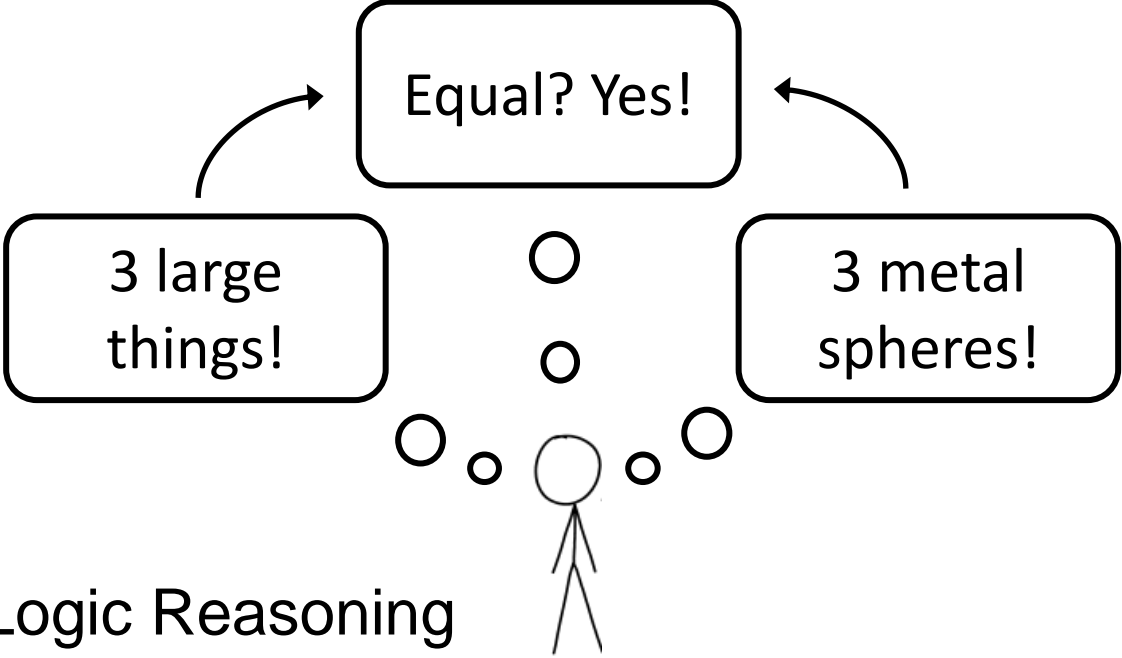
How Do Human Reason From A Visual Scene?



Visual Perception

Question Understanding

Question: *Are there an equal number of large things and metal spheres?*



Incorporate Concepts and Symbolic Programs



ID	Size	Shape	Material	Color	x	y	z
1	Small	Cube	Wood	Purple	0.45	-1.10	0.35
2	Large	Cube	Wood	Green	0.45	1.31	0.35
3	Large	Cylinder	Metal	Green	1.58	-1.60	0.70

Neural networks parse images in **symbolic** concepts

I. Neural Scene Parsing

II. Neural Question Parsing



Neural networks parse questions into **symbolic** functional program.

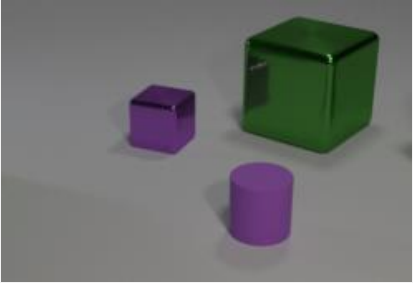
III. Symbolic Program Execution

- 1. filter_cylinder
- 2. relate_behind
- 3. filter_cube
- 4. filter_large
- 5. count

Executing programs on the **symbolic** space.

VQS: Linking Segmentations to Questions and Answers for VQA. Gan et al. ICCV'17
Neural-Symbolic VQA: Disentangling Reasoning from Vision and Language Understanding. Yi, Wu, Gan, et al. NeurIPS'18

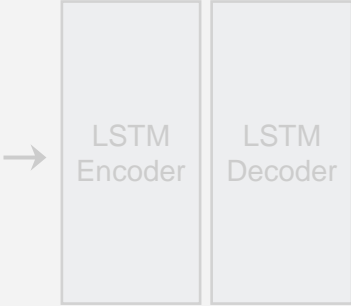
Neural Scene Parsing



I. Neural Scene Parsing

II. Neural Question Parsing

How many cubes that are behind the cylinder are large?



- 1. filter_shape(scene, cylinder)
- 2. relate(behind)
- 3. filter_shape(scene, cube)
- 4. filter_size(scene, large)
- 5. count(scene)

III. Symbolic Program Execution

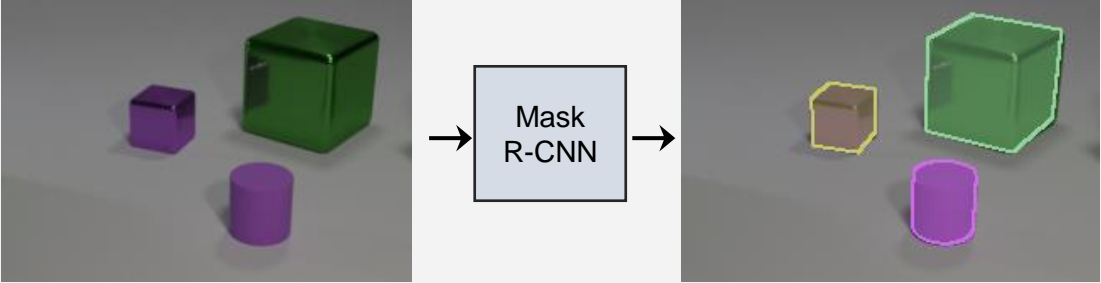
- 1. filter_cylinder
- 2. relate_behind
- 3. filter_cube
- 4. filter_large
- 5. count

ID	Size	Shape	...
1	Small	Cube	...
2	Small	Cylinder	...
3	Large	Cube	...

ID	Size	...
3	Large	...

Answer: 1

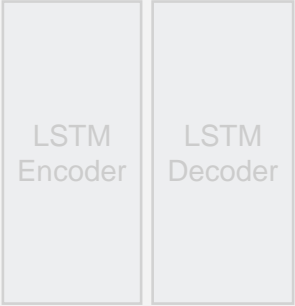
Neural Scene Parsing



I. *Neural* Scene Parsing

II. *Neural* Question Parsing

How many cubes that are behind the cylinder are large?



- 1. filter_shape(scene, cylinder)
- 2. relate(behind)
- 3. filter_shape(scene, cube)
- 4. filter_size(scene, large)
- 5. count(scene)

III. *Symbolic* Program Execution

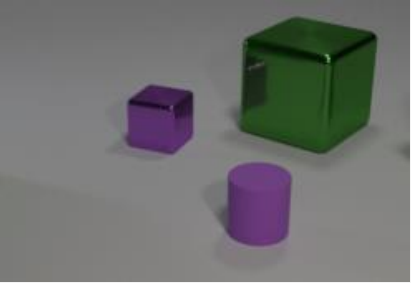
- 1. filter_cylinder
- 2. relate_behind
- 3. filter_cube
- 4. filter_large
- 5. count

ID	Size	Shape	...
1	Small	Cube	...
2	Small	Cylinder	...
3	Large	Cube	...

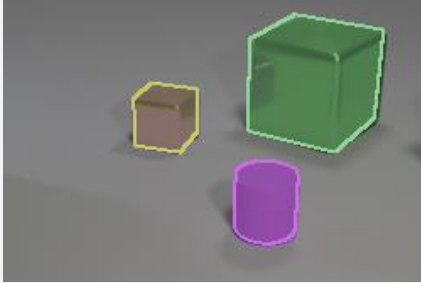
ID	Size	...
3	Large	...

Answer: 1

Neural Scene Parsing



Mask R-CNN



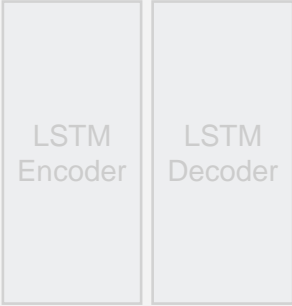
CNN

ID	Size	Shape	Material	Color	x	y	z
1	Small	Cube	Metal	Purple	-0.45	-1.10	0.35
2	Small	Cylinder	Rubber	Purple	0.75	1.31	0.35
3	Large	Cube	Metal	Green	1.58	-1.60	0.70

I. Neural Scene Parsing

II. Neural Question Parsing

How many cubes that are behind the cylinder are large?



- 1. filter_shape(scene, cylinder)
- 2. relate(behind)
- 3. filter_shape(scene, cube)
- 4. filter_size(scene, large)
- 5. count(scene)

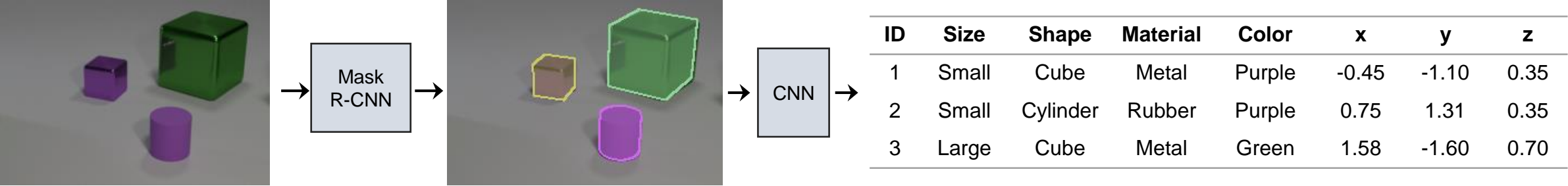
III. Symbolic Program Execution

- 1. filter_cylinder
- 2. relate_behind
- 3. filter_cube
- 4. filter_large
- 5. count

ID	Size	Shape	...	ID	Size	...
1	Small	Cube	...	3	Large	...
2	Small	Cylinder	...			
3	Large	Cube	...			

Answer: 1

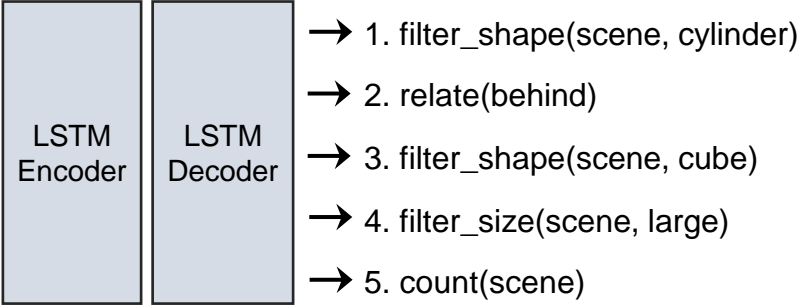
Neural Question Parsing



I. Neural Scene Parsing

II. Neural Question Parsing

How many cubes that are behind the cylinder are large?



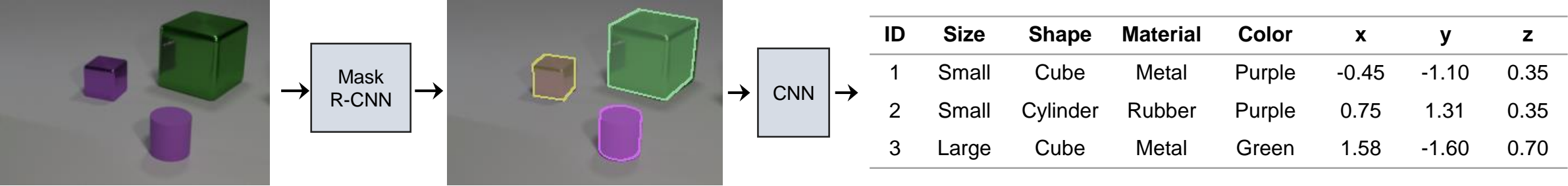
III. Symbolic Program Execution

1. filter_cylinder	3. filter_cube	5. count
2. relate_behind	4. filter_large	

ID	Size	Shape	...
1	Small	Cube	...
2	Small	Cylinder	...
3	Large	Cube	...

Answer: 1

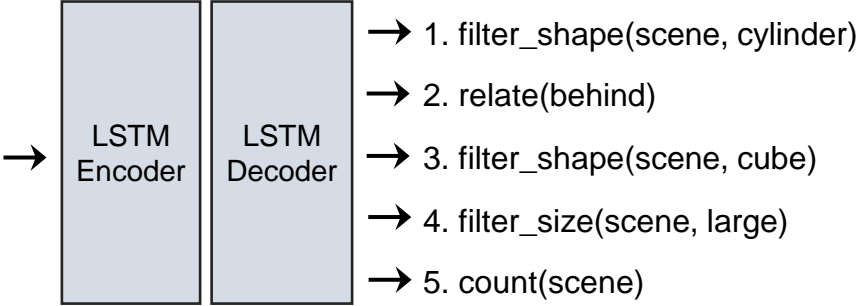
Symbolic Reasoning



I. *Neural* Scene Parsing

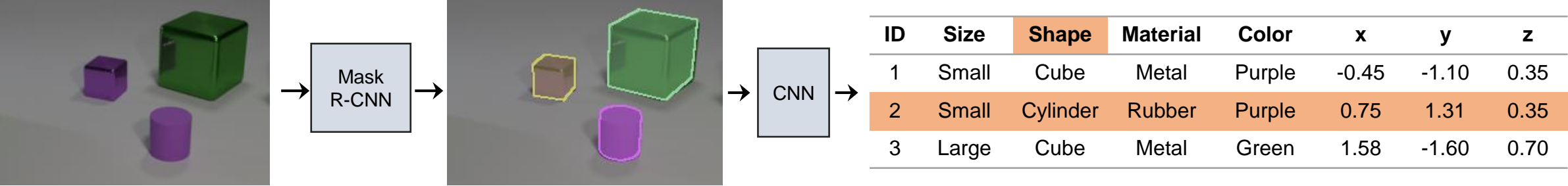
II. *Neural* Question Parsing

How many cubes that are behind the cylinder are large?



III. *Symbolic* Program Execution

Symbolic Reasoning

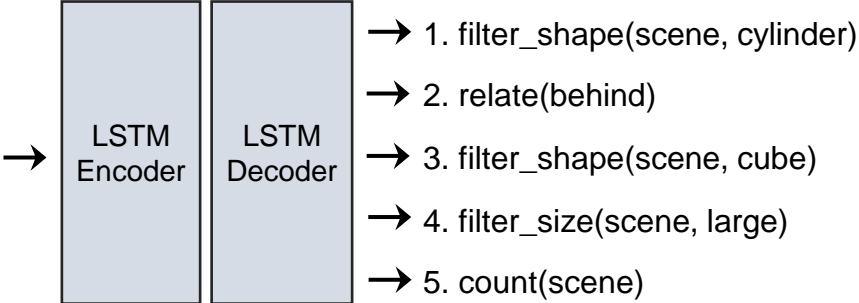


I. Neural Scene Parsing

II. Neural Question Parsing

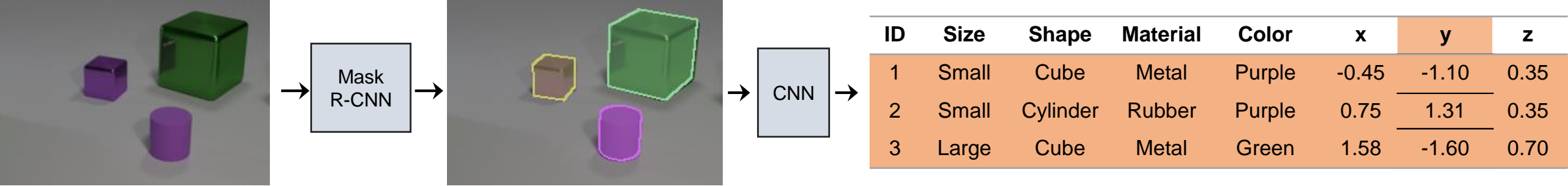
III. Symbolic Program Execution

How many cubes that are behind the cylinder are large?



1. filter_cylinder

Symbolic Reasoning

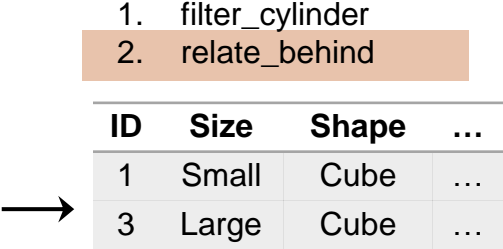
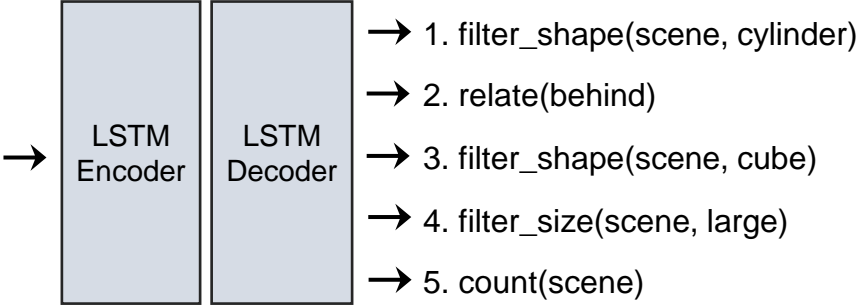


I. Neural Scene Parsing

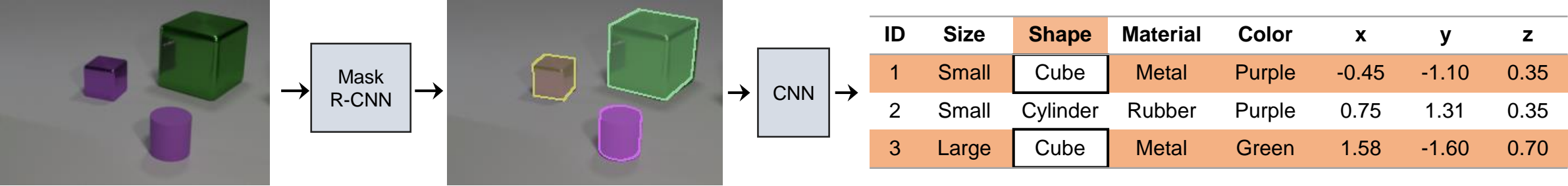
II. Neural Question Parsing

III. Symbolic Program Execution

How many cubes that are behind the cylinder are large?



Symbolic Reasoning

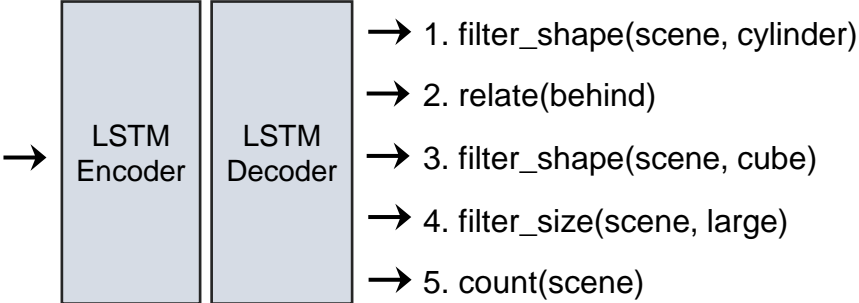


I. Neural Scene Parsing

II. Neural Question Parsing

III. Symbolic Program Execution

How many cubes that are behind the cylinder are large?

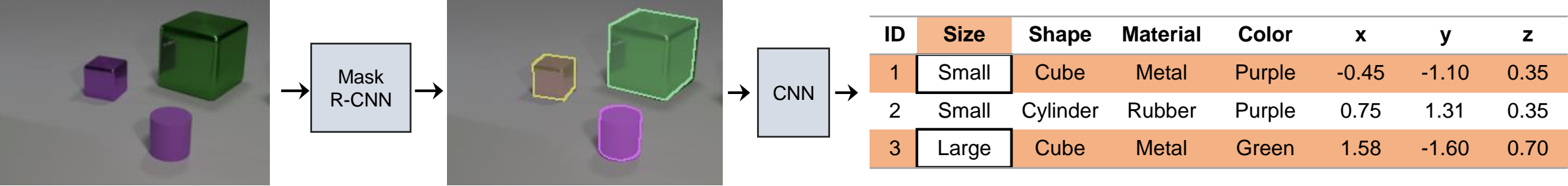


1. filter_cylinder
2. relate_behind

3. filter_cube

ID	Size	Shape	...
1	Small	Cube	...
3	Large	Cube	...

Symbolic Reasoning

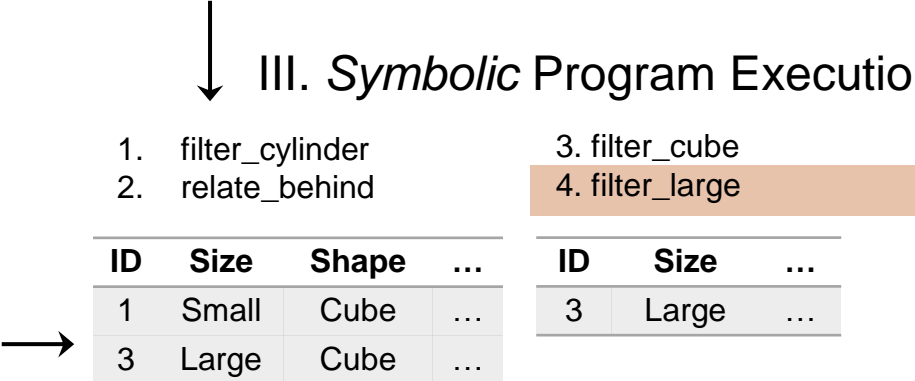
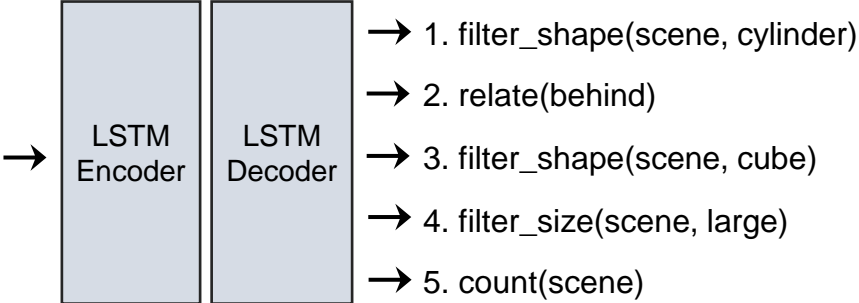


I. Neural Scene Parsing

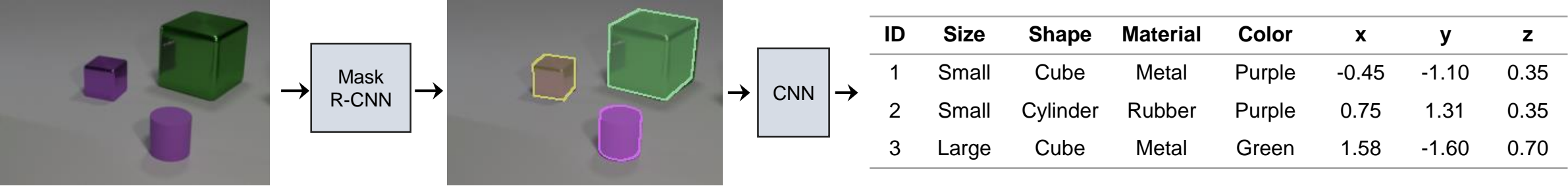
II. Neural Question Parsing

III. Symbolic Program Execution

How many cubes that are behind the cylinder are large?



Symbolic Reasoning

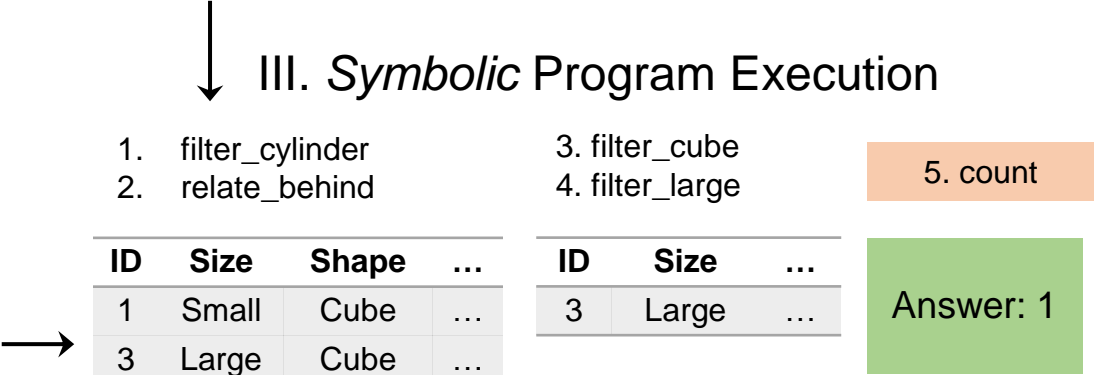
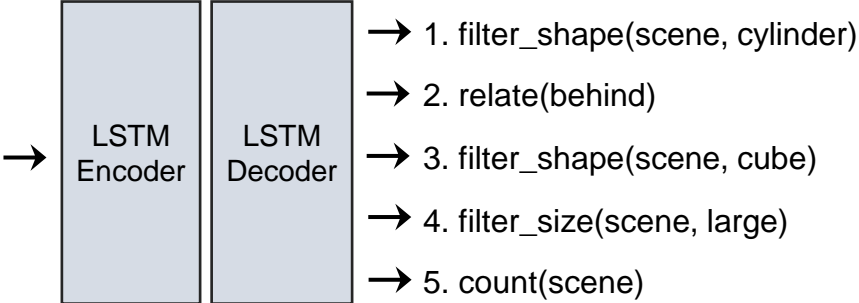


I. Neural Scene Parsing

II. Neural Question Parsing

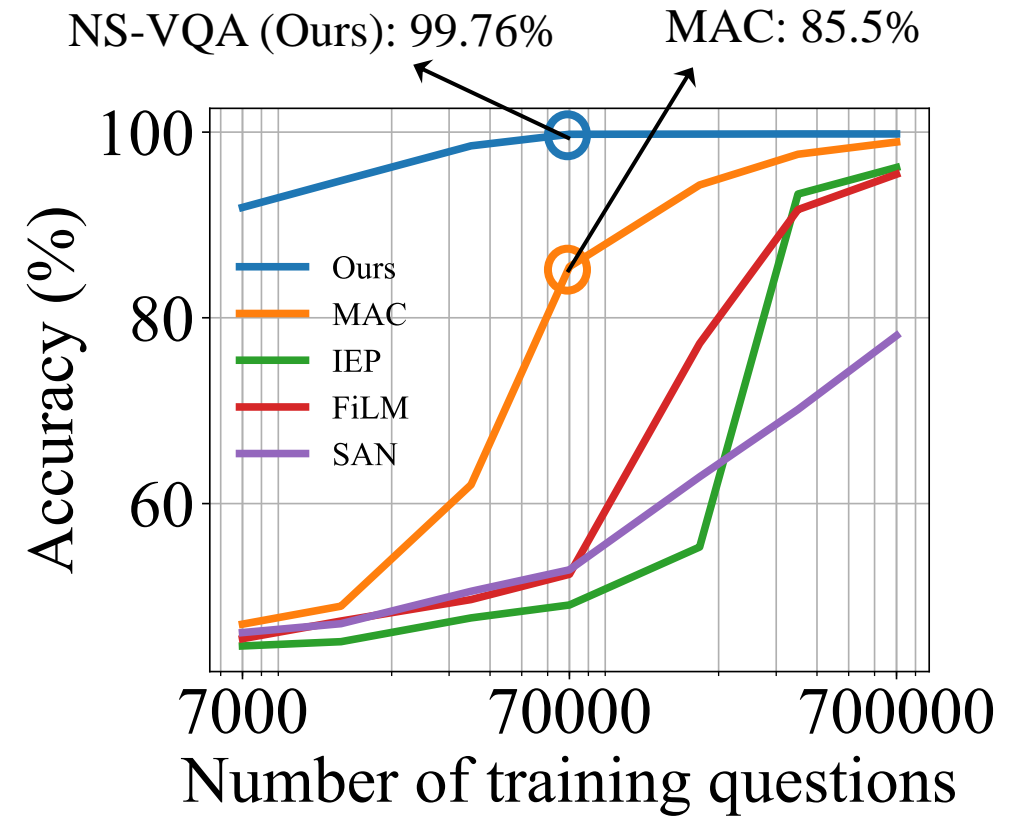
III. Symbolic Program Execution

How many cubes that are behind the cylinder are large?



Evaluation on CLEVR

Method	Accuracy (%)
Human	92.6
RN	95.5
IEP	96.9
FiLM	97.6
MAC	98.9
NS-VQA (Ours)	99.8



Part I: Summary

Incorporate symbolic programs for reasoning

Compositional Reasoning

[GLLSG] VQS, ICCV, 17.
[YWGTKT] NS-VQA, NeurIPS, 18.

Novel Concepts Learning

[MGKTW] NS-CL, ICLR, 19.
[HMGTW] Meta-NL, NeurIPS, 19.

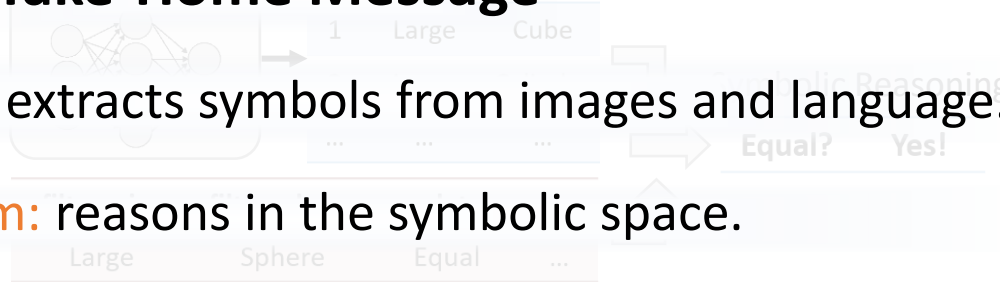
Compositional Reasoning



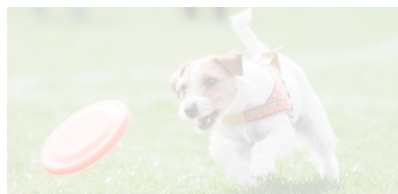
Take-Home Message

- **Neural Network:** extracts symbols from images and language.
- **Symbolic Program:** reasons in the symbolic space.

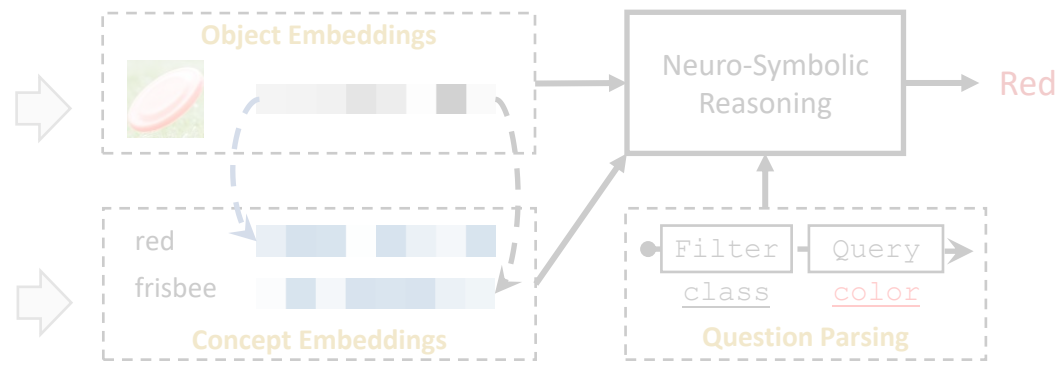
Are there large things and metal spheres?



Novel Concepts Learning



What is the color of the frisbee?



Part I: Incorporate Symbolic Programs

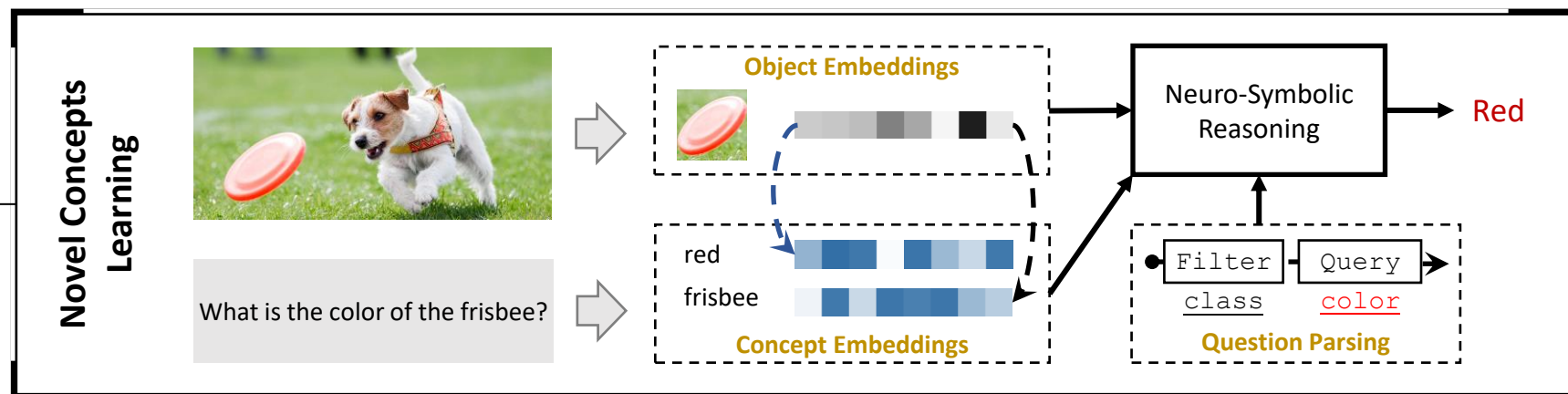
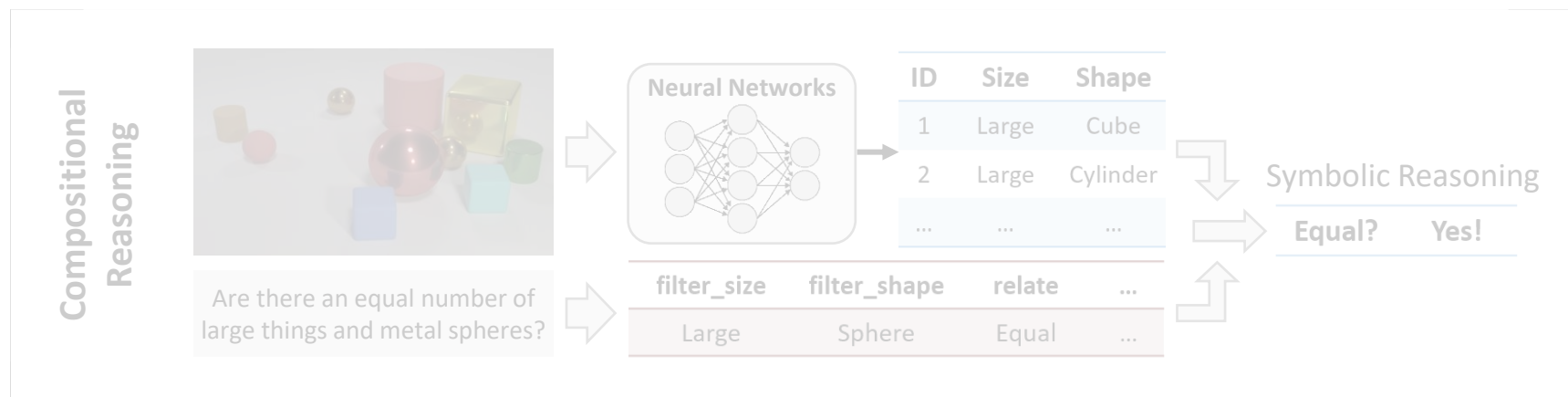
Incorporate symbolic programs for reasoning

Compositional Reasoning

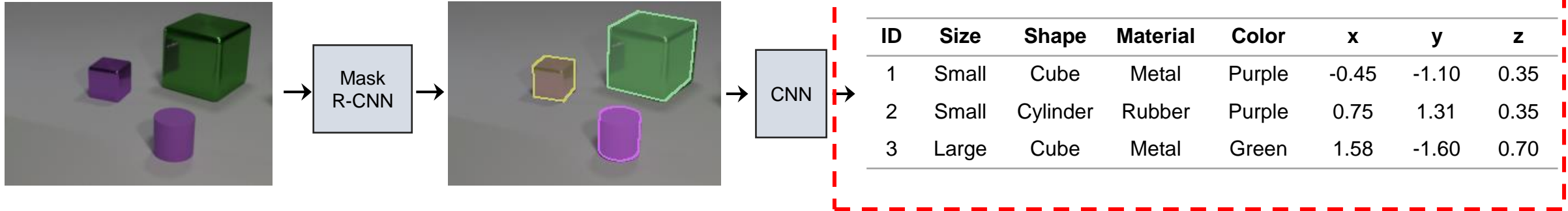
[GLLSG] VQS, ICCV, 17.
[YWGTKT] NS-VQA, NeurIPS, 18.

Novel Concepts Learning

[MGKTW] NS-CL, ICLR, 19.
[HMGTW] Meta-NL, NeurIPS, 19.



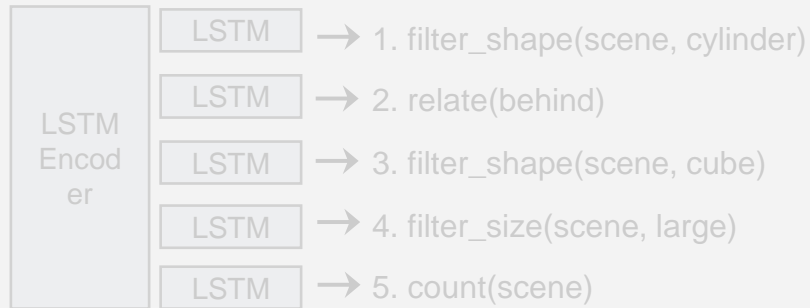
Limitation: Strong Requirement for Labeled Images



I. Neural Scene Parsing

II. Neural Question Parsing

How many cubes that are behind the cylinder are large?



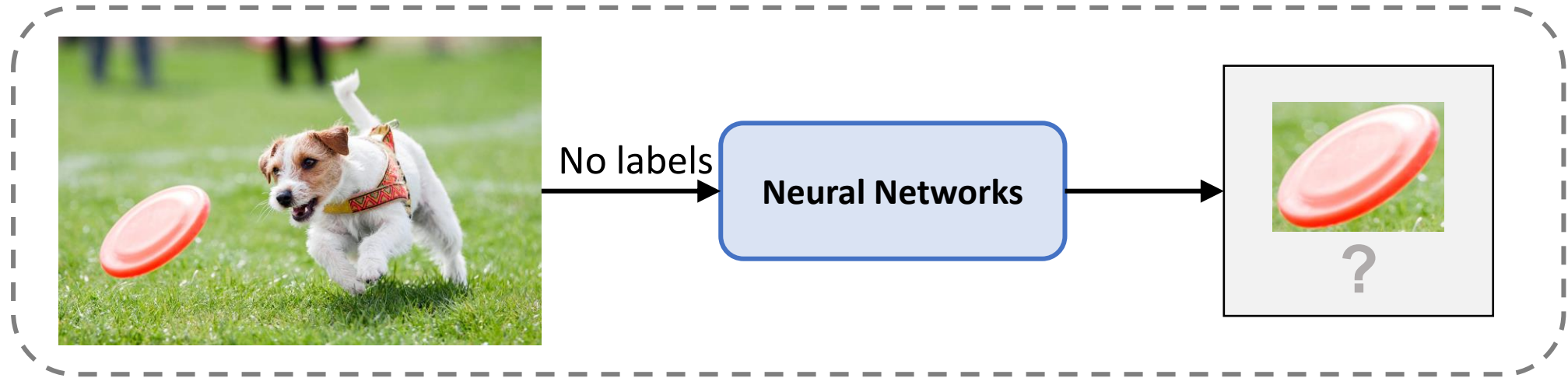
III. Symbolic Program Execution



VQS: Linking Segmentations to Questions and Answers for VQA. **Gan** et al. ICCV'17

Neural-Symbolic VQA: Disentangling Reasoning from Vision and Language Understanding. Yi, Wu, **Gan**, et al. NeurIPS'18

How About Images Without Concept Labels?



Our Idea: Learning Concepts From Weak Supervisions



No labels

Neural Networks

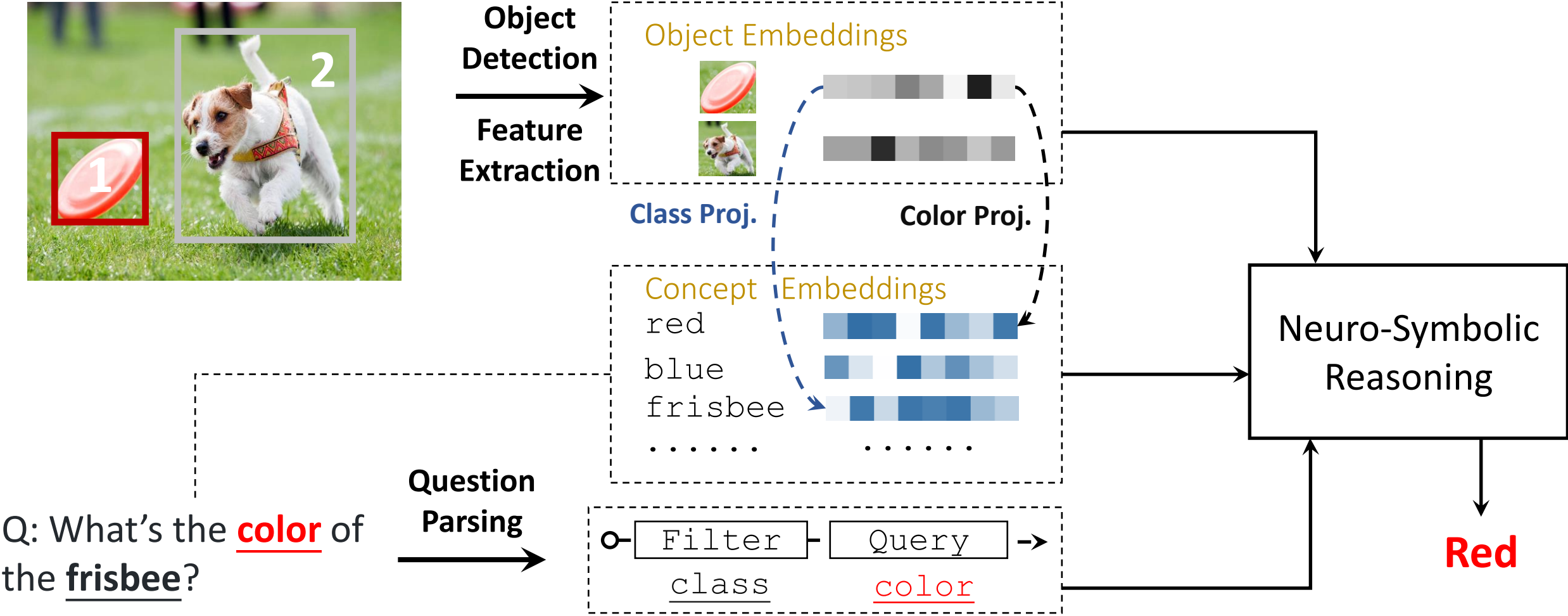


Neuro-Symbolic Reasoning



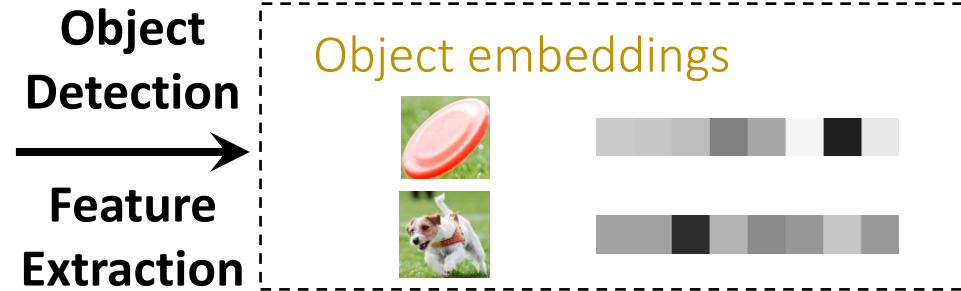
Q: What's the color of the Frisbee? A: Red.

Neuro-Symbolic Concepts Learner



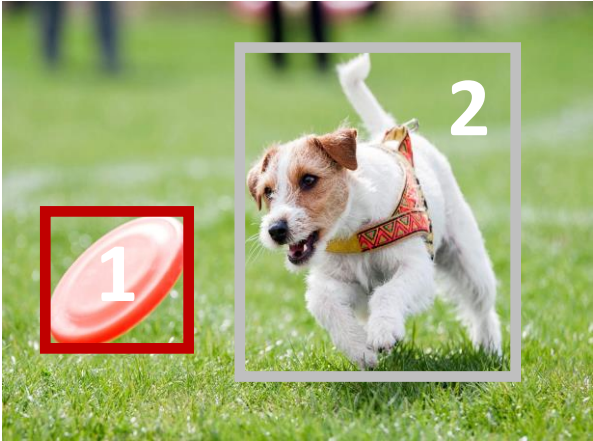
The neuro-symbolic concept learner: Interpreting scenes, words, and sentences from natural supervision. Mao, Gan, et al. ICLR'19

Object Embeddings

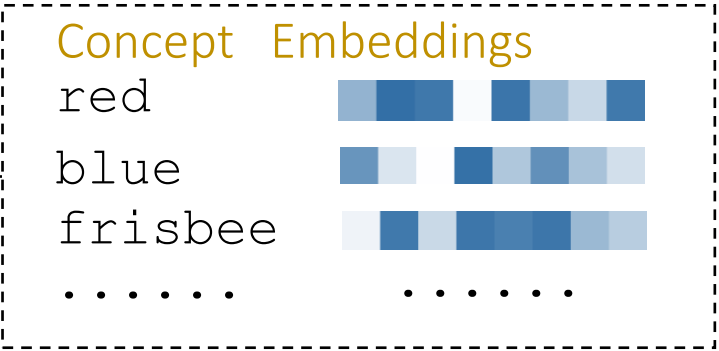


Q: What's the color of the frisbee?

Concept Embeddings

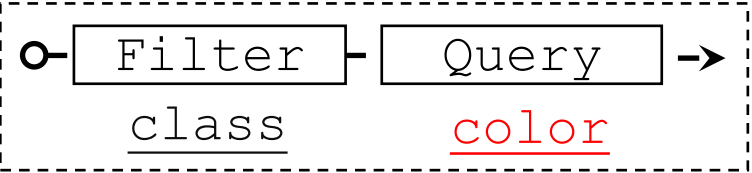


Object
Detection
→
Feature
Extraction

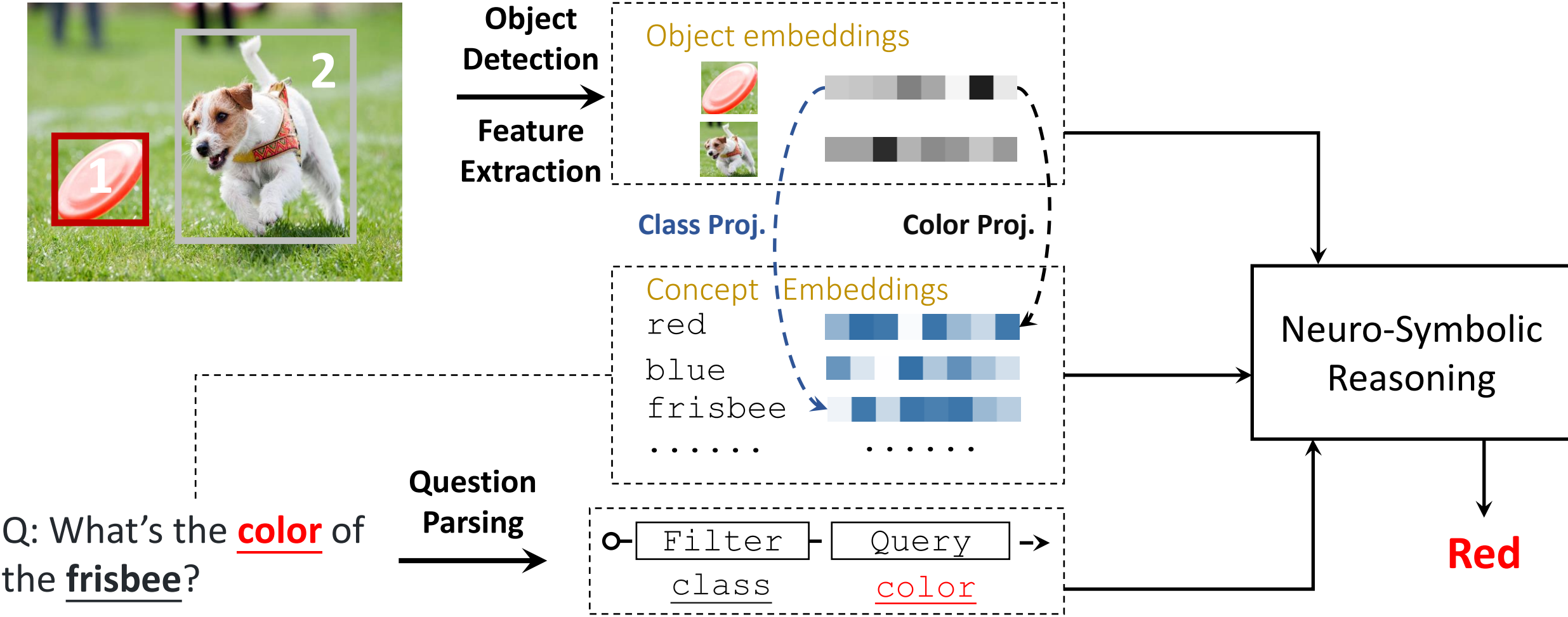


Q: What's the color of the frisbee?

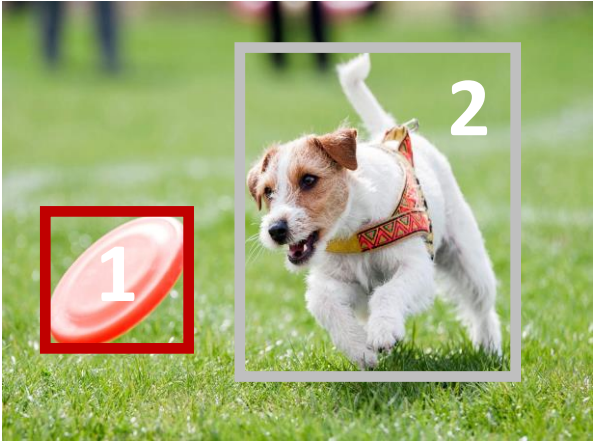
Question
Parsing
→



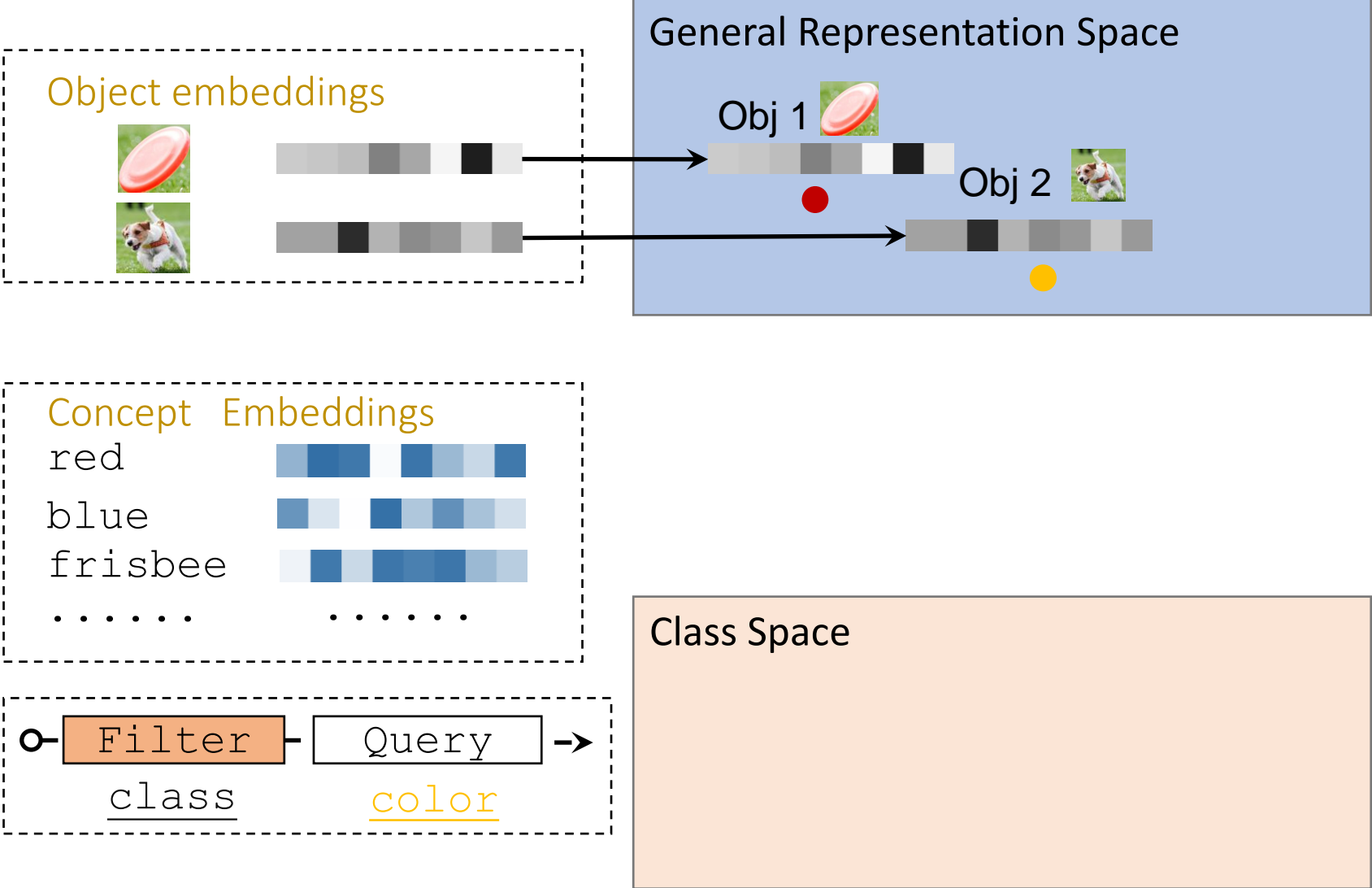
Neuro-symbolic Reasoning



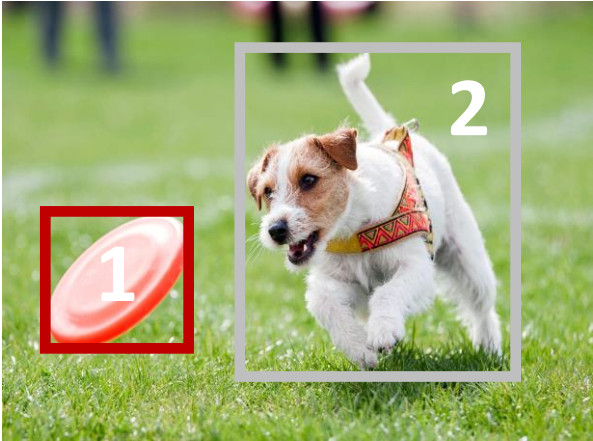
Concept Grounding (Class)



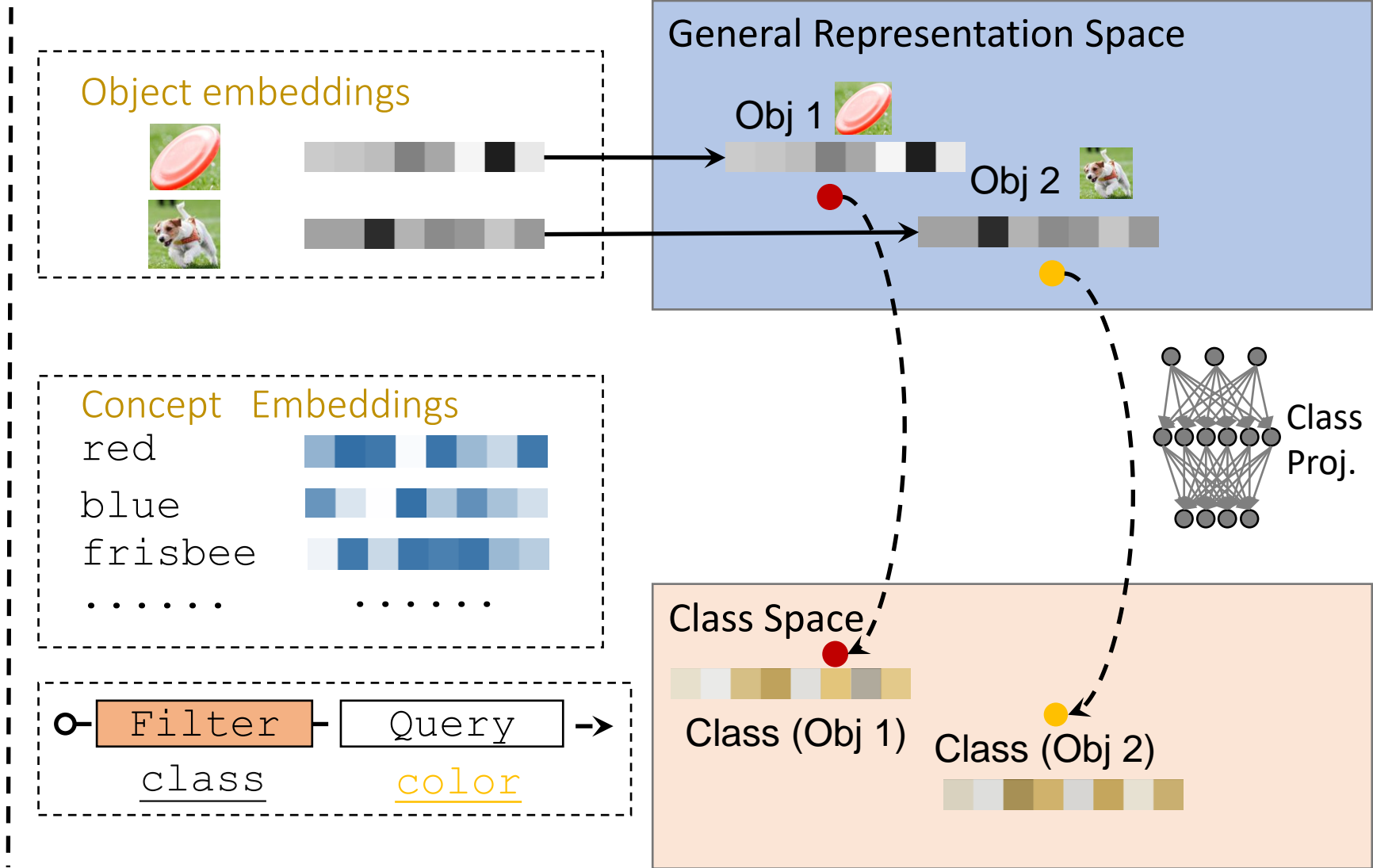
Q: What's the color of the frisbee?



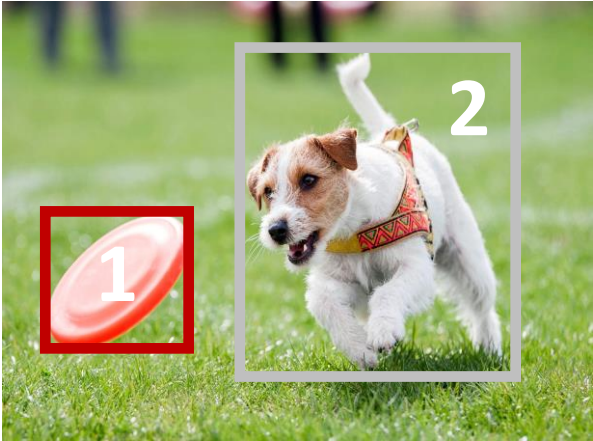
Concept Grounding (Class)



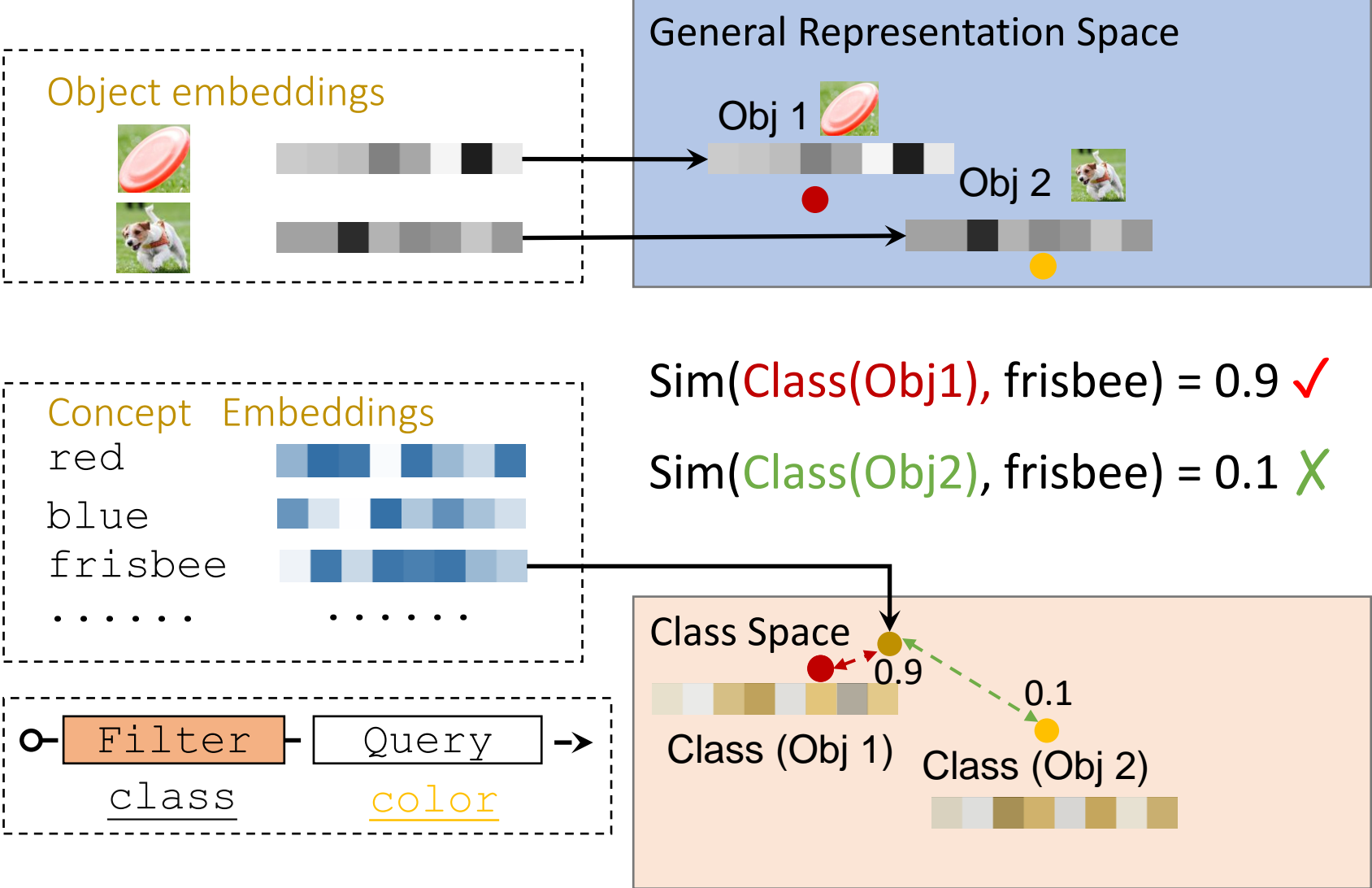
Q: What's the color of the frisbee?



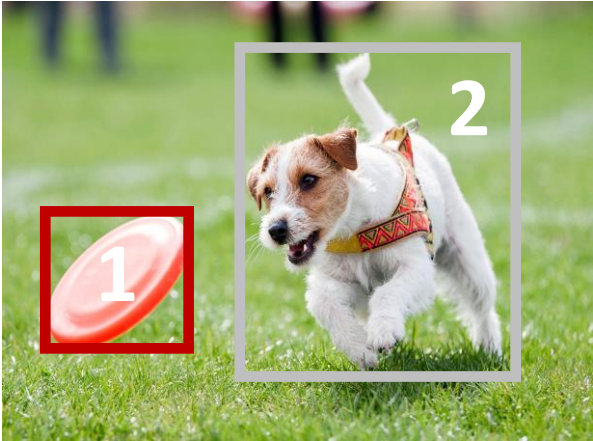
Concept Grounding (Class)



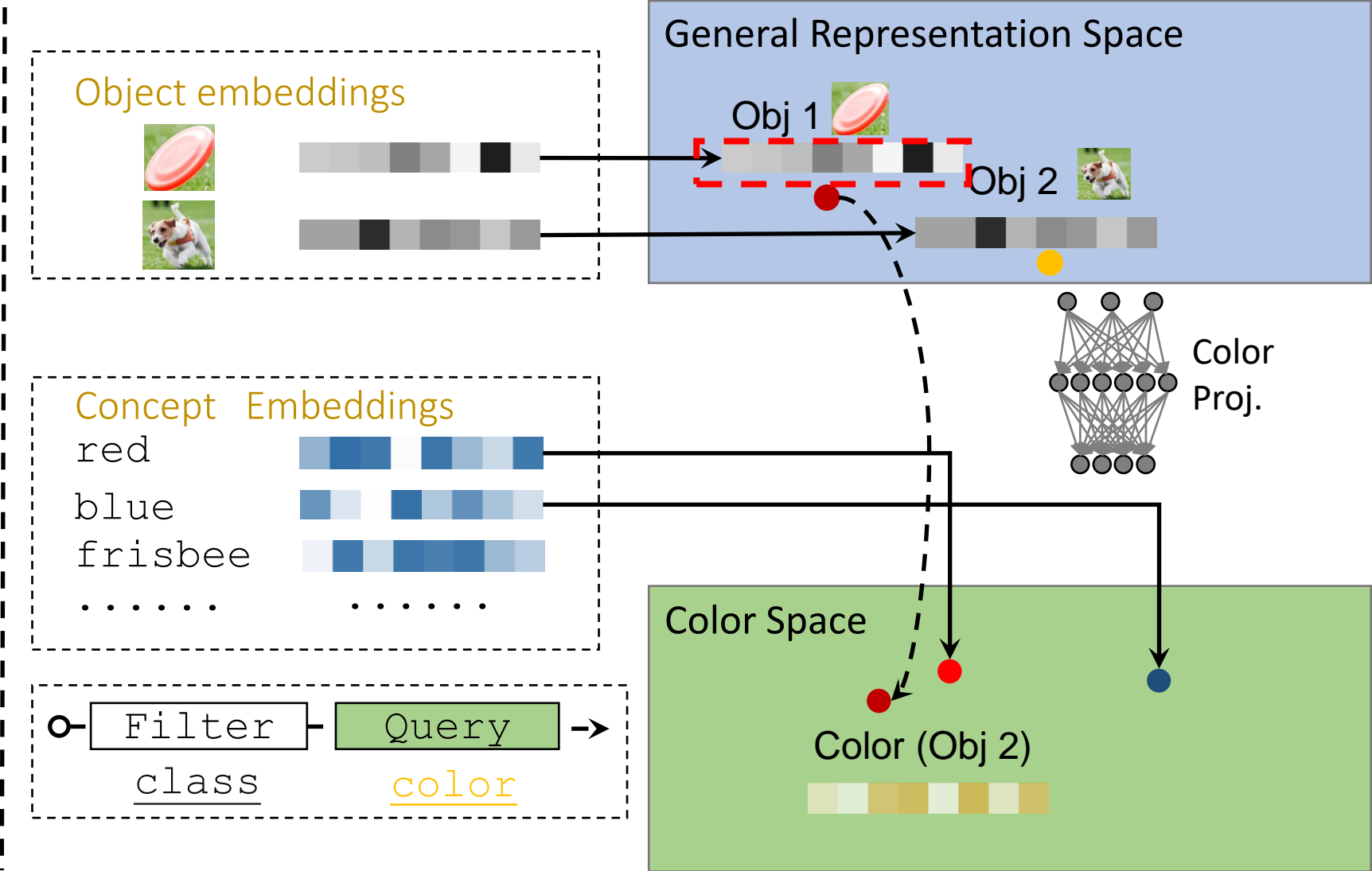
Q: What's the color of the frisbee?



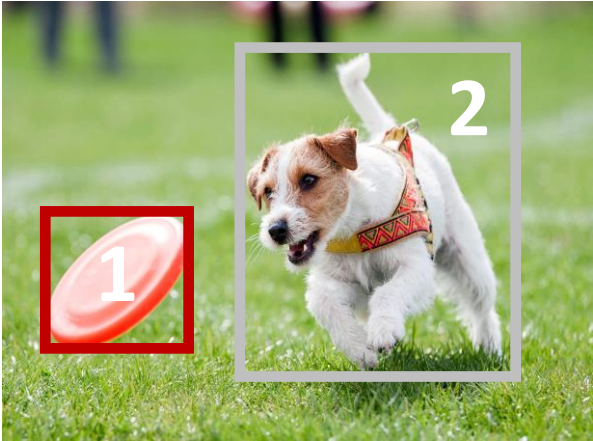
Concept Grounding (Color)



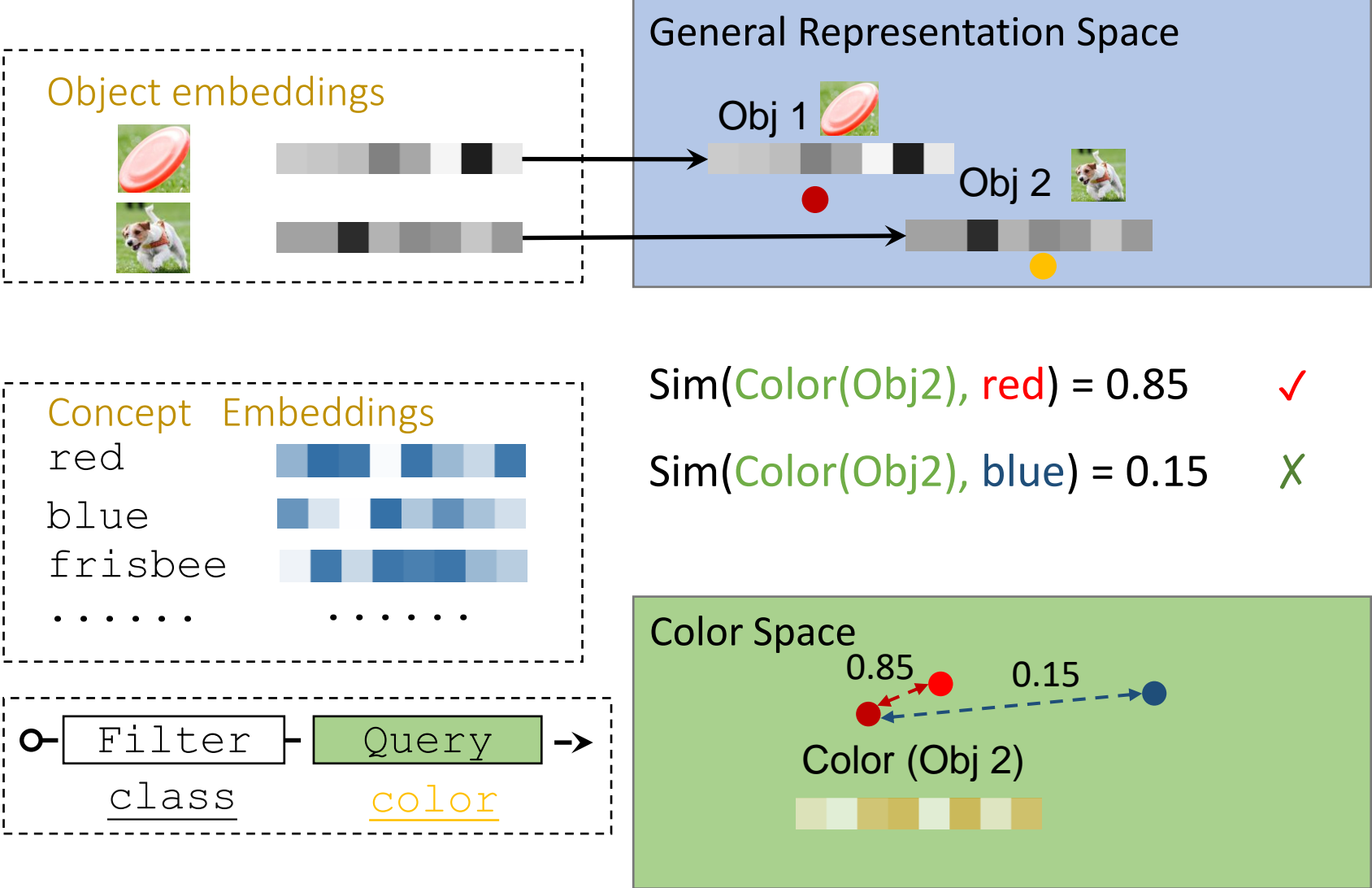
Q: What's the color of the frisbee?



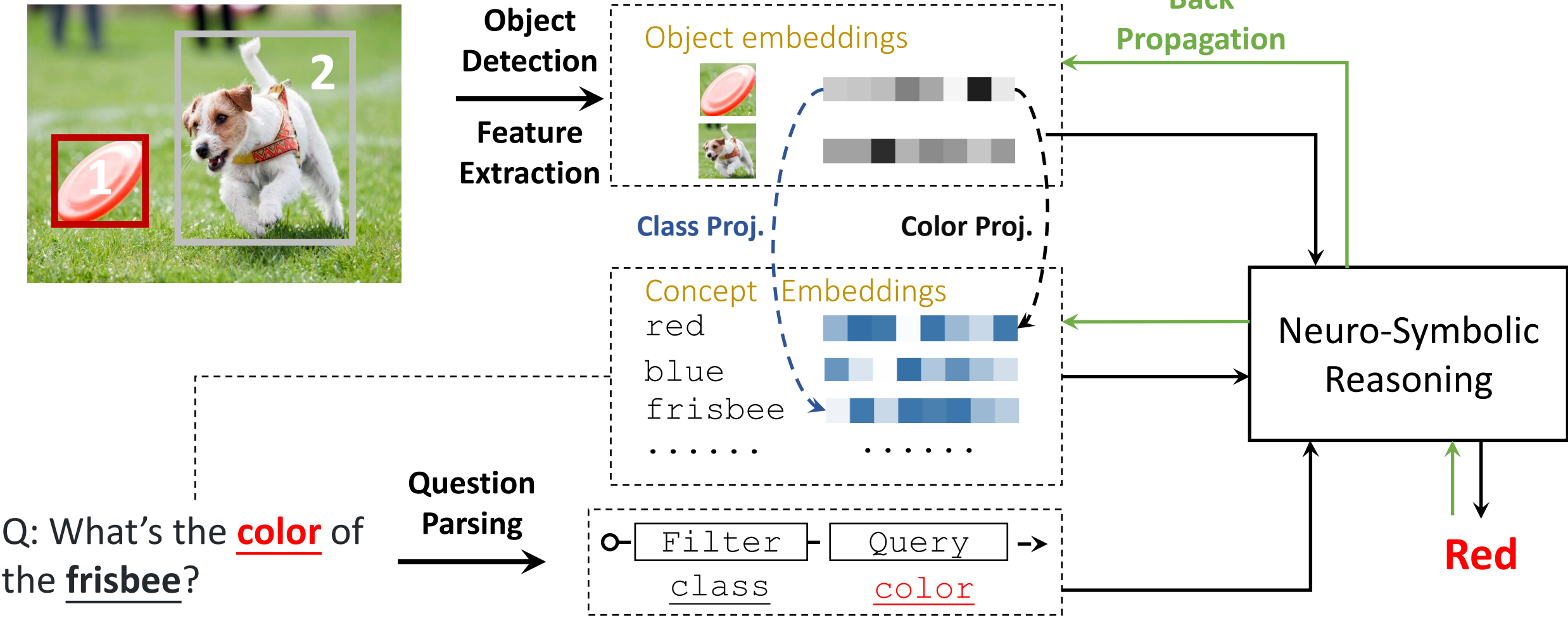
Concept Grounding (Class)



Q: What's the color of the frisbee?

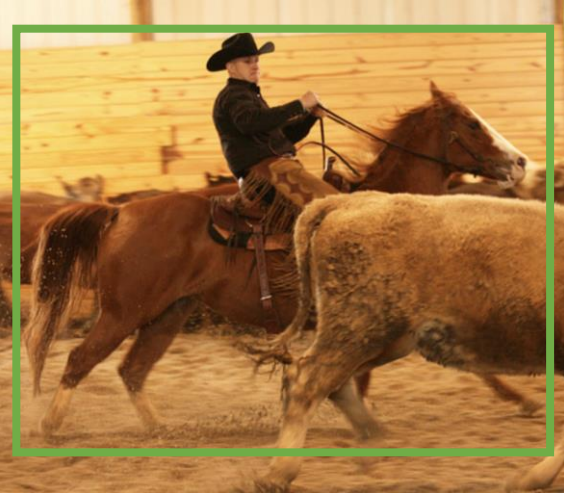


Neuro-symbolic Reasoning

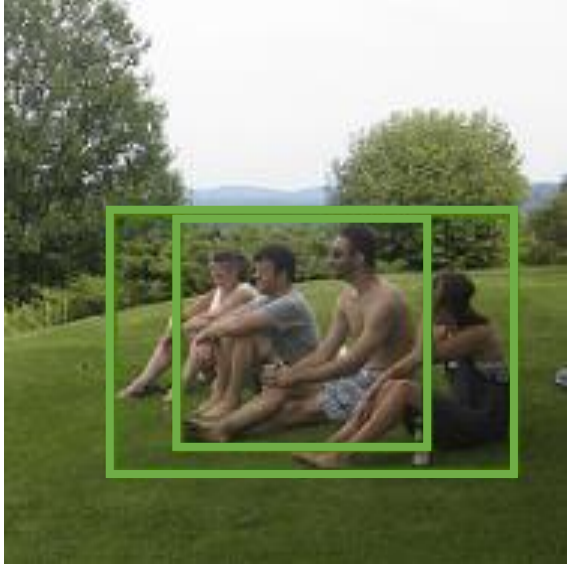


Evaluation Results

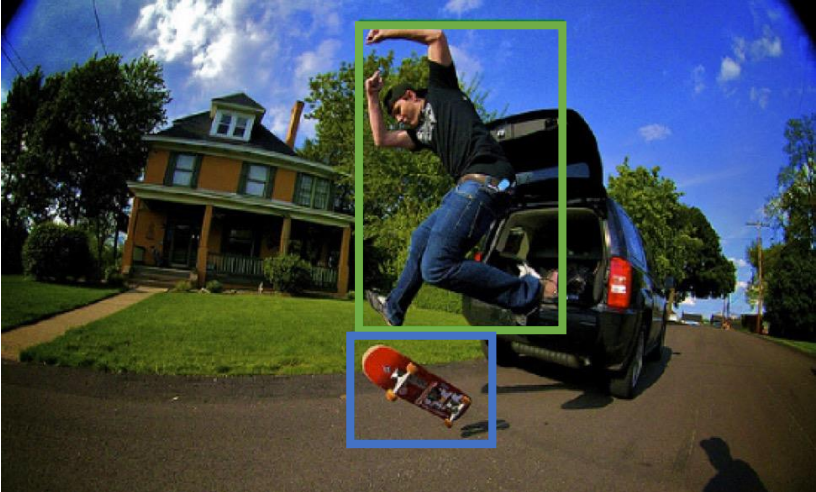
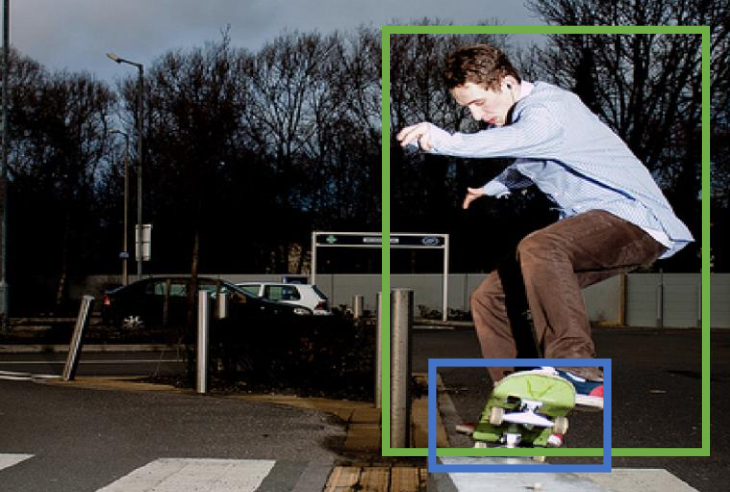
Concept: Horse



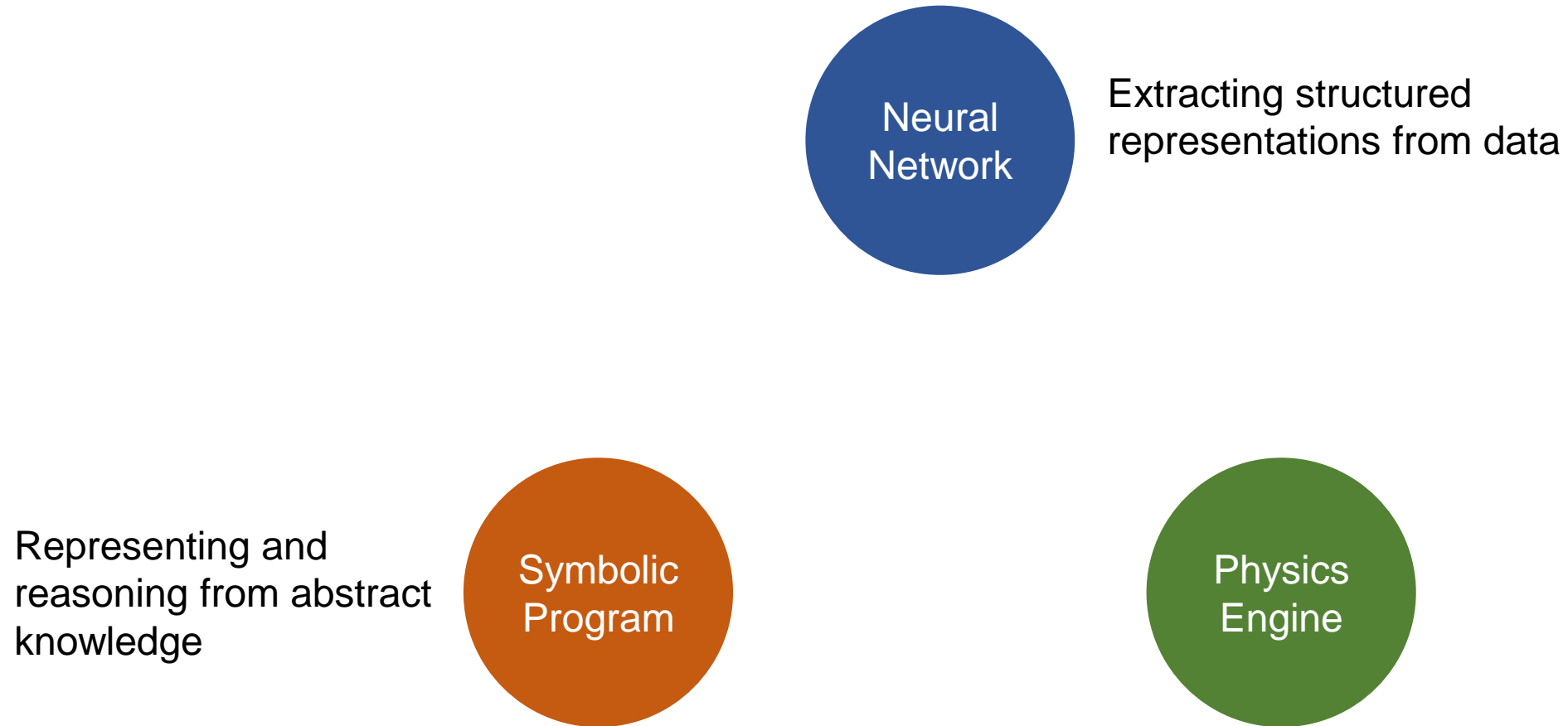
Concept: Person



Concept: Person On a Skateboard



My vision: Neuro-Symbolic AI



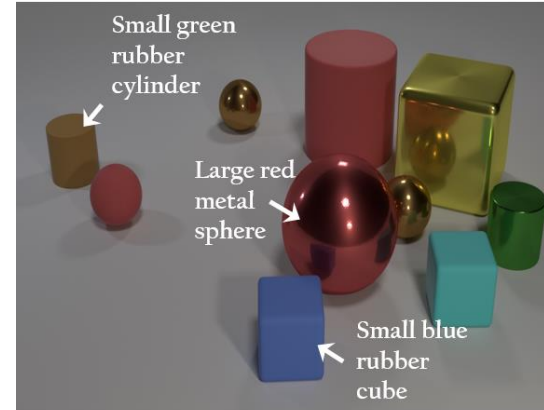
Existing Visual Reasoning Datasets



Q: *What is the mustache made of?*

A: *Banana*

VQA [Antol et al. ICCV 2015]



Q: *Are there an equal number of large things and metal spheres?*

A: *Yes*

CLEVR [Johnson et al. CVPR 2017]



Q: *What does the cat do three times?*

A: *Put head down*

TGIF-QA [Jang et al. CVPR 2017]

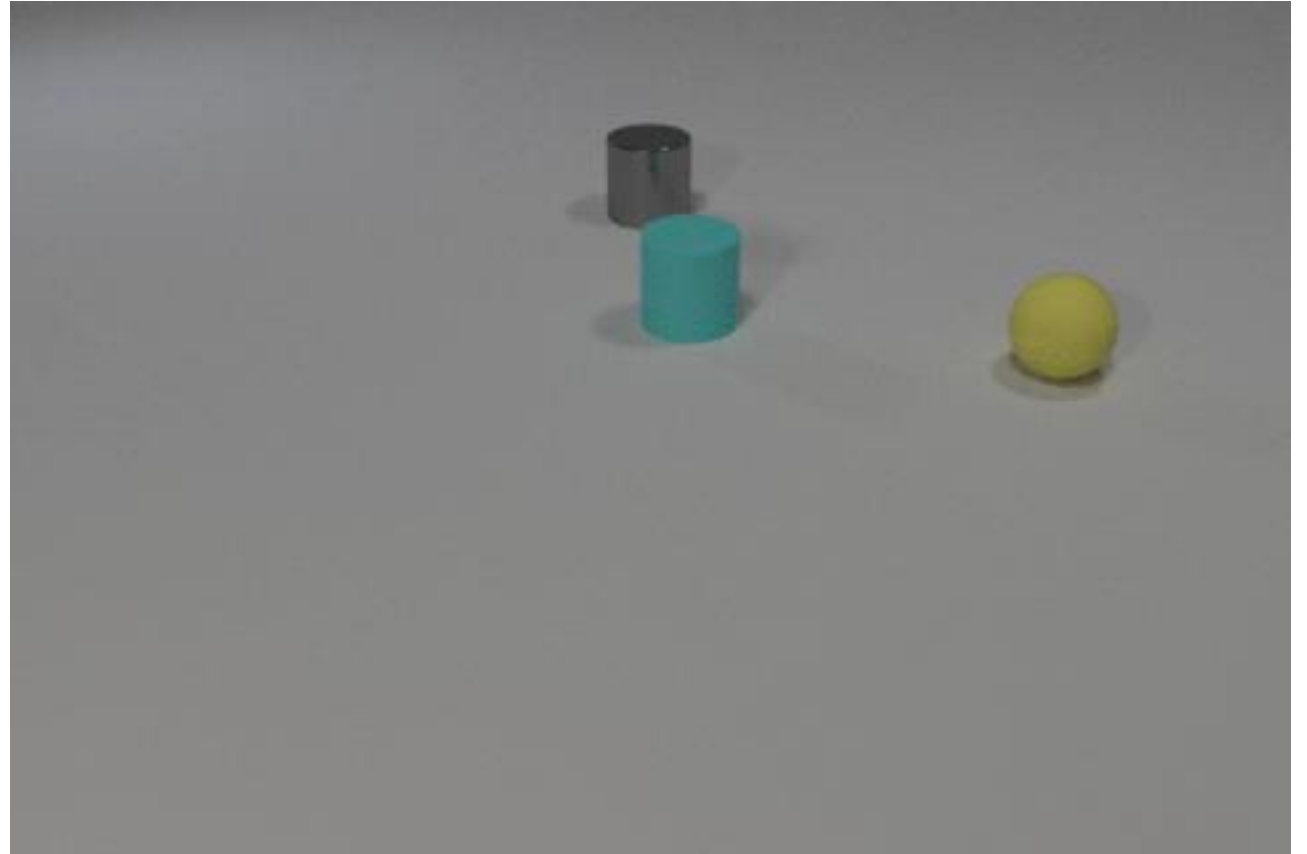
We Also Need to Understand Dynamics



- Describe what has happened
- Explain why it has happened
- Predict what is about to happen

CLEVRER Dataset: From Static Scene to Dynamic Scene

- 20,000 Synthetic videos
- 300,000+ questions
- Controlled biases
- Diagnostic annotations
- Dynamics visual reasoning
 - Descriptive
 - Explanatory
 - Predictive
 - Counterfactual



Question Types

Descriptive

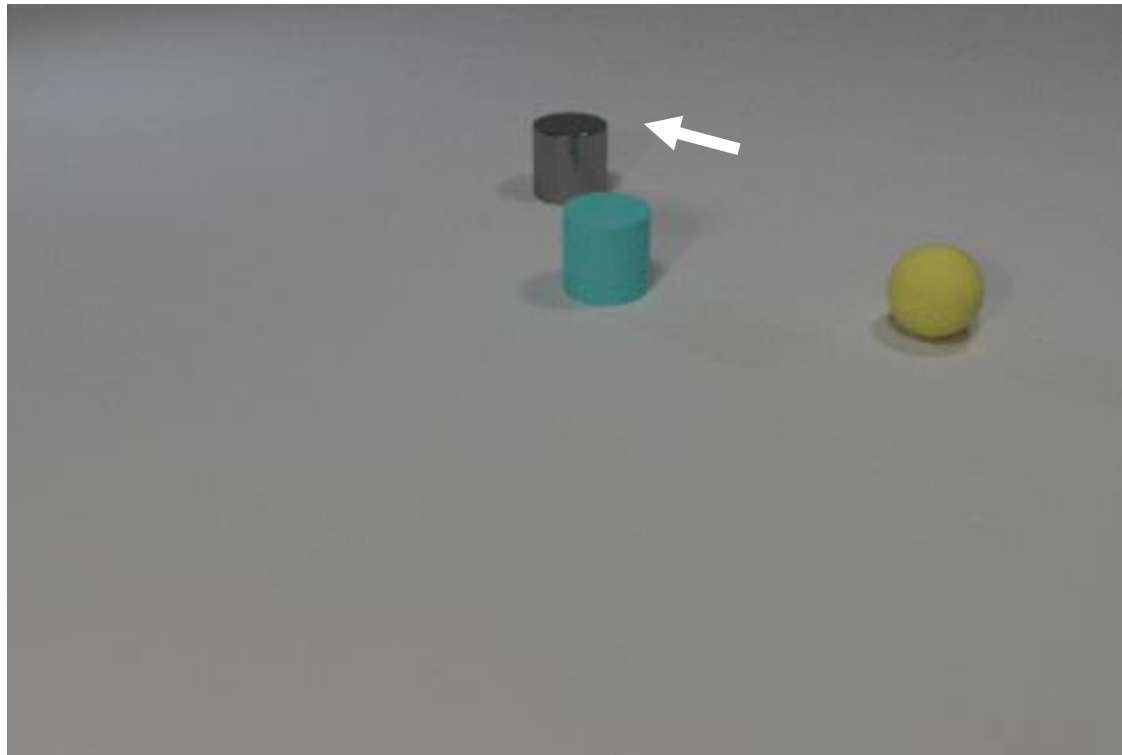
Explanatory

Predictive

Counterfactual

Q: *What is the material of the last object to collide with the cyan cylinder?*

A: *Metal*



Question Types

Descriptive

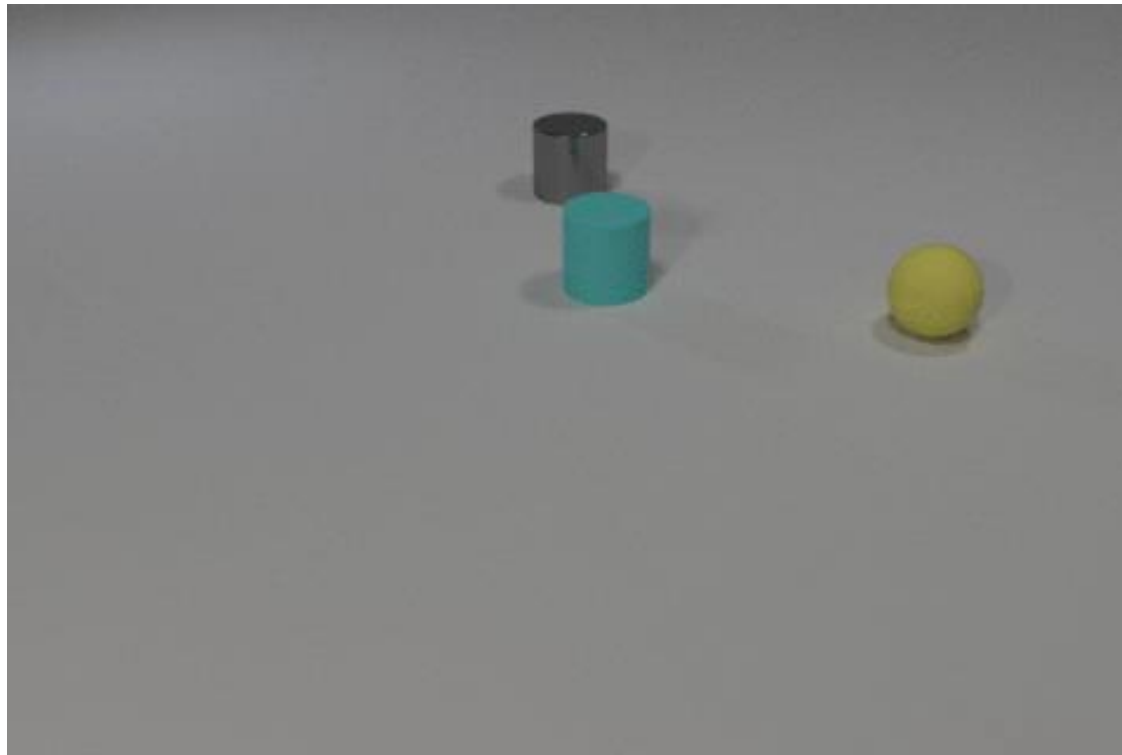
Explanatory

Predictive

Counterfactual

Q: *What is responsible for the collision between the cyan and gray cylinder?*

A: *The collision between the cyan cylinder and the red rubber sphere.*



Question Types

Descriptive

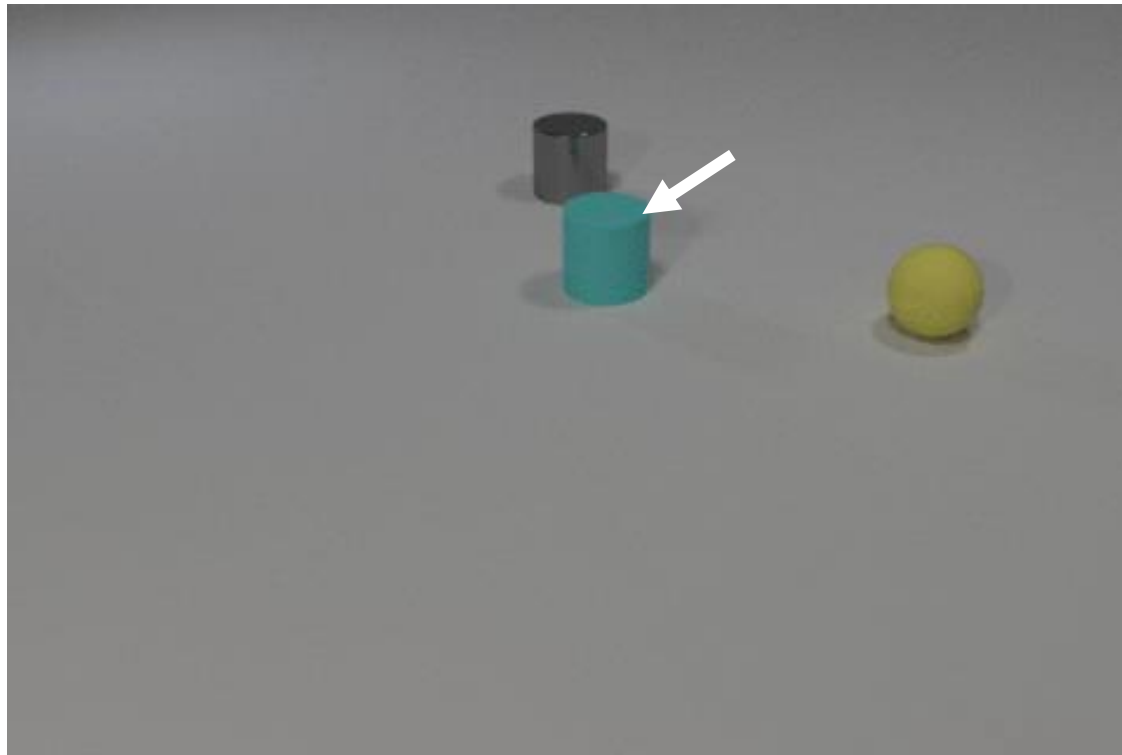
Explanatory

Predictive

Counterfactual

Q: *What will happen next?*

A: *The red rubber sphere collides with the metal sphere.*



Question Types

Descriptive

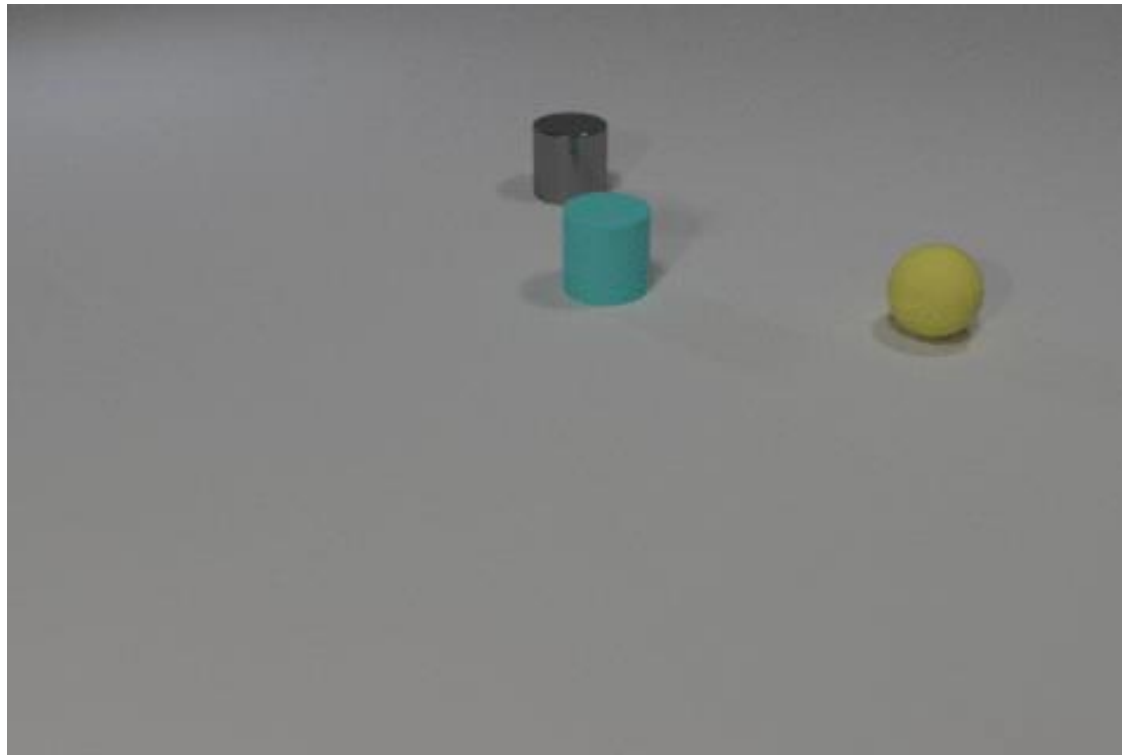
Explanatory

Predictive

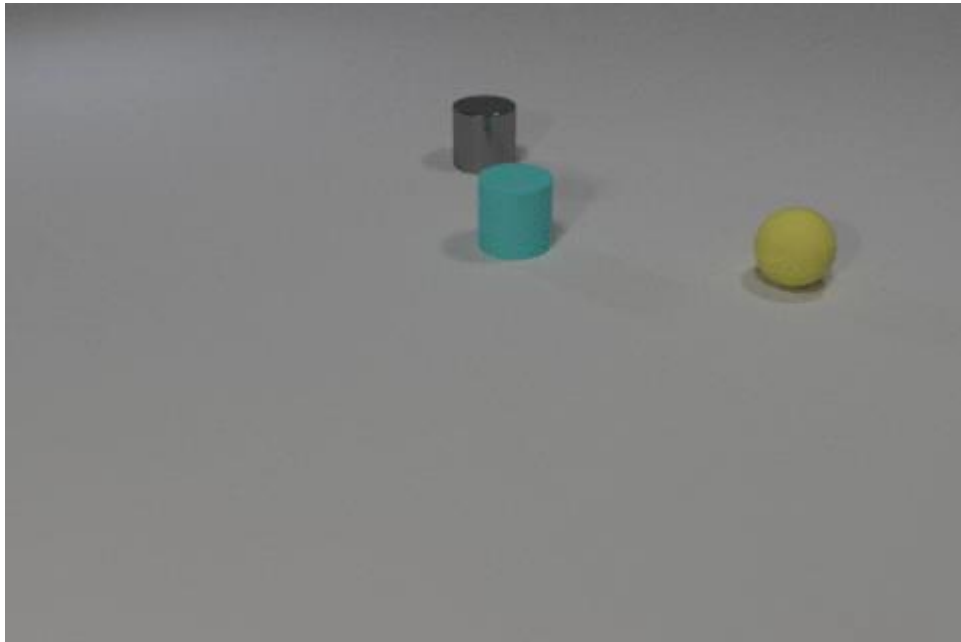
Counterfactual

Q: *What would happen without the **cyan cylinder**?*

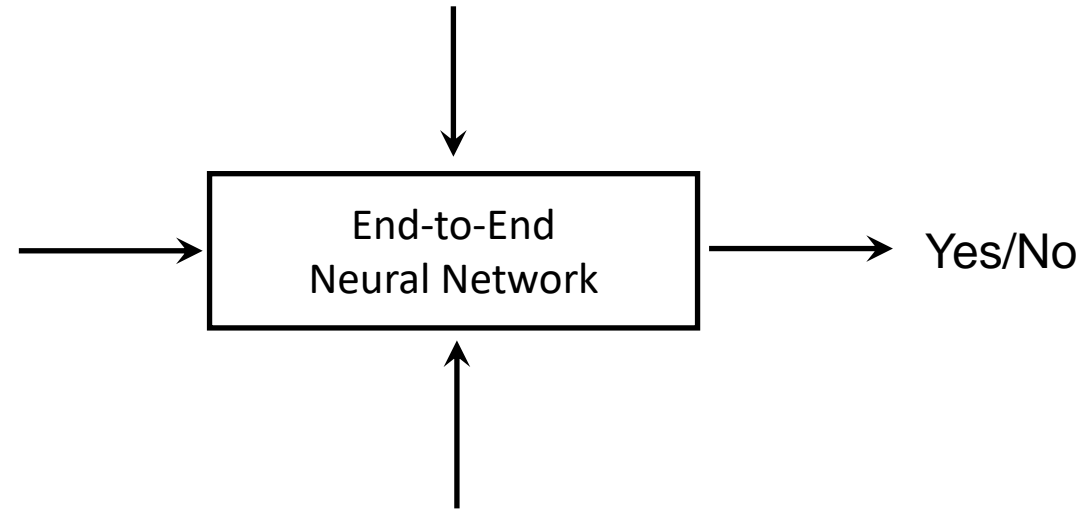
A: *The **red rubber sphere** and the gray object collide.*



Prior Work: End-to-End Video Reasoning

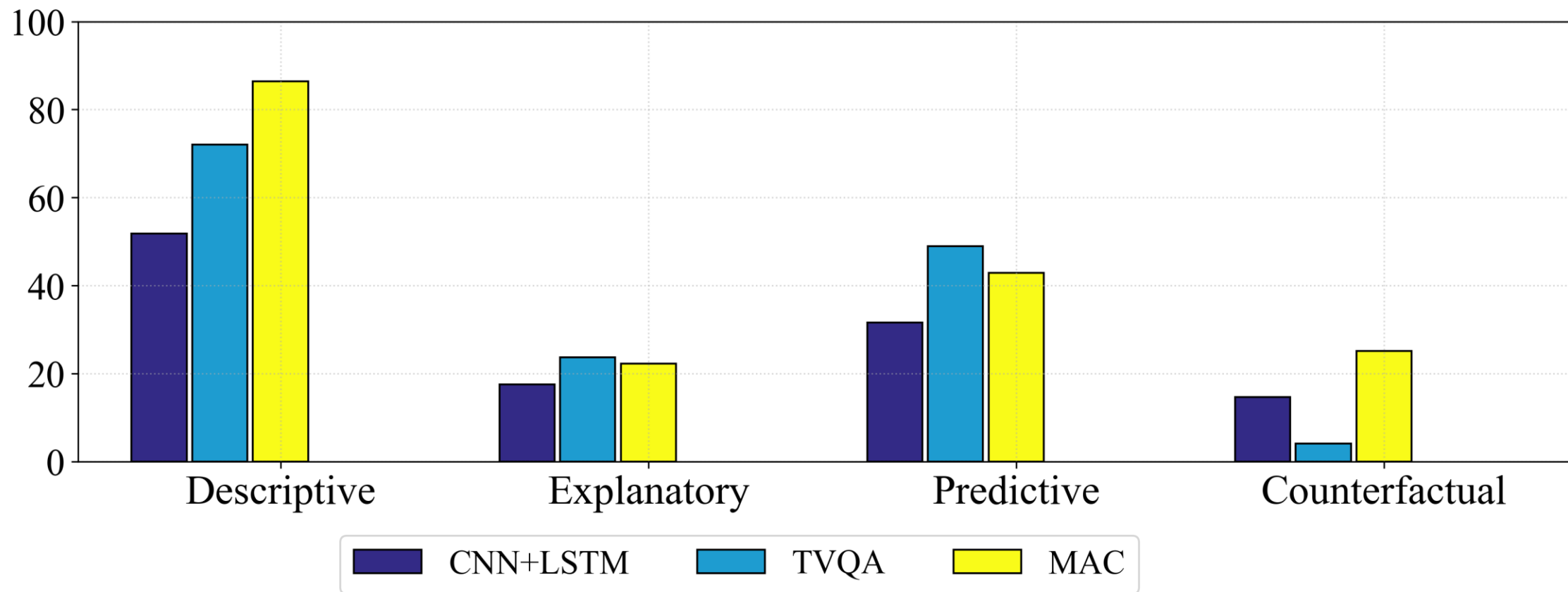


Video Question Answering
Question: *What will happen without the cyan cylinder?*



Candidate Answer: *The red rubber sphere and the gray object collide.*

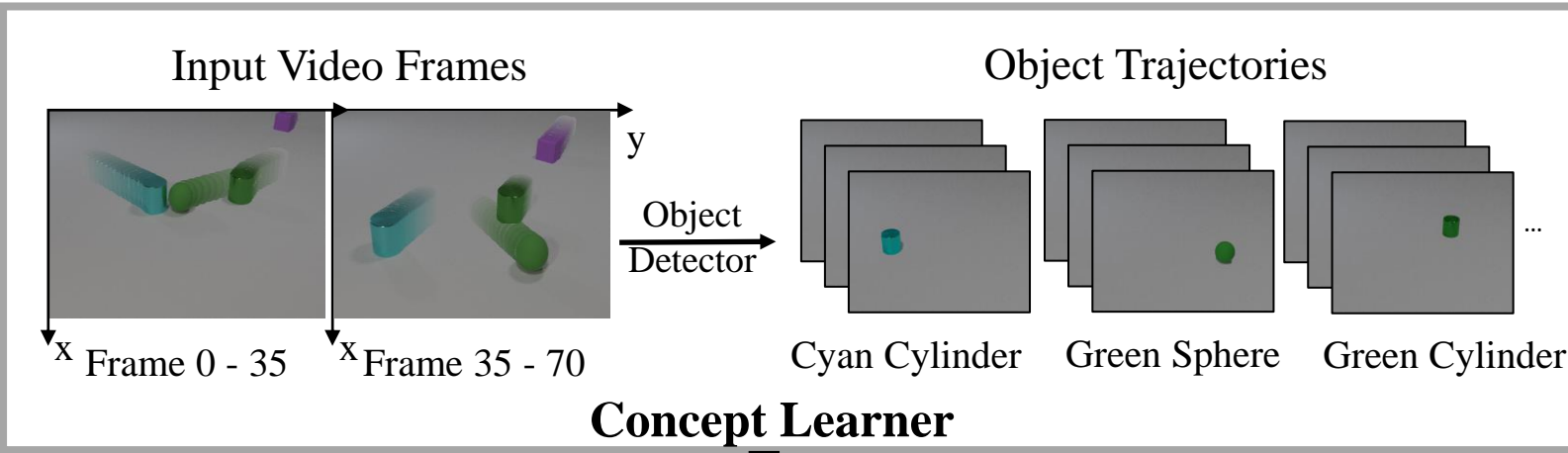
Evaluation of End-to-End Video Reasoning



Physical reasoning requires to understand object dynamics.

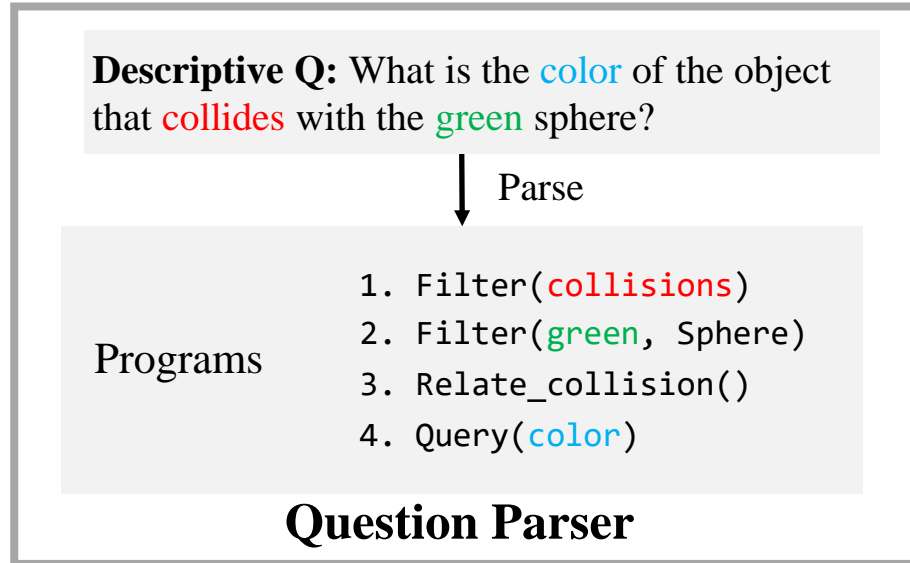
Shi et al. LSTM, 2015. Lei et al. TVQA, 2018. Hudson & Manning. MAC, 2018.

Neuro-symbolic Dynamics Reasoning



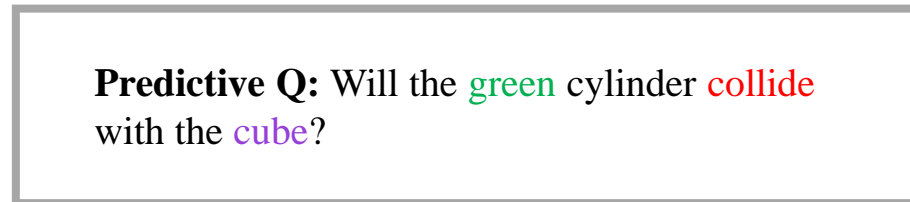
Physical Object Representation

Visual Properties				Collision Events		Physical Parameters			
	Shape	Color	Location	Time		Mass	Restitution	Velocity	...
Obj1	Cylinder	Cyan	(-2.3, -1.7)	1.3s	Obj1	?	?	?	
Obj2	Sphere	Green	(-1.2, 1.1)		Obj2	?	?	?	
Obj3	Cylinder	Green	(-2.2, 2.3)		Obj3	?	?	?	
Obj4	Cube	Purple	(-6.7, 4.1)		Obj4	?	?	?	



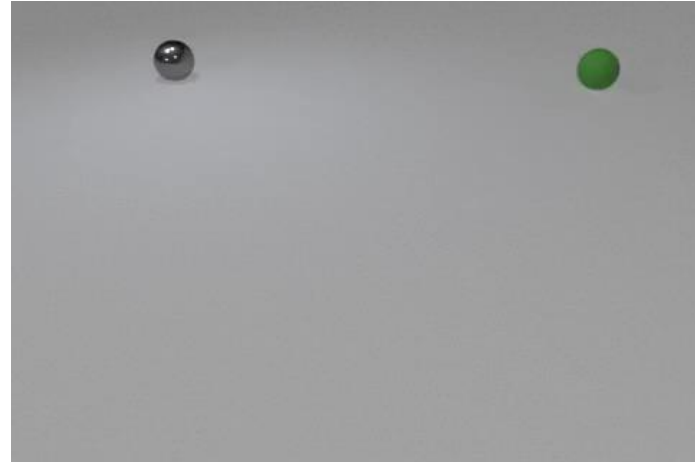
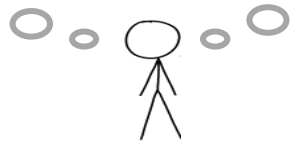
↓

Answer: **Cyan**

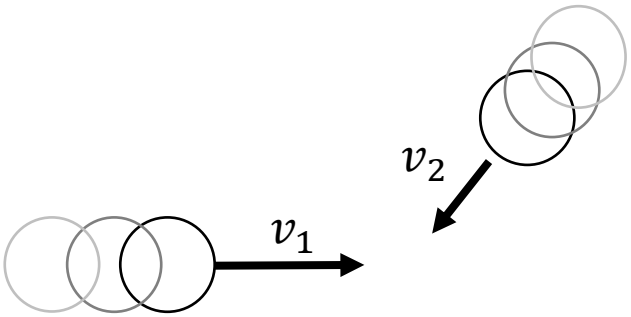


Parameters Estimation from Collision Events

Which object is heavier, the gray sphere or the green sphere?



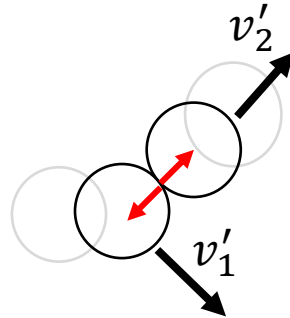
Forward Simulation of Sphere Collisions



	Shape	Location (x)	Velocity (v)	...
Obj1	Sphere	(-1.0,0.0)	(1.0, 0.0)	
Obj2	Sphere	(1.5, 1.5)	(-0.1, -0.1)	
...				

$$x_1^t = x_1^{t-1} + v_1 dt$$

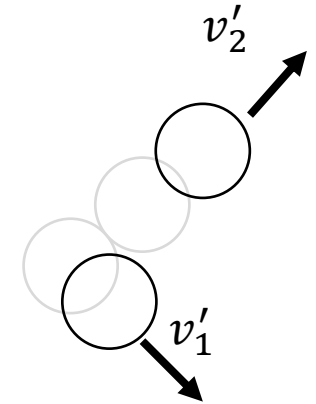
$$x_2^t = x_2^{t-1} + v_2 dt$$



	Shape	Location (x)	Velocity (v)	...
Obj1	Sphere	(0.0,0.0)	(0.5, -0.5)	
Obj2	Sphere	(1.4, 1.4)	(0.6, 0.6)	
...				

$$m_1 v_1 + m_2 v_2 = m_1 v'_1 + m_2 v'_2$$

Coefficient of Restitution: $e = \frac{v'_2 - v'_1}{v_2 - v_1}$



	Shape	Location (x)	Velocity (v)	...
Obj1	Sphere	(0.5, -0.5)	(0.5, -0.5)	
Obj2	Sphere	(2.0, 2.0)	(0.6, 0.6)	
...				

$$x_1^{t+1} = x_1^t + v'_1 dt$$

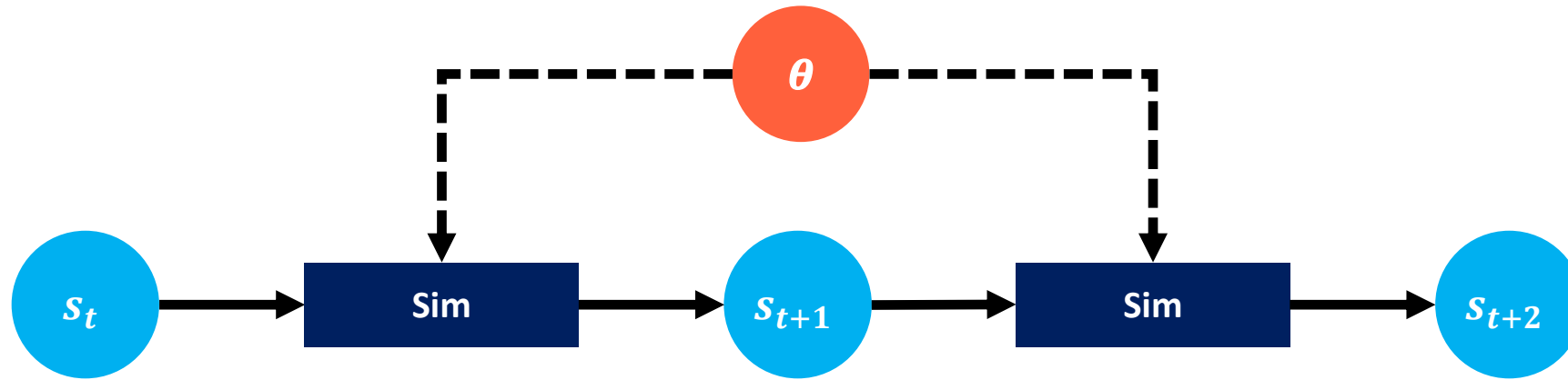
$$x_2^{t+1} = x_2^t + v'_2 dt$$

Physical parameters to be optimized: m, e

	Mass (m)	Restitution (e)	...
Obj1	0.94	0.78	
Obj2	0.76	0.95	
...			

Forward Simulations Using Differentiable Physics

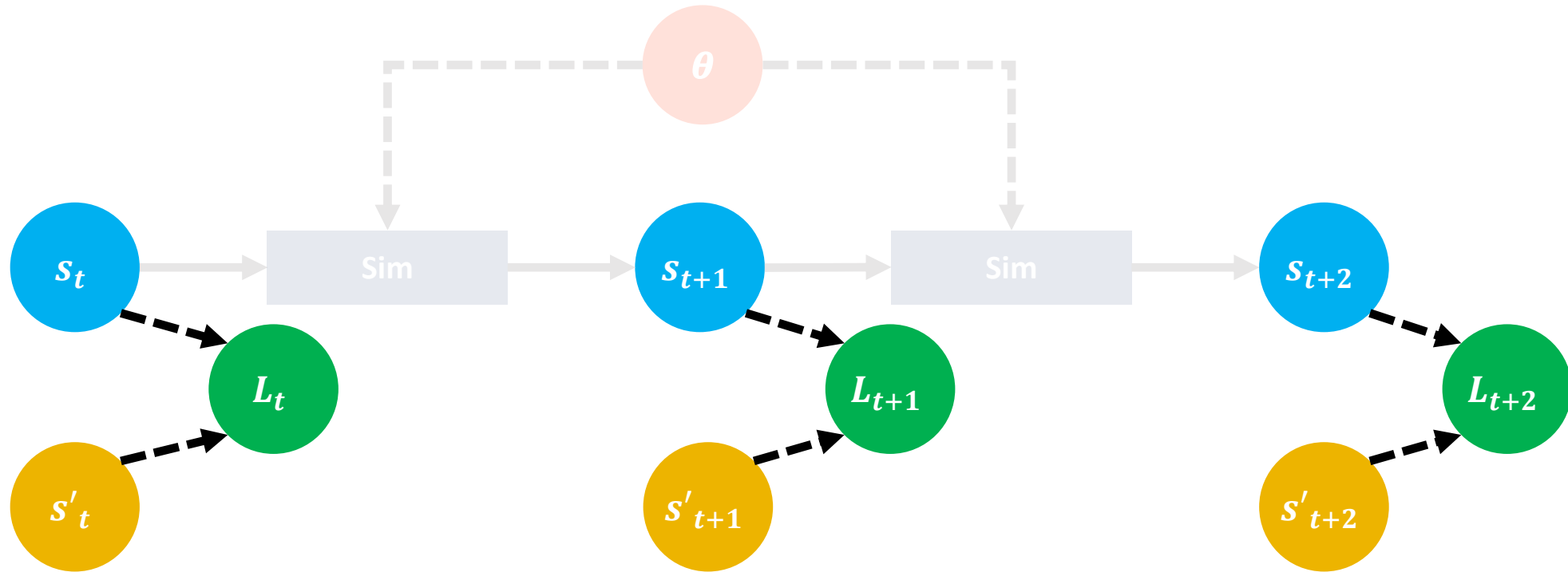
- Forward simulation



- Physical parameters θ : e.g., mass and restitution of rigid bodies
- State s : e.g., positions and velocities of rigid bodies

Parameters Estimation with Differentiable Physics

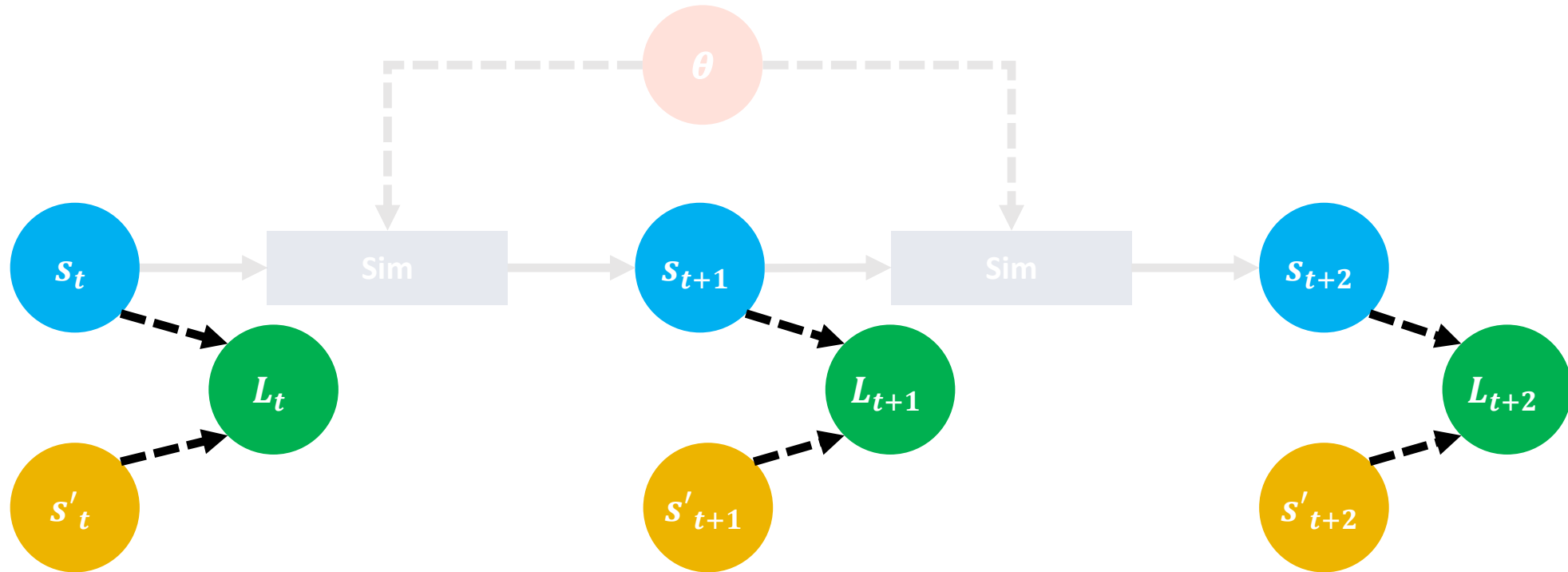
- Loss computation



- Physical parameters θ : e.g., mass and restitution of rigid bodies
- State s : e.g., positions and velocities of rigid bodies
- Observation s' : e.g., raw observations (positions) of rigid bodies
- Loss L : the distance between the state and the visual observation

Parameters Estimation with Differentiable Physics

- Minimize the loss by inferring physical parameter θ

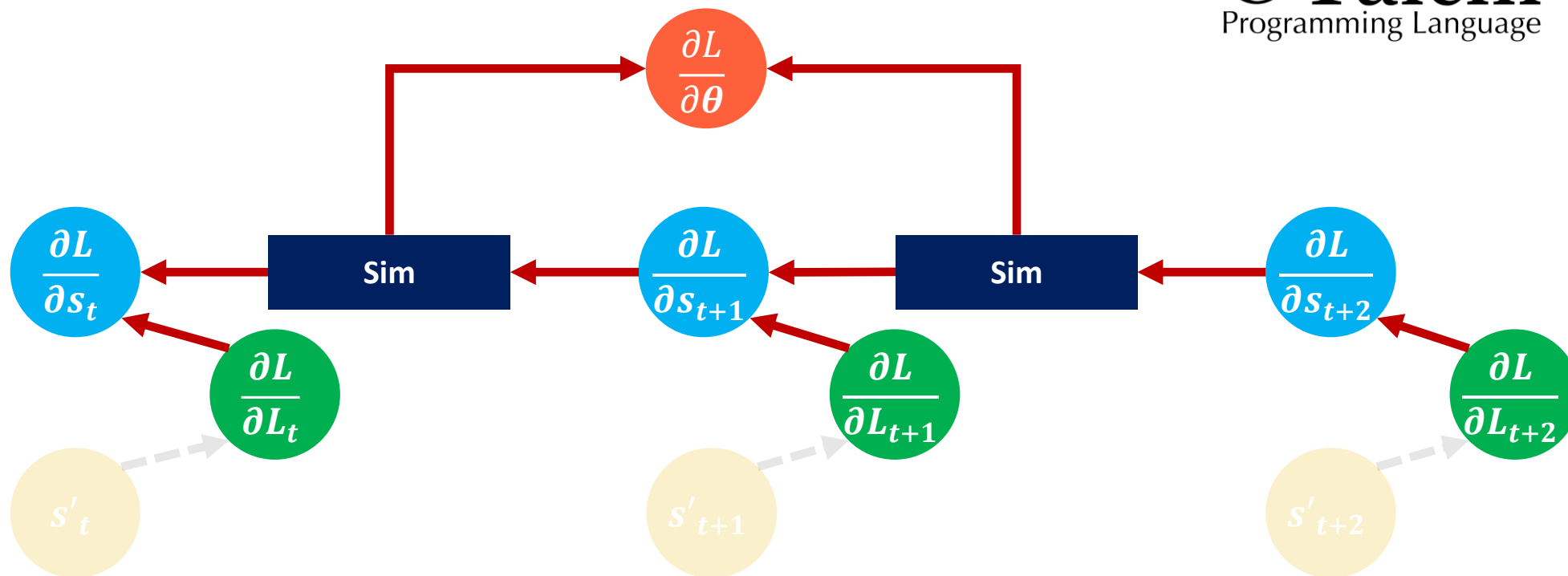


$$L = L_1 + L_2 + \dots + L_T$$

- Physical parameters θ : e.g., mass and restitution of rigid bodies
- State s : e.g., positions and velocities of rigid bodies
- Observation s' : e.g., raw observations (positions) of rigid bodies
- Loss L : the errors between the state and the visual observation

Parameters Estimation with Differentiable Physics

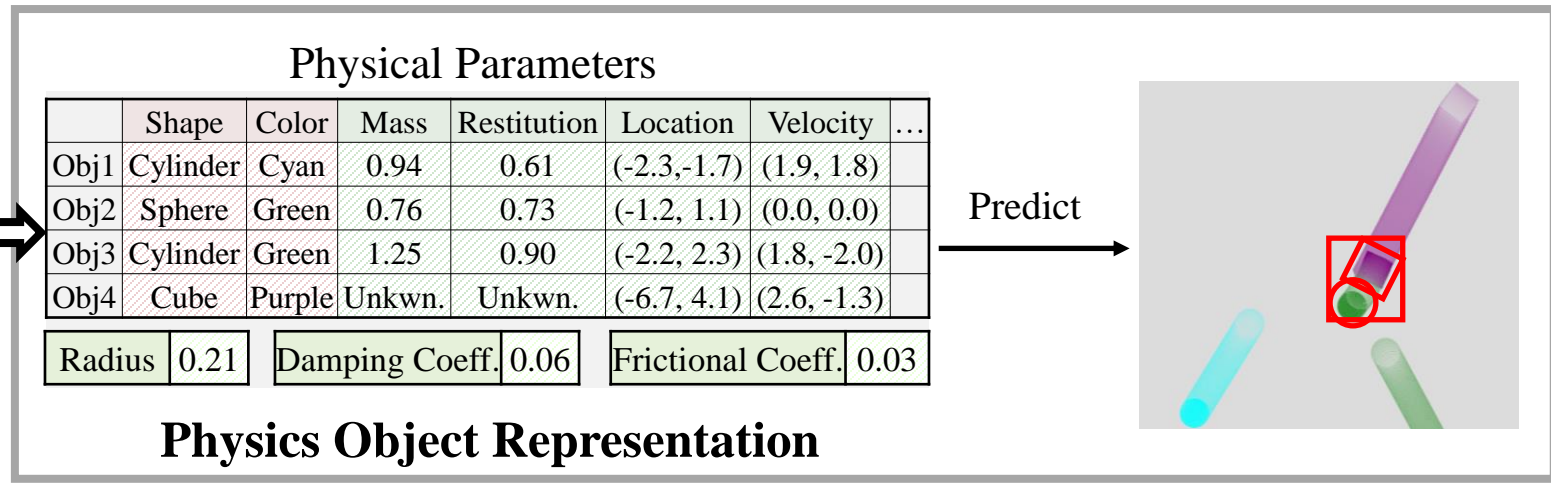
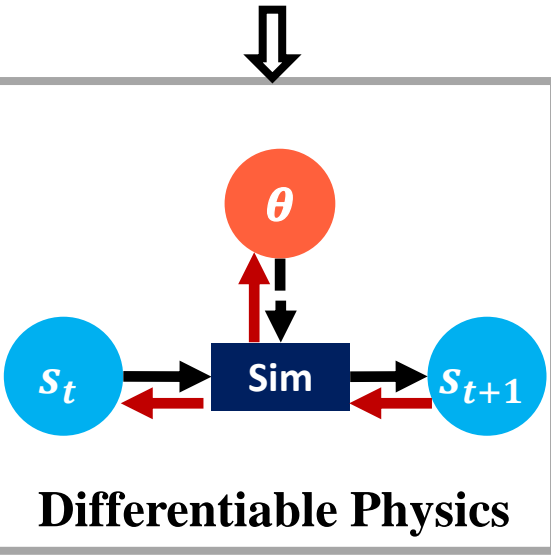
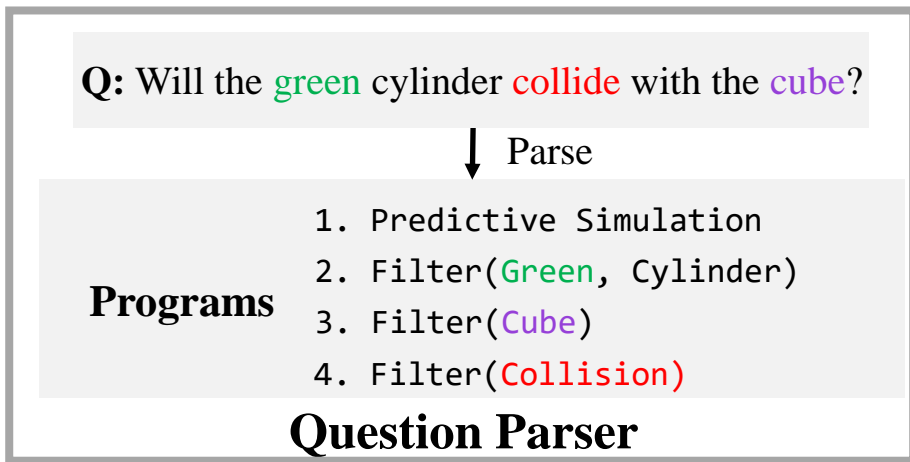
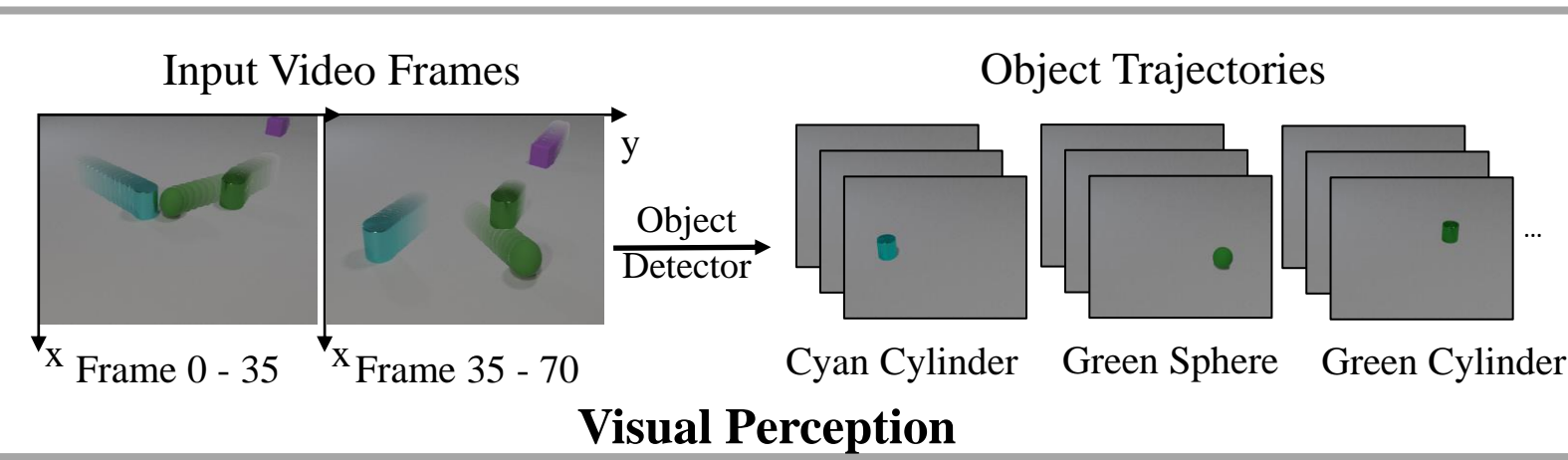
- Backward



$$L = L_1 + L_2 + \dots + L_T$$

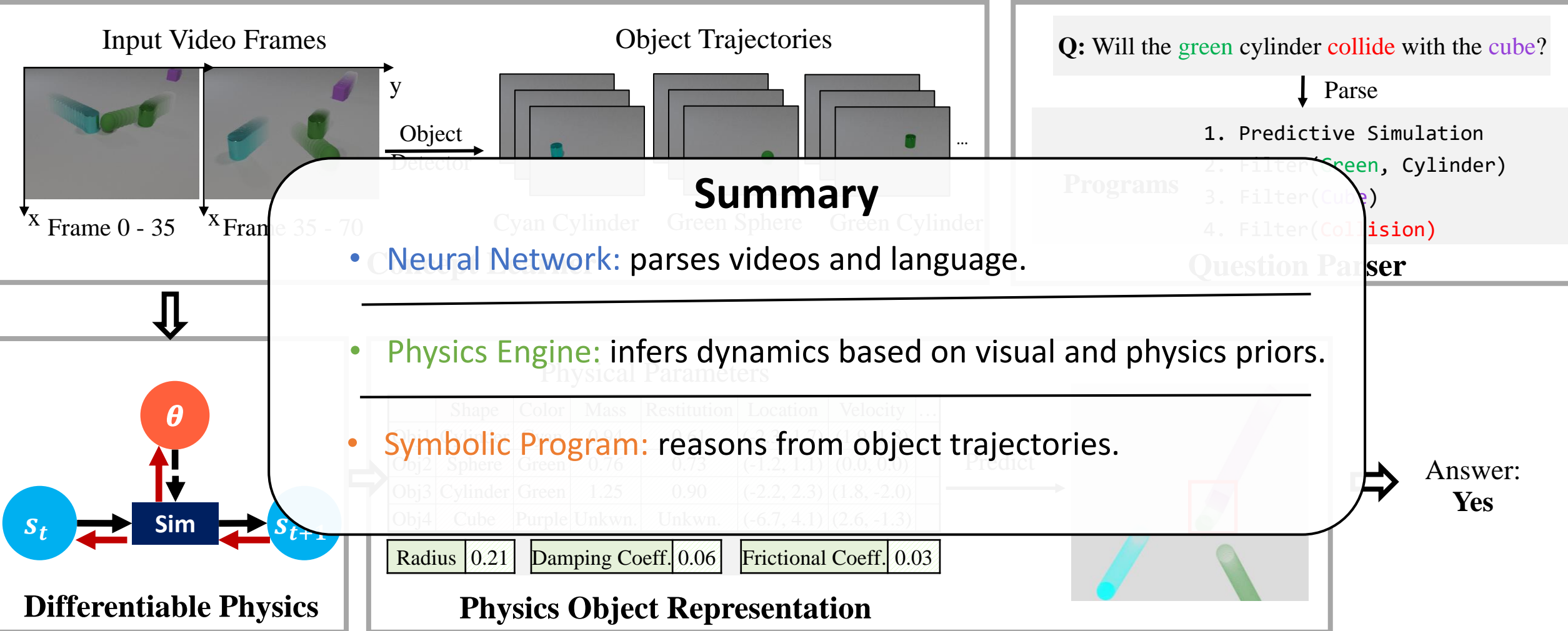
- Physical parameters θ : e.g., mass and restitution of rigid bodies
- State s : e.g., positions and velocities of rigid bodies
- Observation s' : e.g., raw observations (positions) of rigid bodies
- Loss L : the error between the state and the visual observation

Neuro-symbolic Dynamics Reasoning *with* physics



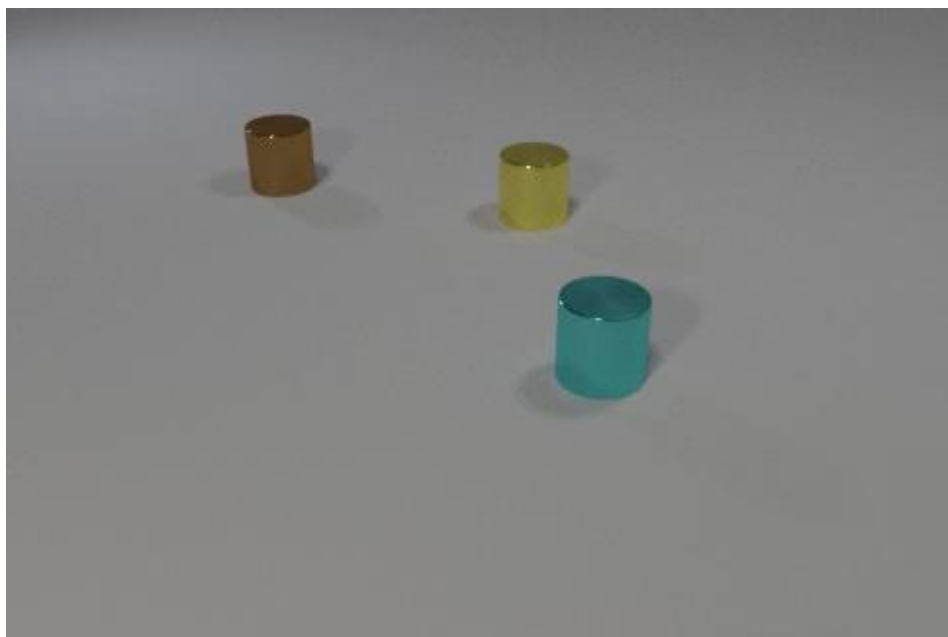
⇒ Answer: Yes

Neuro-symbolic Dynamics Reasoning *with* physics

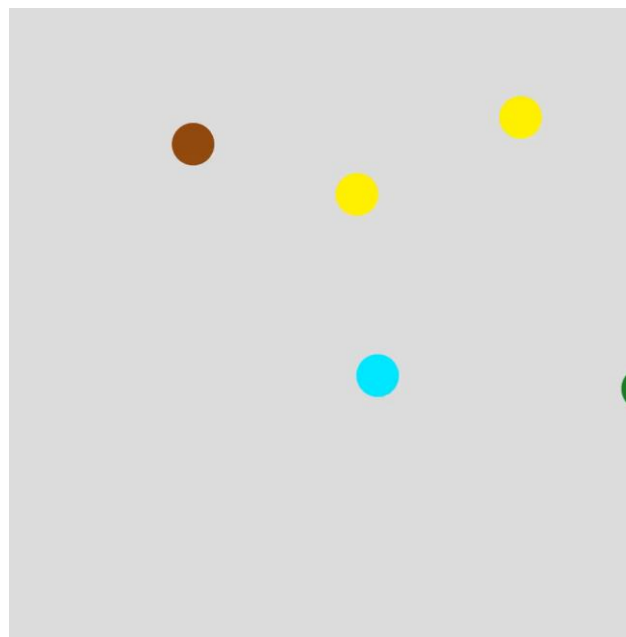


Visualization of the Rconstructed Physics Model

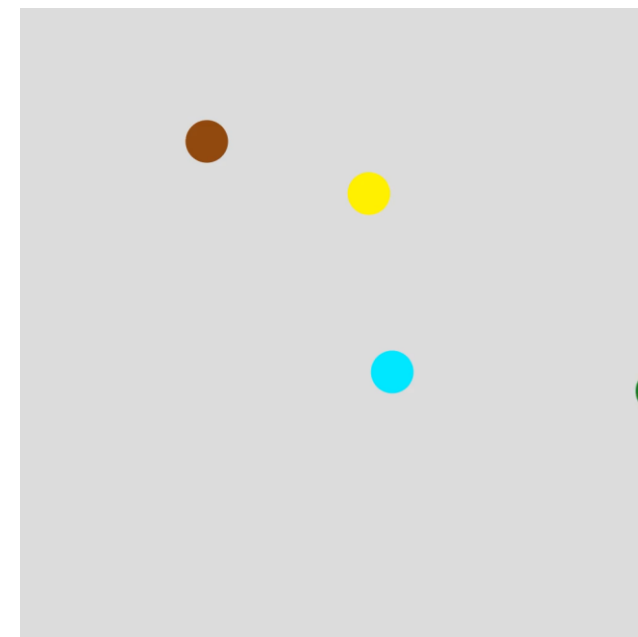
- Results: Fitting the optimized simulation to video observations



Video observations



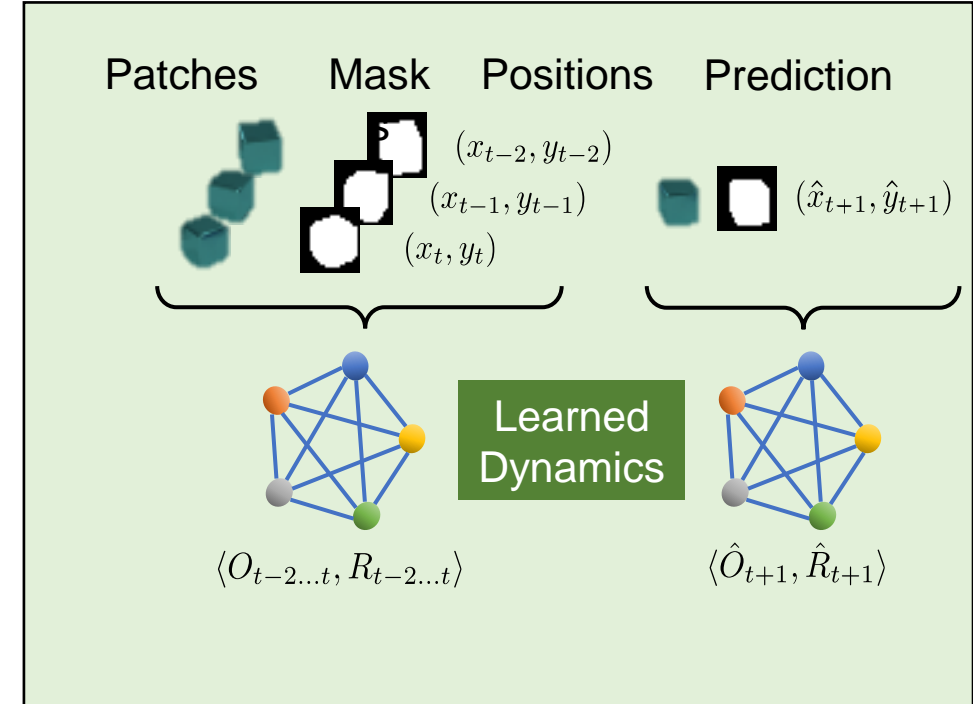
Random initialization



Reconstructed scene

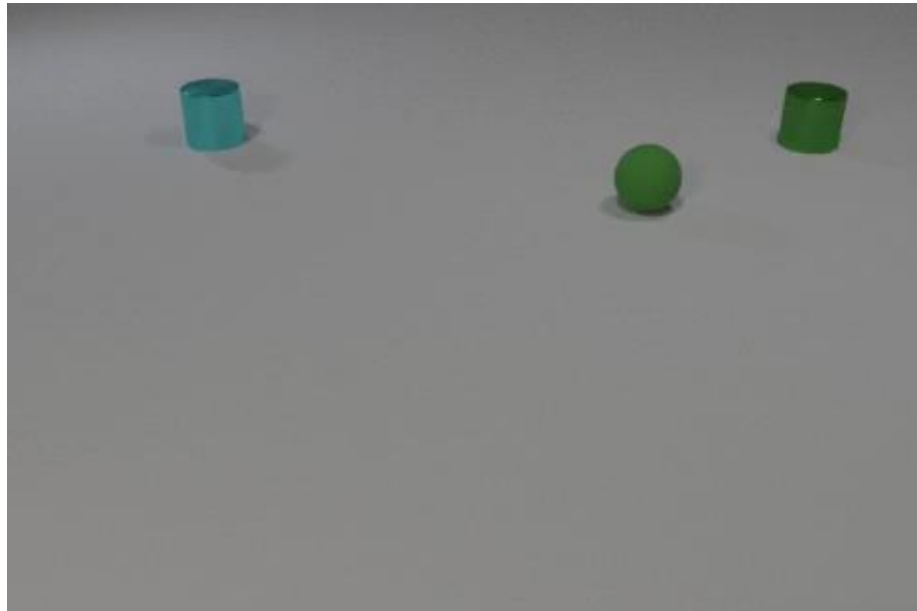
Why Not Learned dynamics Models?

- Input:
 - Nodes: object patches and positions over a time window
 - Edges: collision labels between two objects
- Output:
 - Object positions at the next step
 - Collision labels between two objects



Counterfactual Dynamics Rollout

- Remove the green sphere
 - Q: Would the green cylinder collide with the cyan object if the green sphere is removed?
 - A: True



Video observations

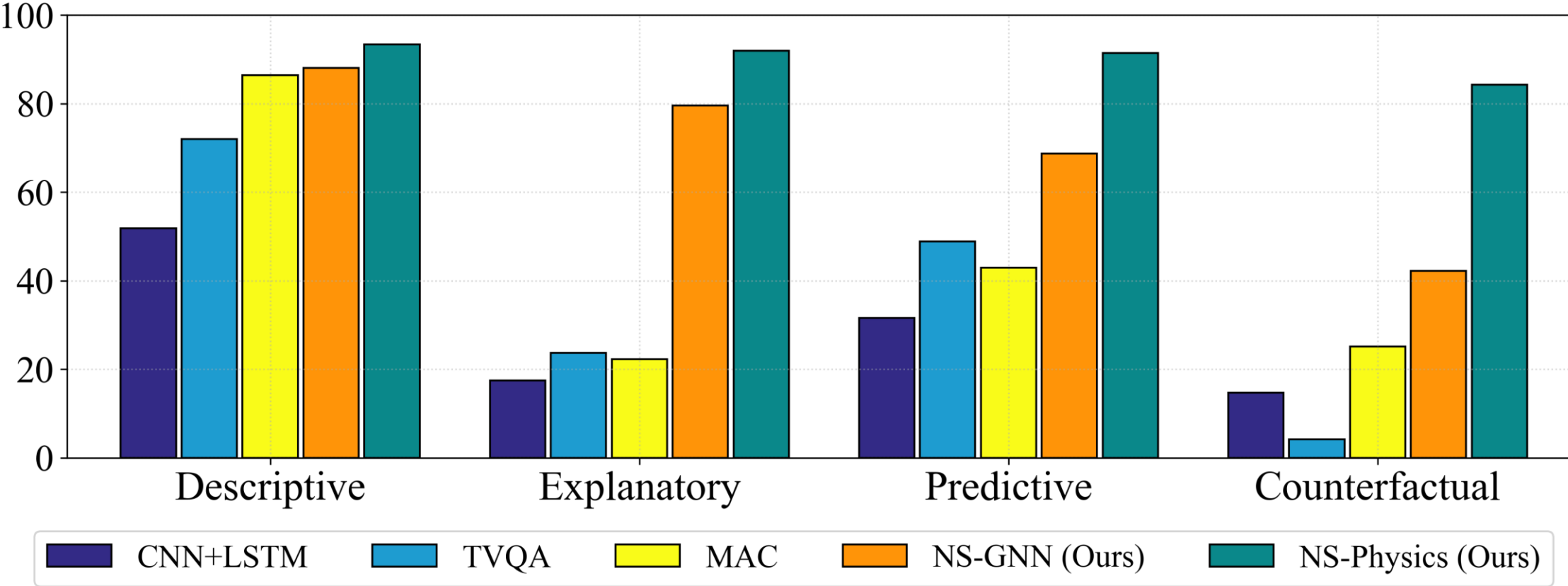


GNNs (✗)



Physics Model (✓)

CLEVRER Test Accuracy



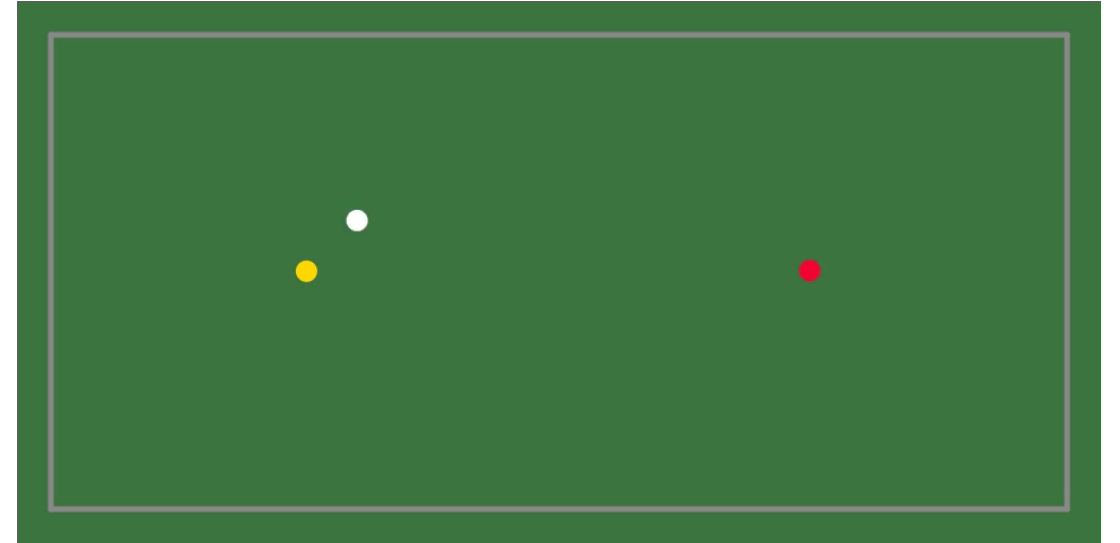
Shi et al. LSTM, 2015. Lei et al. TVQA, 2018. Hudson & Manning. MAC, 2018.
Yi et al. NS-GNN, 2020. Ding et al. NS-Physics, 2021.

Dynamics Visual Reasoning on Real-Billiard

- Estimate the physics models of billiards from **a single video**.



Real world video observations



Predictive Simulation

Application: Invert Physics for Real-world Quadrotor

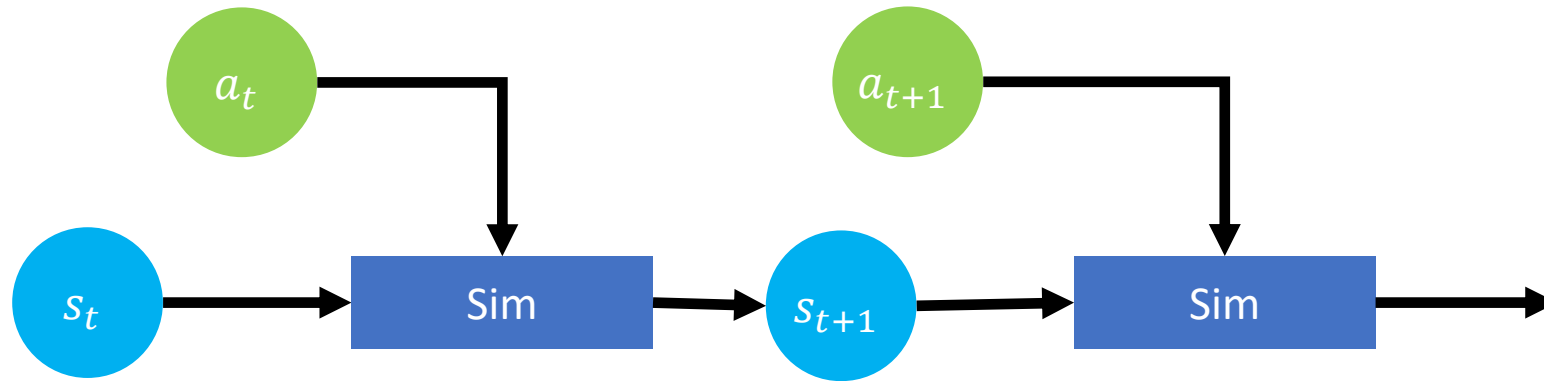
- Imitate the trajectory in a real-world video clip where the **underlying physics are unknown**.



RISP: Rendering-Invariant State Predictor with Differentiable Simulation and Rendering for Cross-Domain Parameter Estimation, Ma*, Du*, Matusik, Tenenbaum, **Gan**. ICLR'22

Differentiable Physics for Planning

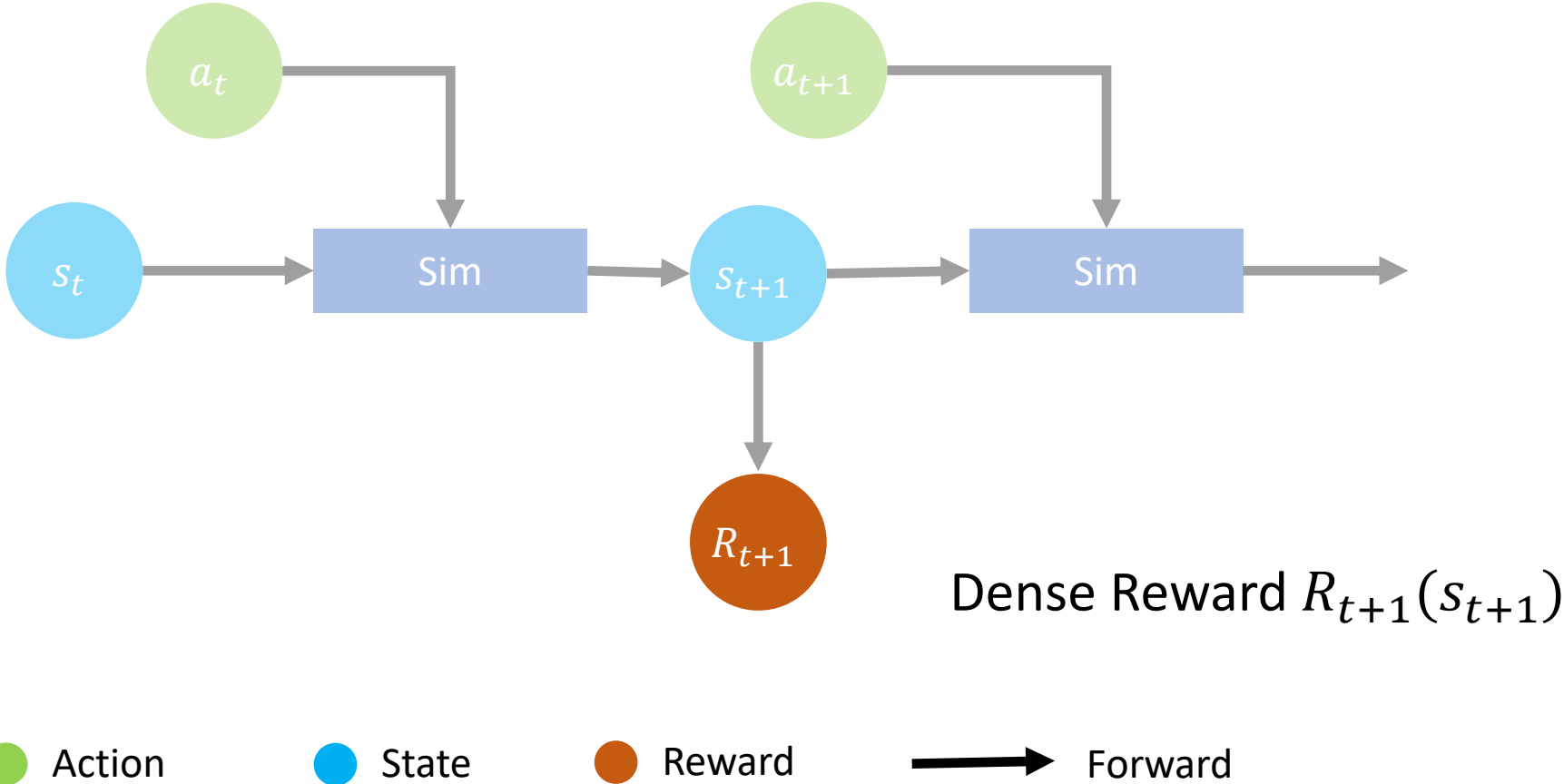
- Forward simulation with action



● Action ● State ● Reward → Forward

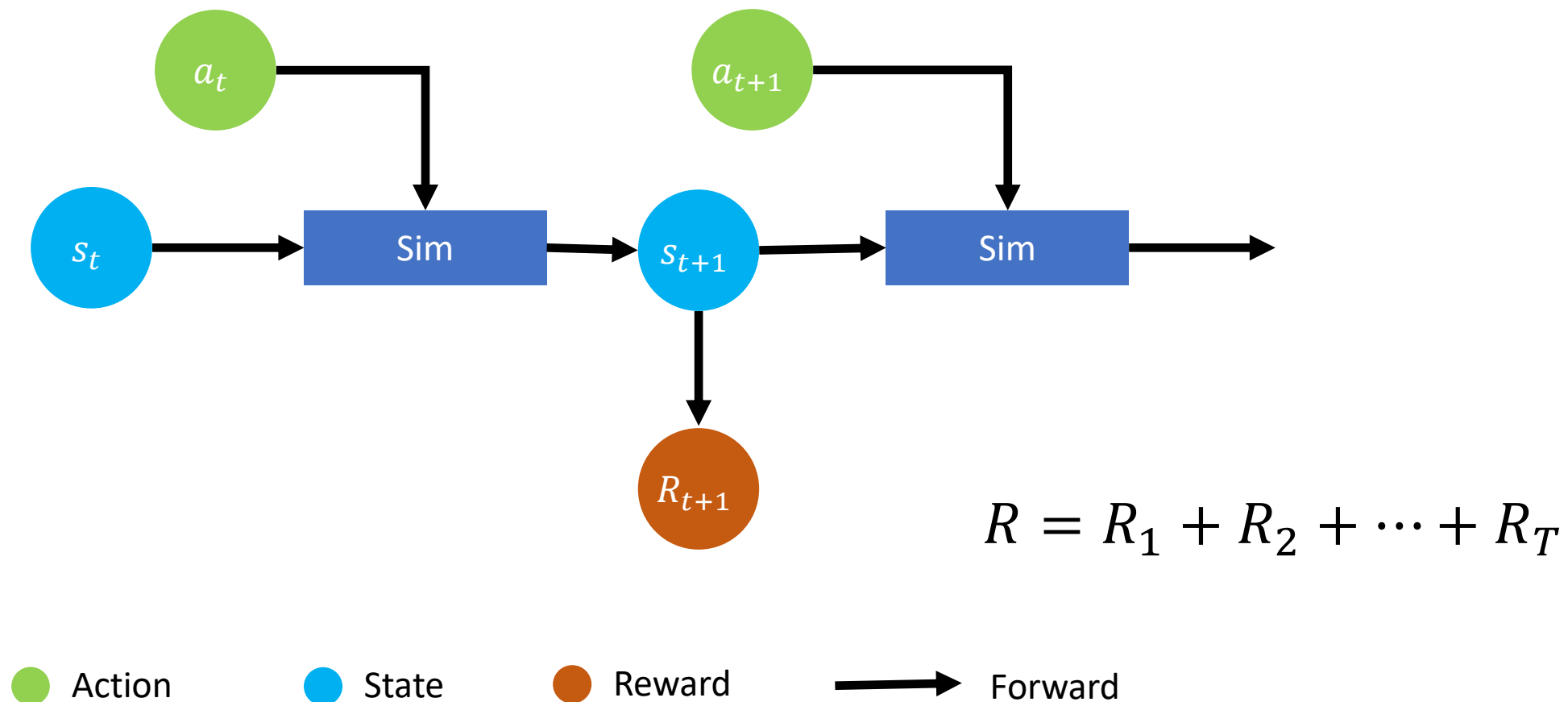
Differentiable Physics for Planning

- Reward



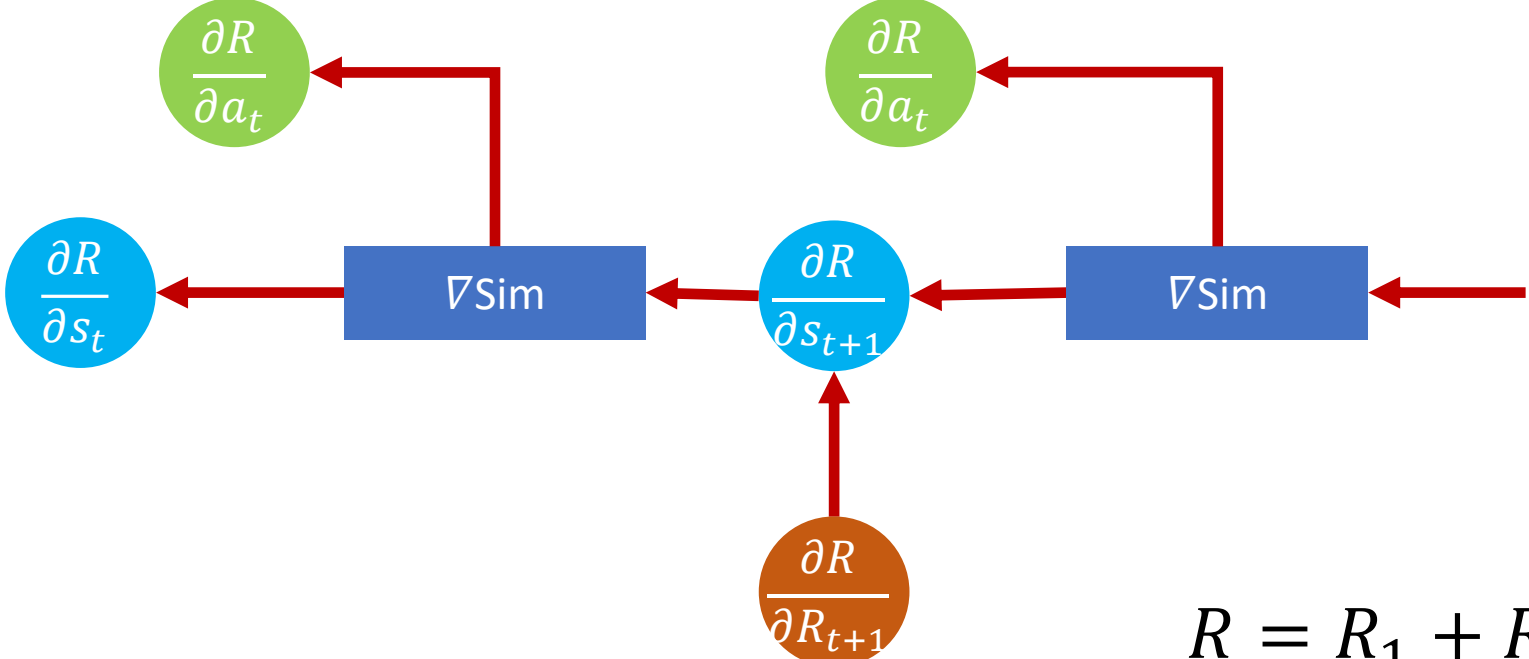
Differentiable Physics for Planning

- Maximize total reward R by choosing a_1, a_2, \dots, a_{T-1}



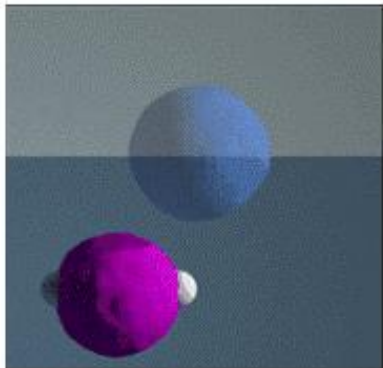
Differentiable Physics for Planning

- Backward

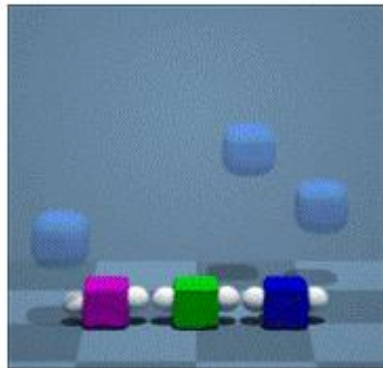


● Action ● State ● Reward ← Backward

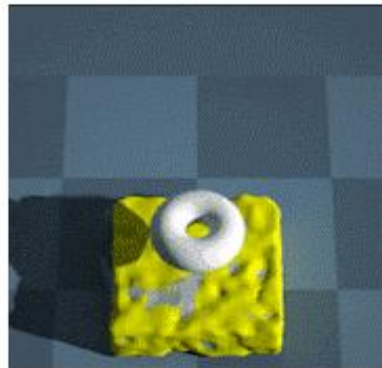
Applications: Soft-body Manipulation



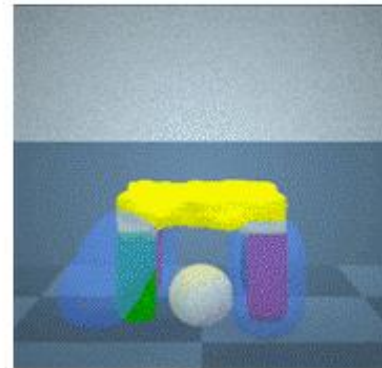
Move



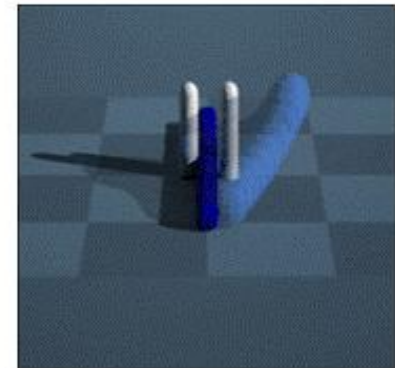
Triple Move



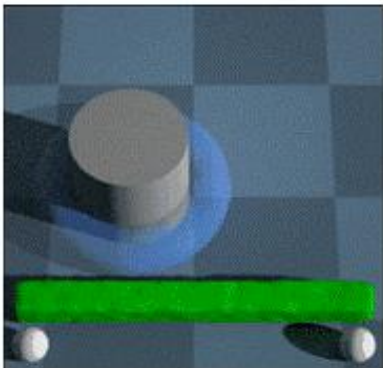
Torus



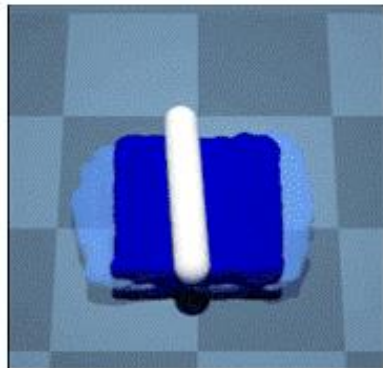
Table



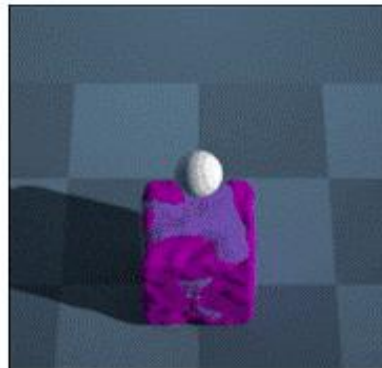
Chopsticks



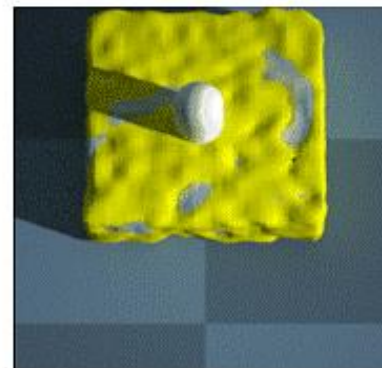
Rope



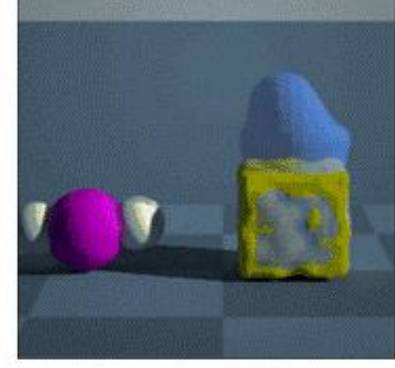
Rolling Pin



Pinch

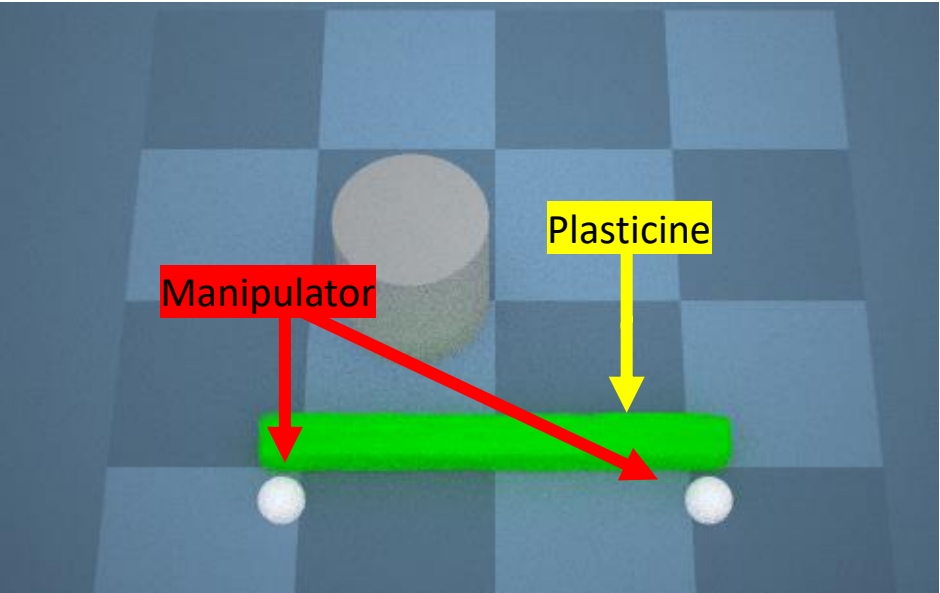


Writer

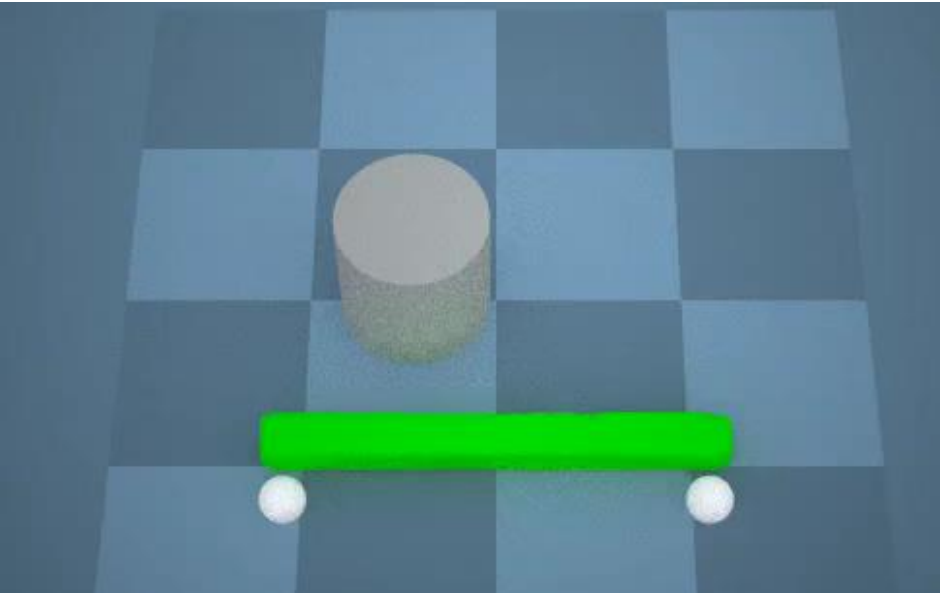


Assembly

Example: Rope Task



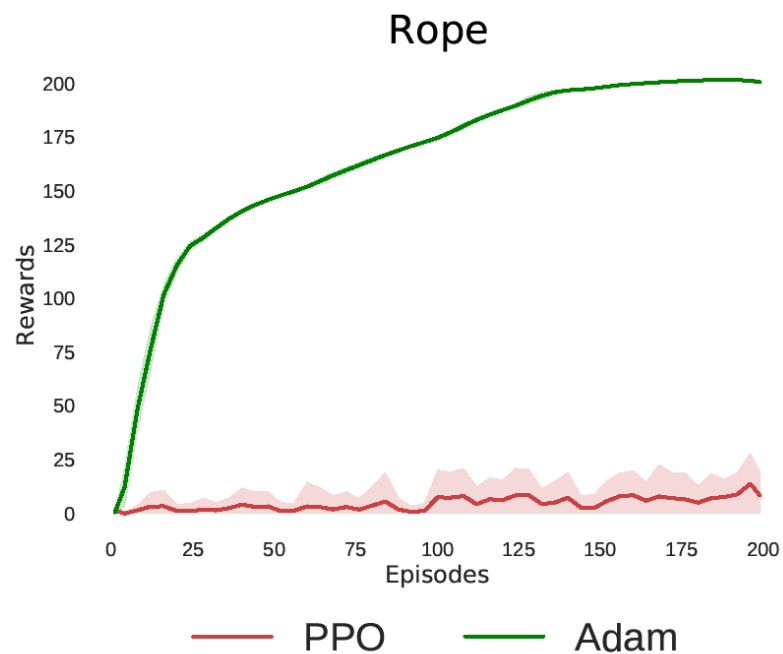
Initial State



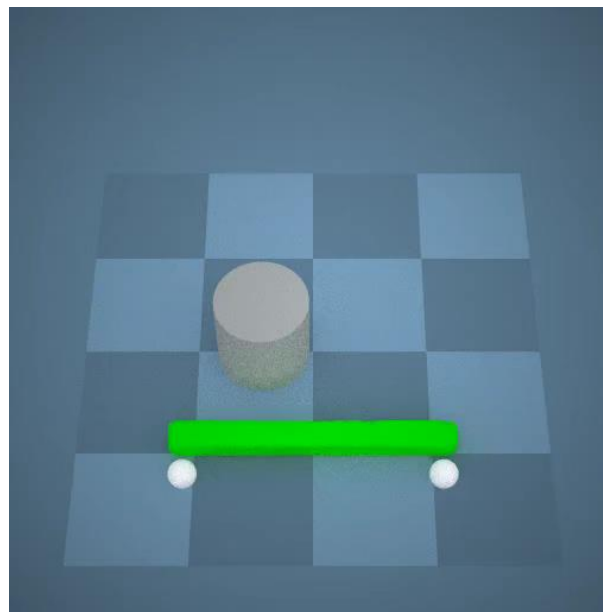
Deform

Evaluation Results

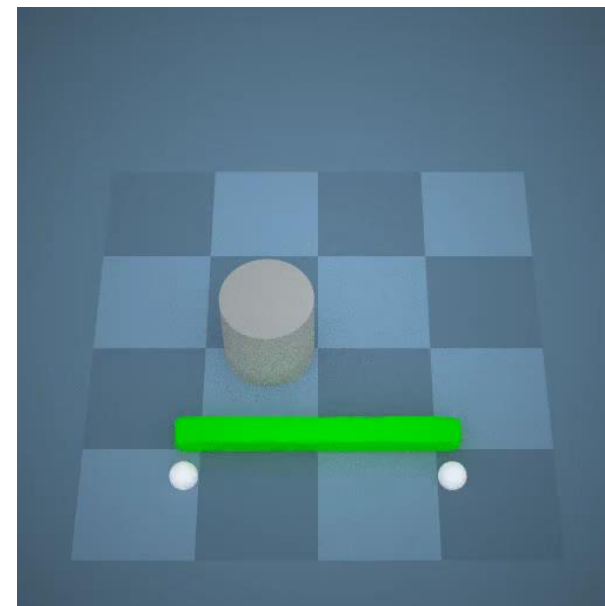
- Gradient-based optimization is much efficient than RL.



Adam Episode 200

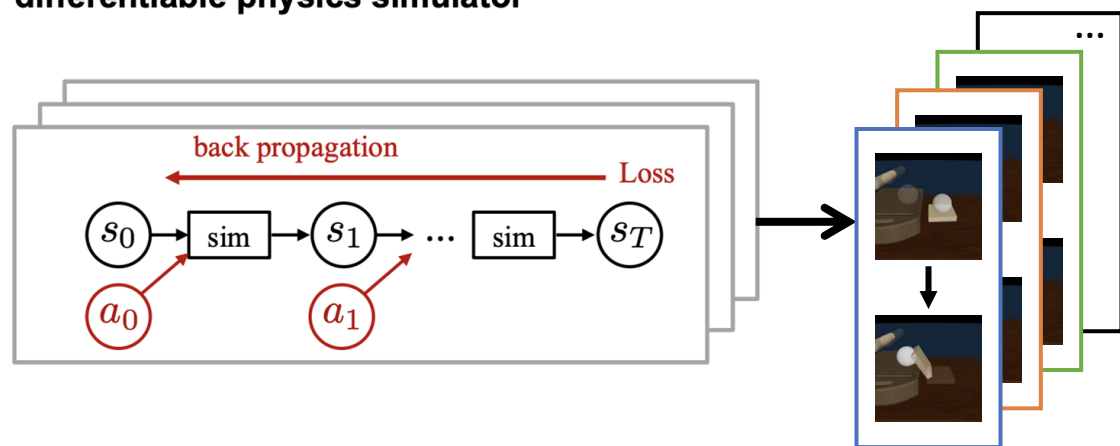


PPO Episode 10K

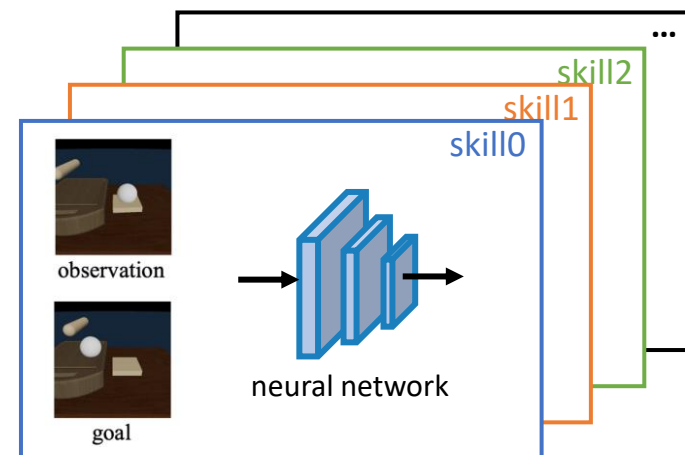


Abstract Skills Using Neural Networks

(a) Collecting demonstration trajectories in a differentiable physics simulator



(b) Neural skill abstraction



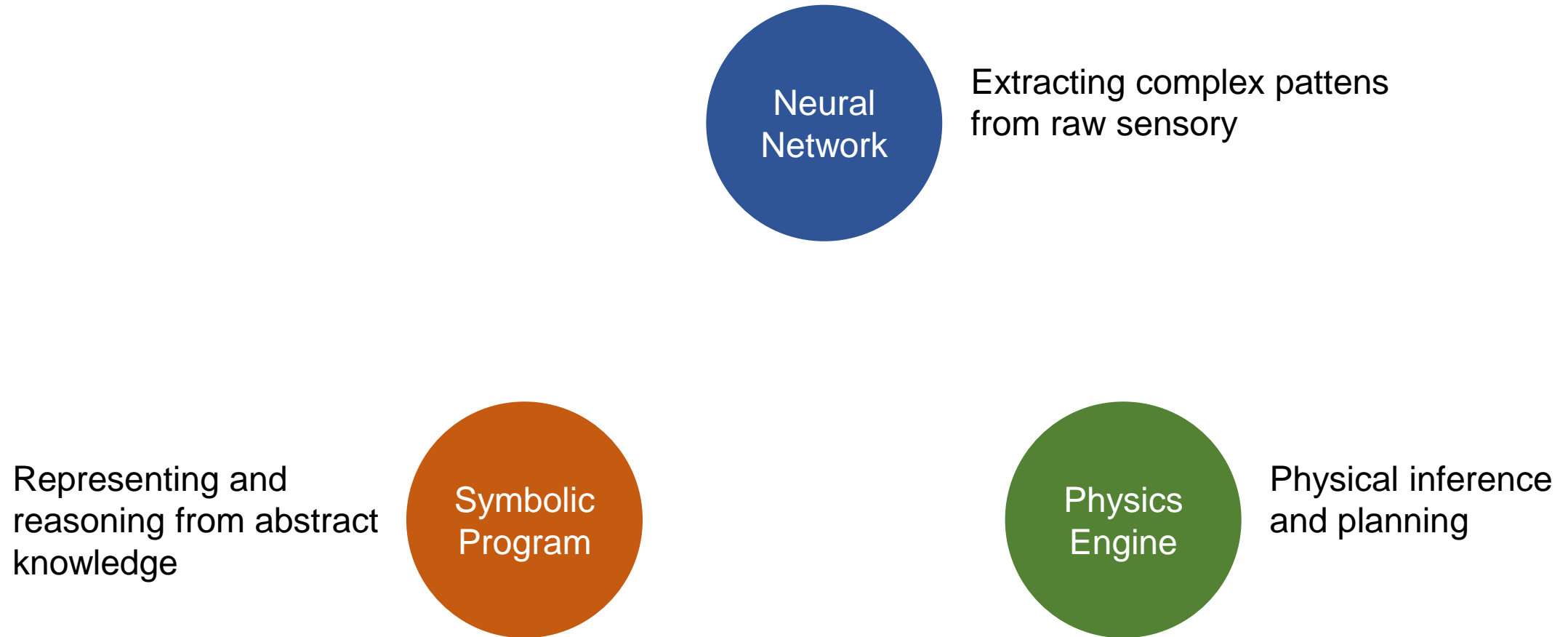
Rearrange and lift



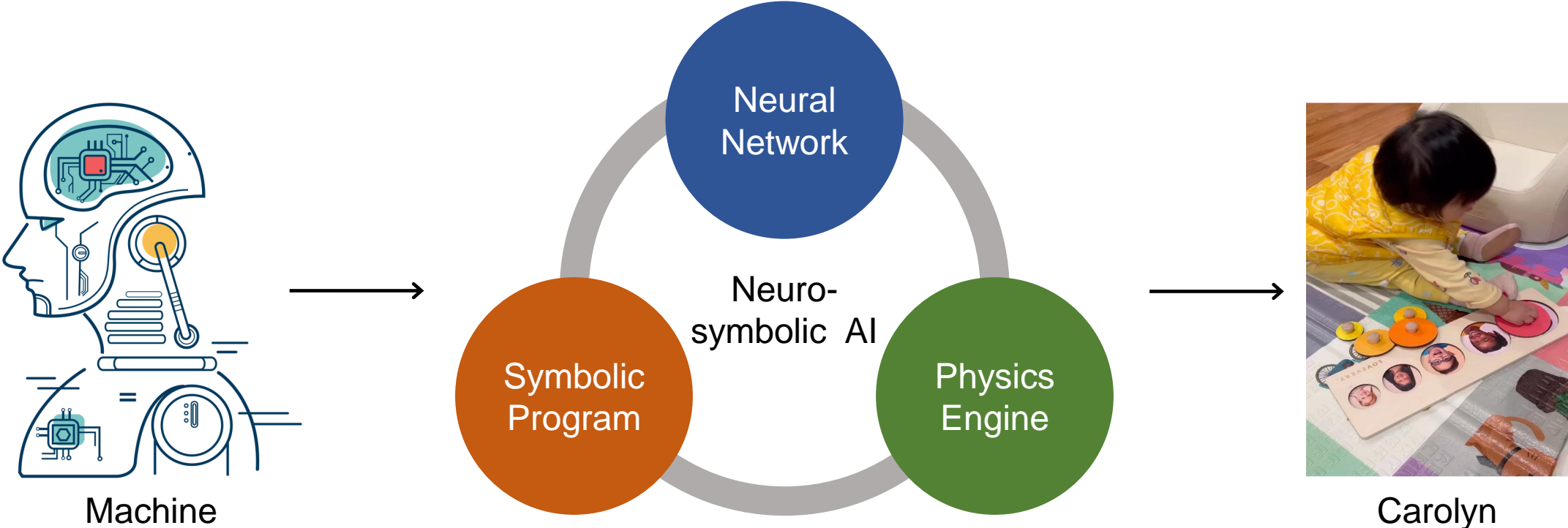
Lift and spread



My Vision: Neuro-Symbolic Embodied AI



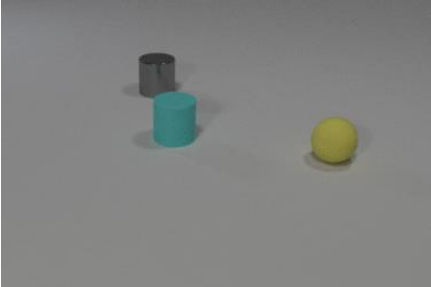
My Vision: Neuro-Symbolic AI



Scene Understanding



Dynamics Reasoning



Physical Interaction

