

Neuro-Symbolic Concept Learner*

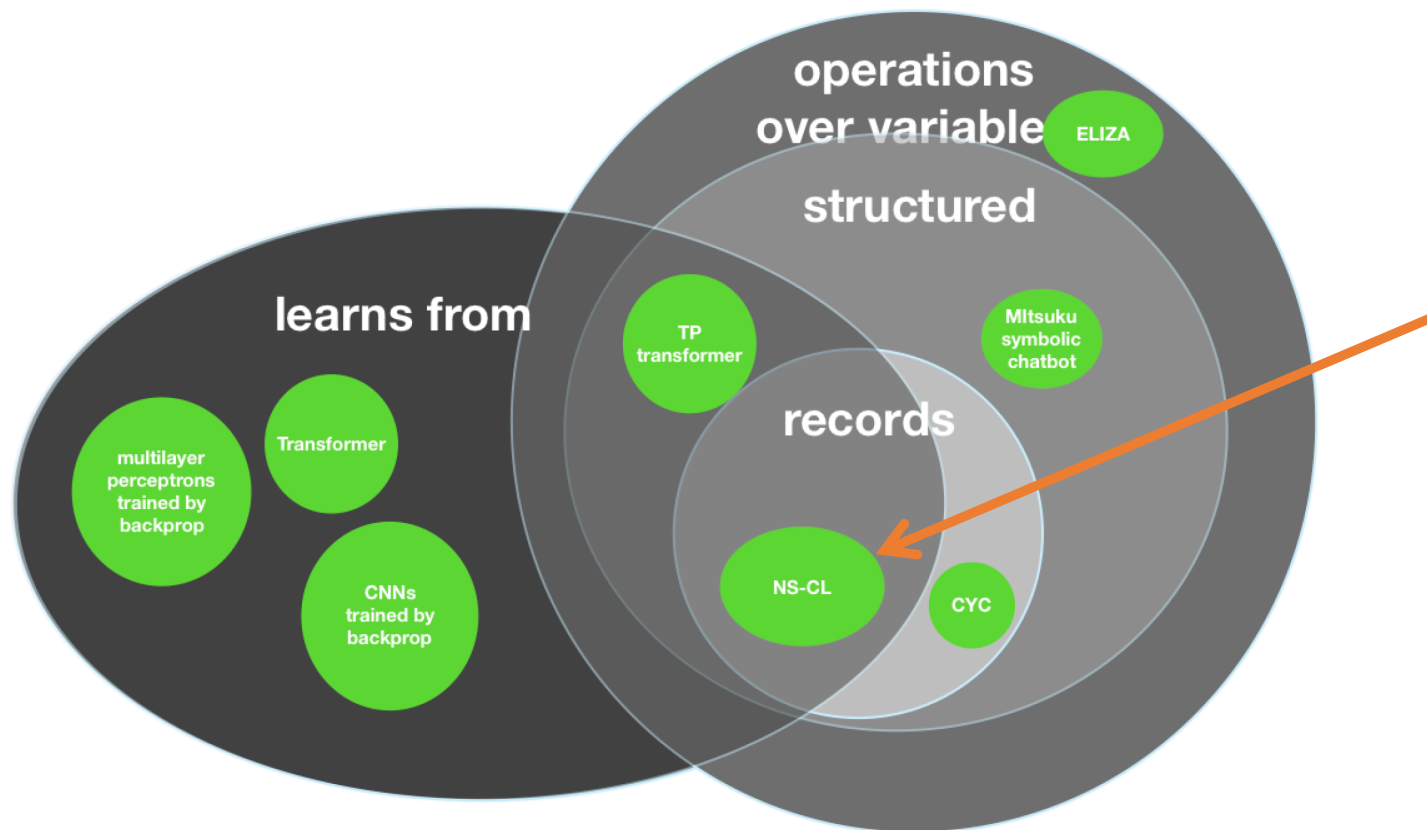
Sargur N. Srihari
srihari@buffalo.edu

*Based on J. Mao, C. Gan, P. Kohli, J. Tenenbaum, J. Wu, "The Neuro-symbolic Concept Learner: Interpreting Scenes, Words and Sentences from Natural Supervision, ICLR 2019

Topics

1. NS-CL within Robust AI
2. Natural Supervision and CLEVR Dataset
3. Neuro-Symbolic Concept Learner (NS-CL)
4. Training of NS-CL
5. Performance and Generalization
6. Relationship to Kahneman's Systems 1 and 2

Robust AI resides in the intersection of Learning and Symbol Manipulation



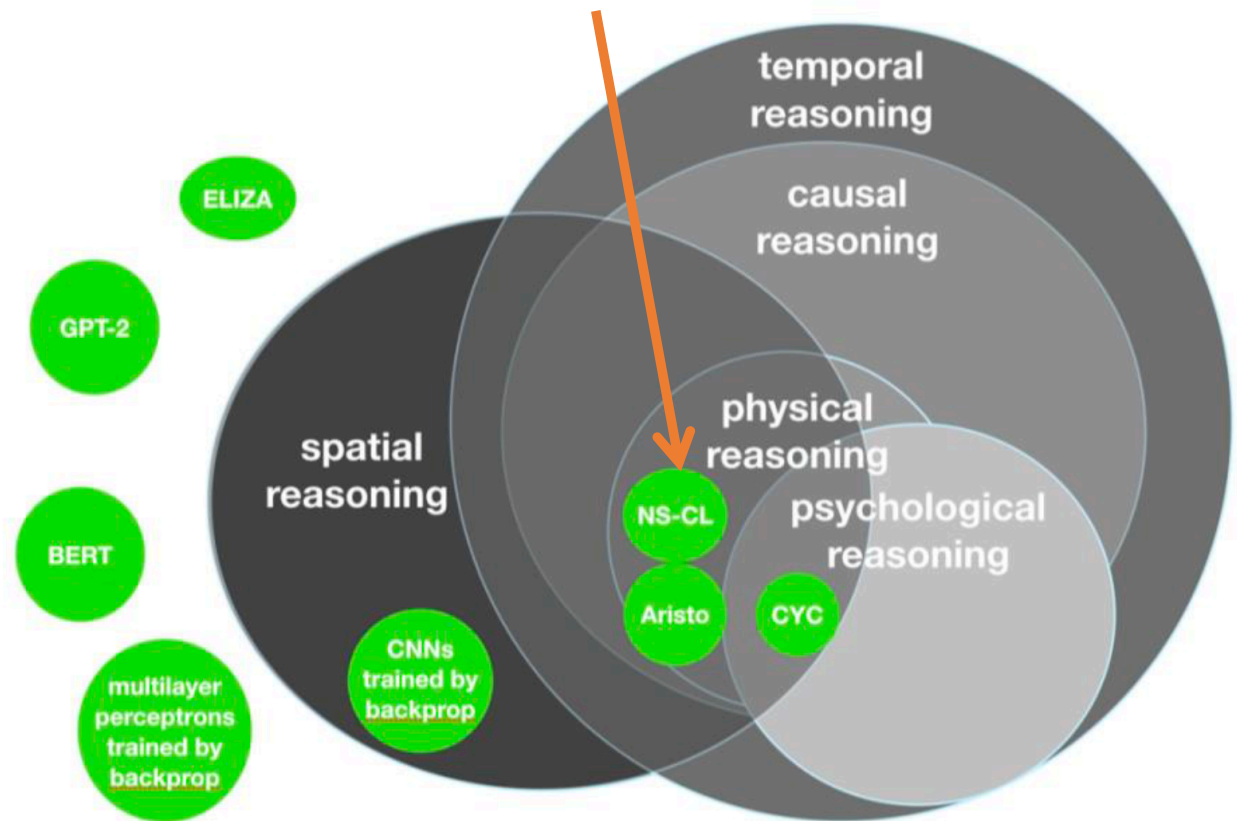
Reasoning Mechanisms

Robust AI needs reasoning machinery for

1. Spatial
2. Physical
3. Psychological
4. Temporal
5. Causal

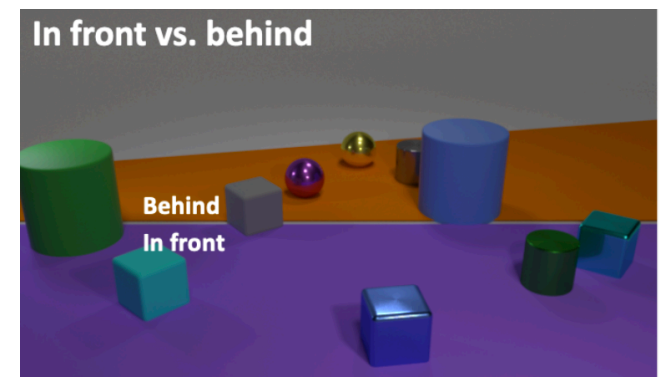
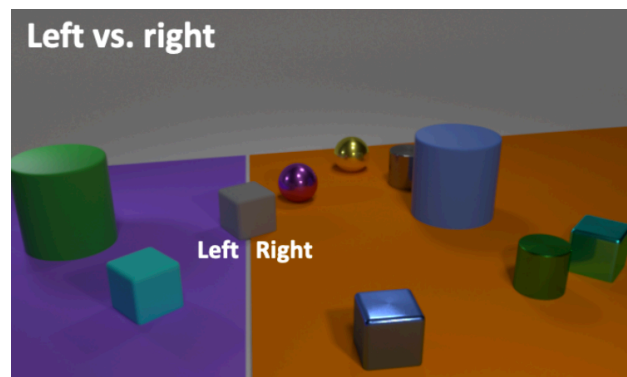
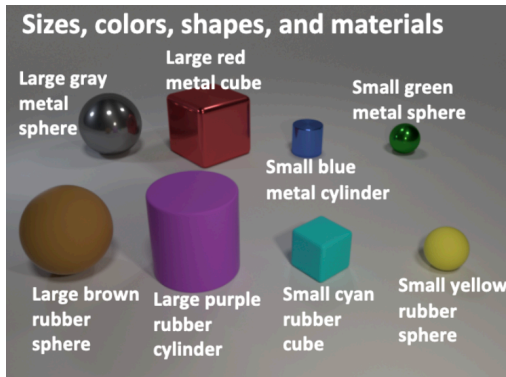
Current neural networks lack

1. Explicit mechanisms for these forms of reasoning
2. Ways of representing and reasoning over such domains



Visual Question Answering(VQA) dataset

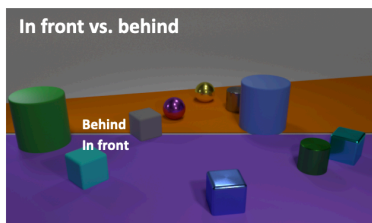
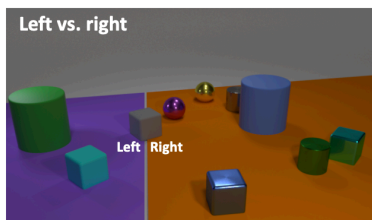
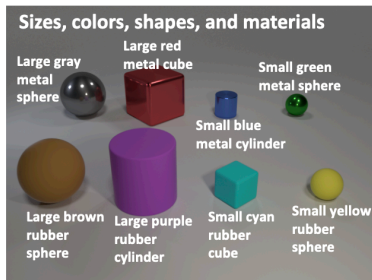
- CLEVR: A Diagnostic Dataset for Compositional Language and Elementary Visual Reasoning
 - J. Johnson, Li F, B. Hariharan, Lawrence Z, L. van der Maaten, R. Girshick



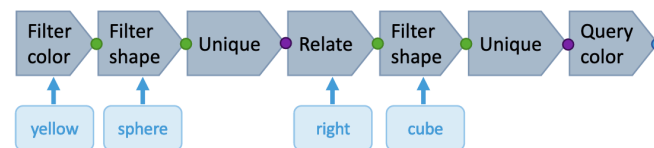
Natural language questions about images

Middle: Examples of questions and their associated functional programs.

Right: Catalog of basic functions used to build questions

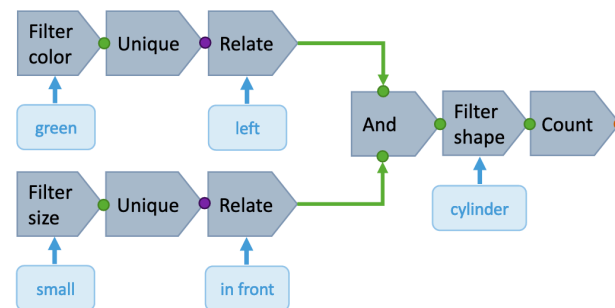


Sample chain-structured question:



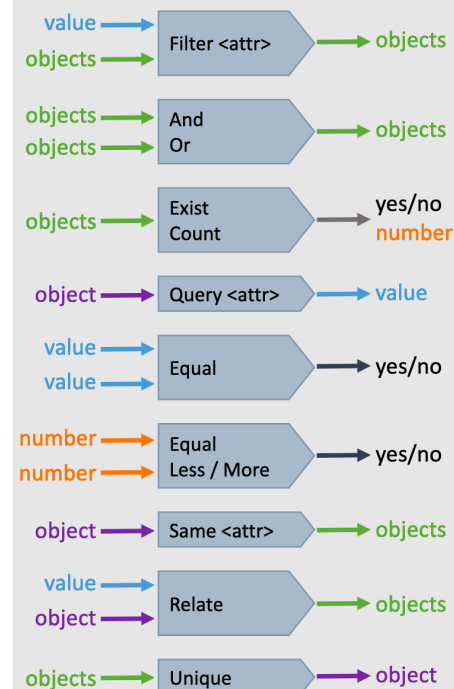
What color is the cube to the right of the yellow sphere?

Sample tree-structured question:



How many cylinders are in front of the small thing and on the left side of the green object?

CLEVR function catalog

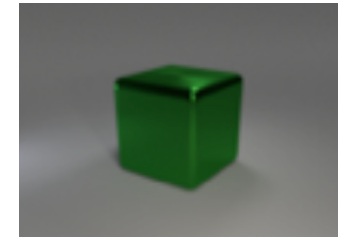
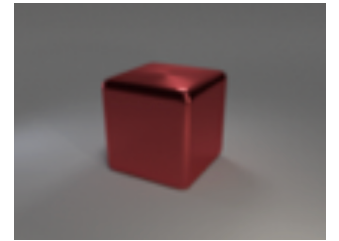


CLEVR Dataset

Split	Images	Questions	Unique questions	Overlap with train
Total	100,000	999,968	853,554	-
Train	70,000	699,989	608,607	-
Val	15,000	149,991	140,448	17,338
Test	15,000	149,988	140,352	17,335

Presentation of image

- Humans are capable of learning visual concepts by jointly understanding vision and language
- System with no prior knowledge of color is presented with images of the red and green cubes

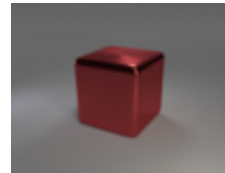


Natural Supervision in NS-CL

- Learns from images and QA pairs
 - Requires no annotations on images or semantic programs for sentences
 - Like human concept learning, it learns via curriculum learning
- Learns representations/concepts of objects from questions
 - (e.g., What's the color of the cylinder?) on simple scenes (3 objects)
 - By doing so, it learns object-based concepts such as colors and shapes.
- Learns relational concepts by leveraging object-based concepts to interpret object referrals
 - (e.g., Is there a box right of a cylinder?)
- Model adapts to more complex scenes and compositional questions

Presentation of image and QA

- System identifies the difference in object visual appearance (color), and aligns it to the corresponding words in the questions and answers (Red and Green)

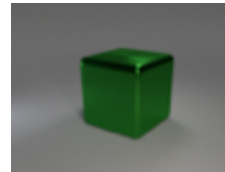


Q: What's the color of the object?

A: Red.

Q: Is there any cube?

A: Yes.



Q: What's the color of the object?

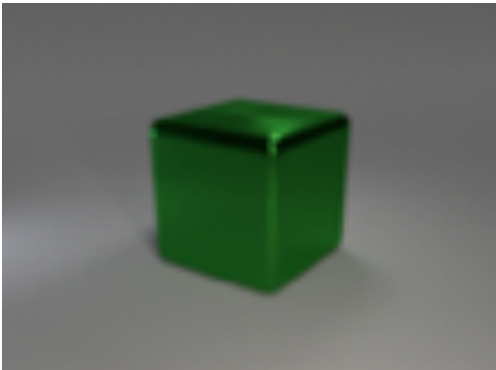
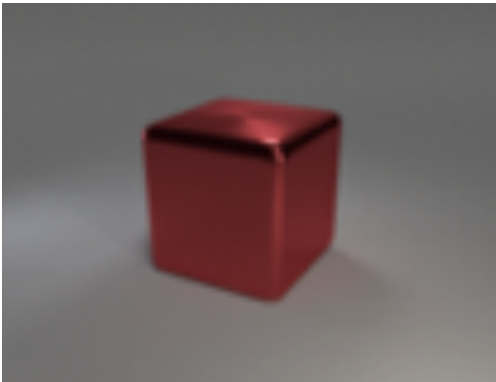
A: Green.

Q: Is there any cube?

A: Yes.

Learning basic, object-based concepts

Images

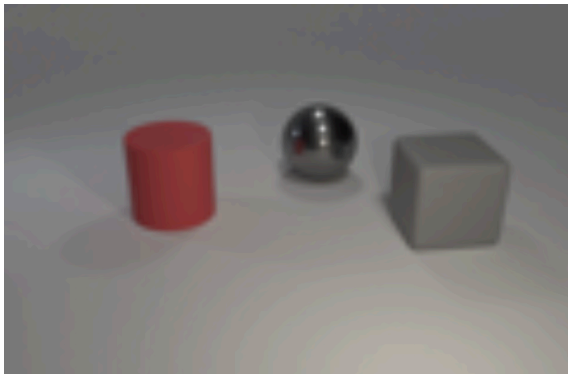


Paired Questions and Answers

- Q: What's the color of the object?
 - A: Red.
 - Q: Is there any cube?
 - A: Yes.
-
- Q: What's the color of the object?
 - A: Green.
 - Q: Is there any cube?
 - A: Yes.

Learning relational concepts based on referential expressions

Image

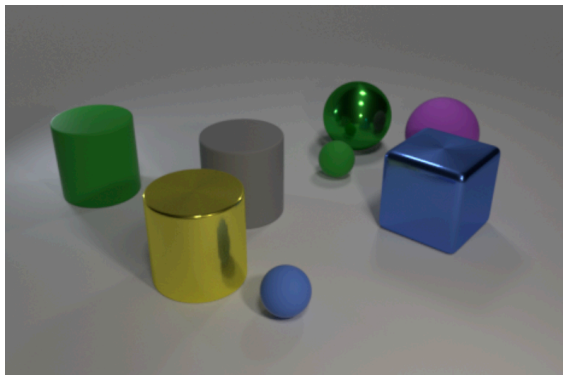


Paired Questions and Answers

- Q: How many objects are right of the red object?
- A: 2
- Q: How many objects have the same material as the cube?
- A: 2

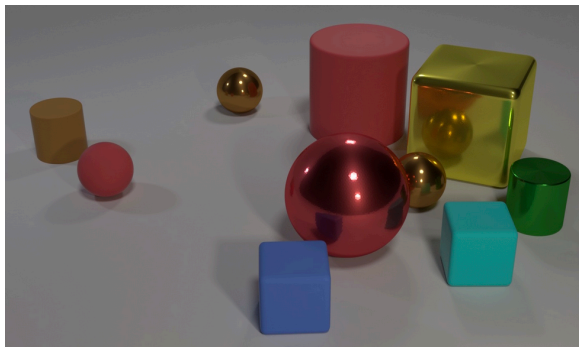
Interpret complex questions from visual cues

Image



Paired Question and Answer

- Q: How many objects are both right of the green cylinder and have the same material as the small blue ball?
- A: 3



NeuroSymbolic Concept Learner Model

- Model learns by
 1. Looking at images and
 2. Reading paired questions and answers.
- Model learns
 1. Visual concepts
 2. Words
 3. Semantic parsing of sentenceswithout explicit supervision on any of them

NS-CL has three modules

1. Neural Perception Module

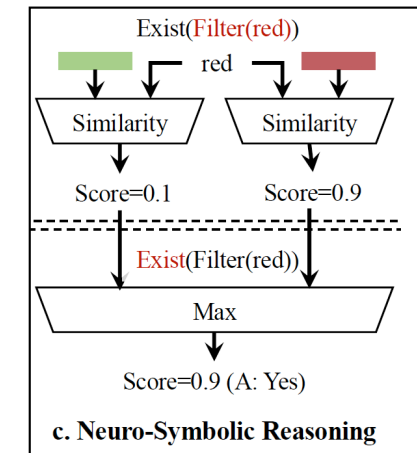
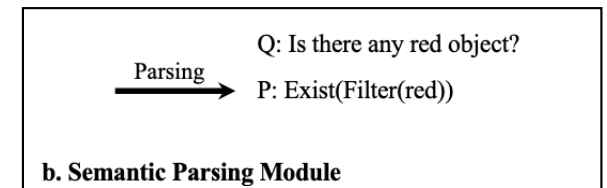
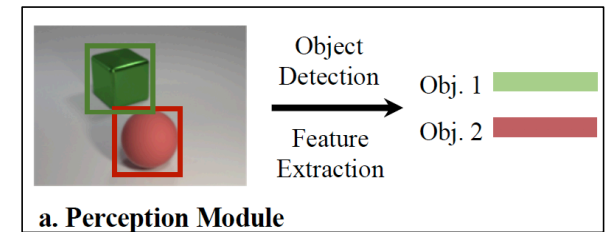
- Extracts object-based representations of scene
 - Each object represented as a vector embedding of fixed dimension

2. Semantic parsing Module

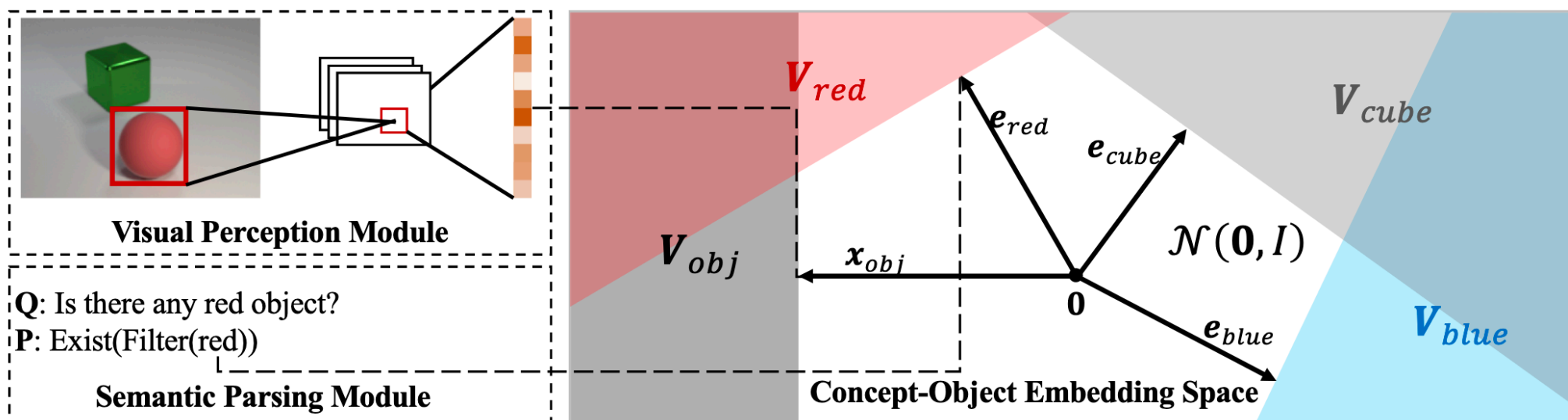
- Translates questions into executable programs
 - Program consists of hierarchical primitive functional modules
 - Is there any red object translated into two-step program: $\text{Exist}(\text{Filter}(\text{red}))$
 - Concepts represented as vector embeddings in a latent vector space

3. Neuro-Symbolic Reasoning Module

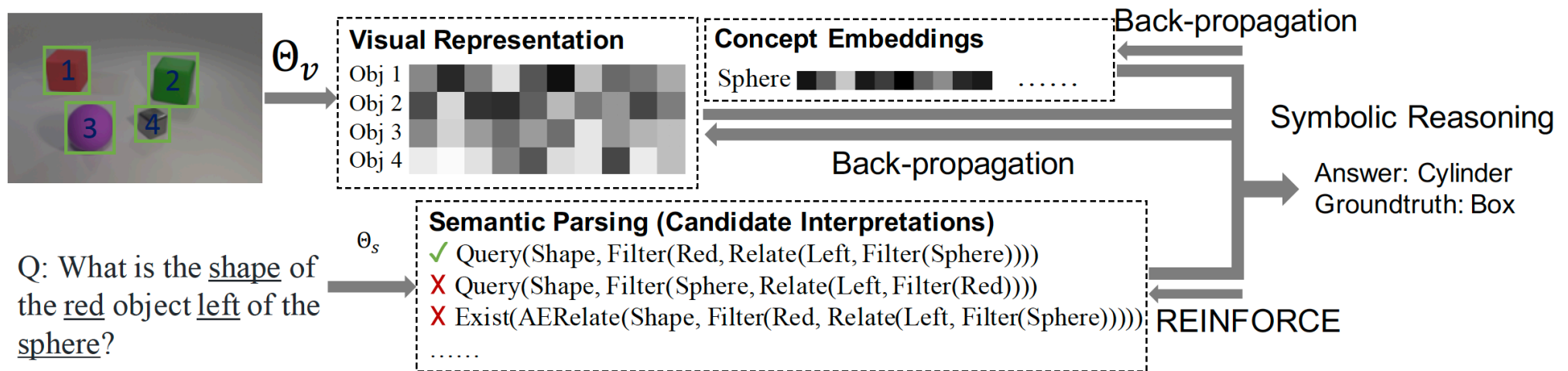
- Executes program to answer the question
 - It first computes the similarities between the concept embedding red and object embeddings to classify both objects.
 - Then, it answers the question by checking whether a red object has been filtered out
 - During training, it receives ground-truth answer as supervision and back-propagates training signals
 - Reads out the perceptual representation of objects, classifies their attributes/relations, and executes the program to obtain an answer



Embedding Space



Joint learning of visual concepts, words, and parsing



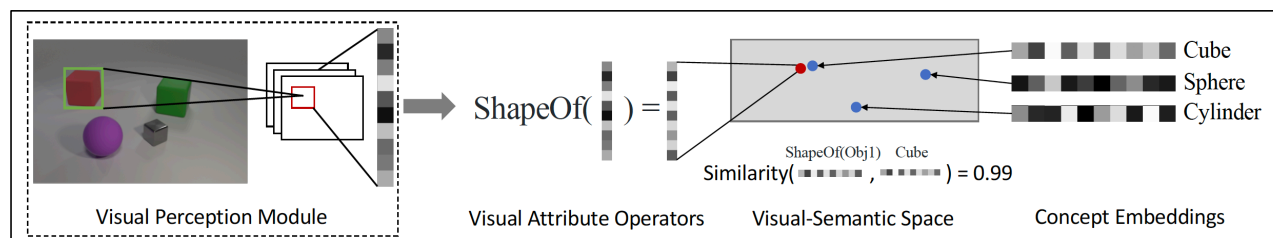
Visual perception module detects objects and extracts a latent representation for each of them

Semantic parsing module translates input question into an executable program given a domain specific language (DSL)

Concept Quantization

- Visual reasoning determines object attributes (e.g., color, shape)
 - Each attribute contains a set of visual concepts (e.g., Cube)
 - Attributes are neural mapping of object representation into an embedding

Inference of
Object shape:



- Concepts, e.g., Cube, Sphere, Cylinder, are learned vectors
 - Cosine distances between vectors $[\cdot, \cdot]$ determine the shape of the object
 - Probability that object o_i is a cube is $\sigma([\text{ShapeOf}(o_i); \mathbf{v}^{\text{Cube}}] - \gamma)/\tau$, where
 - $\text{ShapeOf}()$ denotes the neural operator,
 - \mathbf{v}^{Cube} the concept embedding of Cube
 - σ is the Sigmoid function
 - γ and τ are scalar constants for scaling and shifting similarity values
 - Classifying relational concepts (e.g., Left) between a pair of objects is done similarly,
 - except concatenate visual representations for both objects to form the representation of their relation

Visual perception details

- A pretrained Mask R-CNN generates object proposals for all objects
- Bounding box of object paired with its image is sent to a ResNet-34 to
 1. Extract the region-based (by RoI Align)
 2. Image-based features
- Concatenate them to represent each object
- Inclusion of the representation of the full scene adds context
 - which is essential for inference of relative attributes, e.g., size or position

Semantic Parsing and DSL

- Semantic parsing translates question into executable program
 - With a hierarchy of primitive operations
 - Represented in a domain-specific language (DSL) designed for VQA
- DSL has a set of operations for visual reasoning
 - e.g., filtering out objects with concepts or query attribute of an object
- Operations share same input/output interface
 - Thus combined to form programs of any complexity

GRU implementation

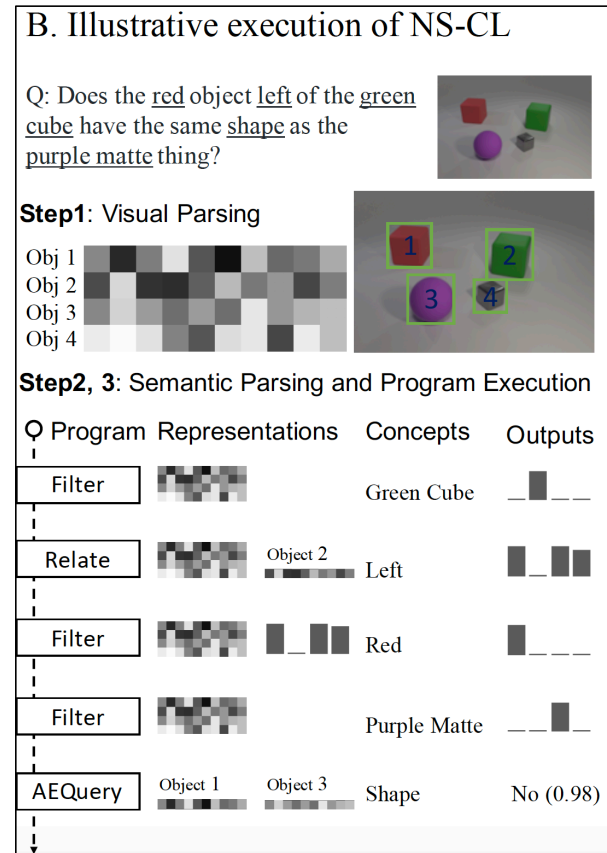
- Semantic parser generates hierarchies of latent programs in a sequence to tree manner
- A bidirectional GRU encodes the input question, which outputs a fixed-length embedding of the question
- A decoder based on GRU cells is applied to the embedding, and recovers the hierarchy of operations as the latent program
- Some operations take concepts as their parameters, such as Filter(Red) and Query(Shape)
 - These concepts are chosen from all concepts in the input question

Quasi-symbolic program execution

- Given the latent program recovered from the question in natural language, a symbolic program executor executes the program and derives the answer based on the object-based visual representation
- Program executor is a collection of deterministic functional modules to realize logic operations specified in the DSL

Differentiability wrt Visual Operations

- Intermediate results are probabilistic
 - Set of objects is a vector
 - The attention mask, over all objects
 - $\text{Mask}_i \in [0,1]$ denotes probability that i -th object belongs to the set
 - First Filter outputs a mask of length 4
 - Output “mask” on objects is fed to next module (Relate in this case) as input
- The execution of programs continues



Training Paradigm

- Two part objective:
 - Concept learning
 - Language understanding
- Find optimal parameters
 - θ_v of the visual perception module Perception (including the ResNet-34 for extracting object features, attribute operators. and concept embeddings) and
 - θ_s of the semantic parsing module SemanticParse, to maximize the likelihood of answering the question Q correctly

$$\Theta_v, \Theta_s \leftarrow \arg \max_{\Theta_v, \Theta_s} \mathbb{E}_P[\Pr[A = \text{Executor}(\text{Perception}(S; \Theta_v), P)]]$$

- P denotes the program,
- A the answer,
- S the scene, and
- Executor the quasi-symbolic executor.
- The expectation is taken over $P \sim \text{SemanticParse}(Q; s)$

Optimization

- Compute the gradient w.r.t. v as

$$\nabla_{\Theta_v} \mathbb{E}_P[D_{\text{KL}}(\text{Executor}(\text{Perception}(S; \Theta_v), P) \| A)]$$

- REINFORCE (Williams, 1992) optimizes the semantic parser θ_s :

$$\nabla_{\Theta_s} = \mathbb{E}_P[r \cdot \log \text{Pr}[P = \text{SemanticParse}(Q; \Theta_s)]]$$

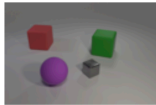
- where the reward $r = 1$ if the answer is correct and 0 otherwise
- Off-policy search is used to reduce the variance of REINFORCE

Curriculum Visual Concept Learning

A. Curriculum concept learning

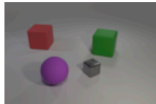
○ **Initialized** with DSL and executor.

□ **Lesson1:** Object-based questions.



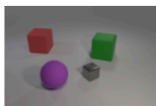
Q: What is the shape of the red object?
A: Cube.

□ **Lesson2:** Relational questions.



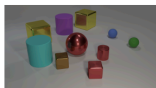
Q: How many cubes are behind the sphere?
A: 3

□ **Lesson3:** More complex questions.



Q: Does the red object left of the green cube have the same shape as the purple matte thing?
A: No

□ **Deploy:** complex scenes, complex questions



Q: Does the matte thing behind the big sphere have the same color as the cylinder left of the small matte cube?
A: No.

Split the training samples into four stages

1. Learning object-level visual concepts;
2. Learning relational questions;
3. Learning complex questions with perception modules fixed
4. Joint fine-tuning of all modules.

This was found essential to the learning of NS-CL

Dataset

- Train NS-CL on 5K images (<10% of CLEVR's 70K training images)
- Generate 20 questions for each image for the entire curriculum learning process
- The Mask R-CNN module is pretrained on 4K generated CLEVR images with bounding box annotations

Performance: NS-CL vs. Baseline

- Evaluation on simple questions like “How many red objects are there?”
- Evaluated performance on all concepts in the CLEVR dataset
- NS-CL outperforms both convolutional and attentional baselines
- NC-CL has better interpretation of visual concepts

	Visual	Mean	Color	Mat.	Shape	Size
IEP	Conv.	90.6	91.0	90.0	89.9	90.6
MAC	Attn.	95.9	98.0	91.4	94.4	94.2
TbD (hres.)	Attn.	96.5	96.6	92.2	95.4	92.6
NS-CL	Obj.	98.7	99.0	98.7	98.1	99.1

Summary of NS-CL Performance

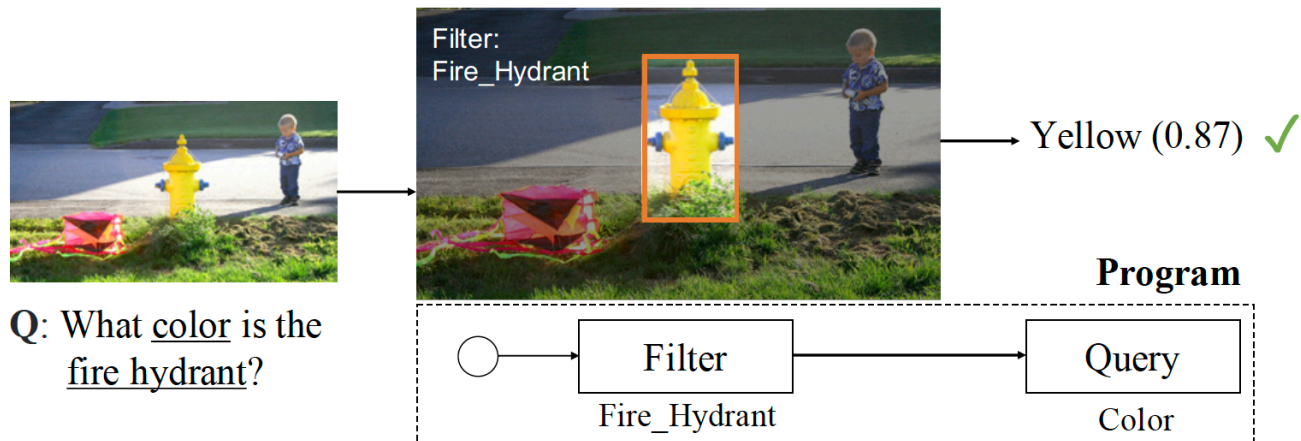
- Interpretable and robust visual reasoning:
 - State-of-the-art performance on CLEVR dataset
 - Naturally learns disentangled visual and language concepts,
 - Enabling combinatorial generalization w.r.t. both visual scenes and semantic programs

NS-CL Generalization

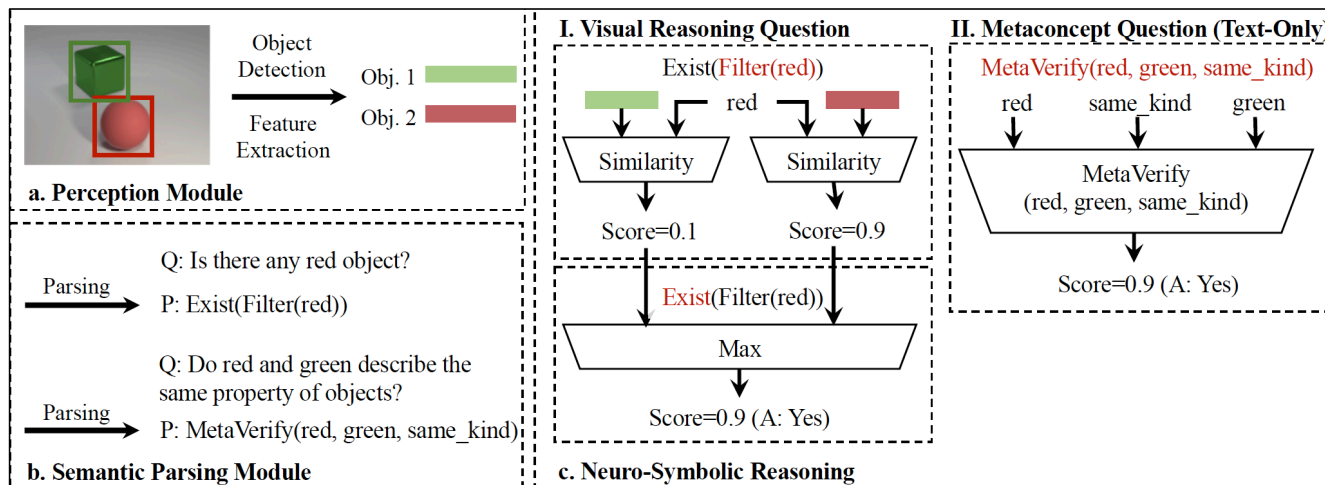
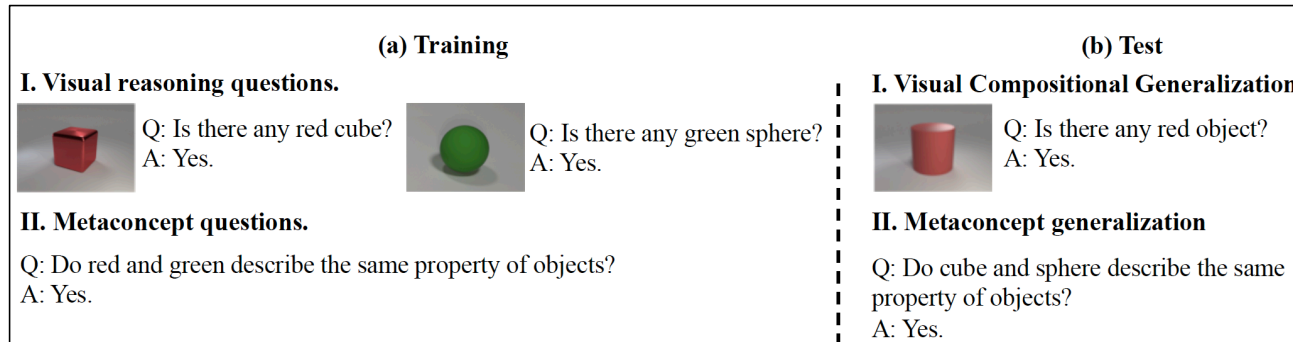
1. Generalizes to scenes with more objects and longer semantic programs than those in the training set
2. Generalizes to new visual attribute compositions as demonstrated on the CLEVR-CoGenT dataset
3. Adaptation to new visual concepts
 - such as learning a new color
4. Learned visual concepts transfer to new tasks
 - such as image-caption retrieval, without any extra fine-tuning

Extension to Natural Images

- Comparable results with baselines







NS-Meta-concept Learning (NS-ML)







Performance of NS-ML

		Q.Type	GRU (Lang. Only)	GRU-CNN	BERT (Variant I ; Variant II)	NS-CL	VCML
CLEVR	Synonym	50.0	66.3 \pm 1.4	60.9 \pm 10.6	76.2 \pm 10.2 ; 80.2 \pm 16.1	1 . \pm	1 . \pm
	Same-kind	50.0	64.7 \pm 5.1	61.5 \pm 6.6	75.4 \pm 5.4 ; 80.1 \pm 10.0	92.3 \pm 4.9	99.3\pm1
GQ	Synonym	50.0	80.8 \pm 1.0	76.2 \pm 0.8	76.2 \pm 2.4 ; 83.1 \pm 1.5	81.2 \pm 2.8	91.1\pm1.7
	Same-kind	50.0	56.3 \pm 2.3	57.3 \pm 5.3	59.5 \pm 2.7 ; 68.2 \pm 4.0	66.8 \pm 4.1	69.1\pm1.7
CUB	Hypernym	50.0	74.3 \pm 5.2	76.7 \pm 8.8	75.6 \pm 1.2 ; 61.7 \pm 10.3	80.1 \pm 7.3	94.8\pm1.3
	Meronym	50.0	80.1 \pm 5.9	78.1 \pm 4.8	63.1 \pm 3.2 ; 72.9 \pm 9.9	97.7\pm1.1	92.5 \pm 1.0

(a) Synonym Metaconcept Generalization

Training		
	Q: Is there any <i>airplane</i> ? A: Yes.	
	Q: Is there any <i>kid</i> ? A: Yes.	
		Q: Is airplane a <i>synonym</i> of plane? A: Yes.
		Test Q: Is kid a <i>synonym</i> of child? A: Yes.

(b) Same_kind Metaconcept Generalization

Training		
	Q: Is there any <i>red</i> object? A: Yes.	
	Q: Is there any <i>truck</i> ? A: Yes.	
		Q: Do red and yellow <i>describe the same property of objects</i> ? A: Yes.
		Test Q: Do bus and truck <i>describe the same property of objects</i> ? A: Yes.

NS-CL and Kahneman's Systems 1,2

1. Capabilities of Systems 1,2
 - What features do they use?
2. Is System 2 inherently sequential?
3. Metrics to evaluate performance of Systems
4. Introspection in Systems 1, 2
 - Can M-consciousness alone solve the problem?
5. Governance: Switching between Systems
 - Recognize when each system to be used
6. Reasoning when using both Systems
 - Heuristics of System 1 with Careful Reasoning of System 2?
7. Ethics theories tied to Systems 1,2
 - Decide if action is moral or not
8. Abstraction: attention, knowledge graphs
 - Skills transferable to another environment?
9. Combining multi-agent systems
 - Several AI systems communicating
10. Architectures for supporting this vision of AI
 - Contribute to building model of world